

## スパコンを内から見て

技術課 神谷 基司

かみや・もとし／名古屋大学大学院理学研究科で博士（理学）を取得後、京都大学や理化学研究所での研究活動を経て2018年4月に特任専門員として分子科学研究所技術課計算科学技術班に採用。翌年2月に技術職員として採用され、現在に至る。



### はじめに

2018年4月より分子研に入所し、それから計算科学研究センターでスーパーコンピュータの管理者として色々やっています。その前は京都大学や理化学研究所でポスドクとして研究活動を行っていました。その間には現在管理者をしている計算科学研究センターや京コンピュータを含め、いくつかのスパコン（大型計算機）を使ってきました。その一方で、学生時代に所属研究室の計算機クラスタの管理者に祭り上げられてしまい、直接研究には関係無い管理者スキルを磨き上げてしまってもいました。ポスドク生活になってもその辺は極端には変わらず、所属研究室の計算機環境の整備にはなんとなく関わっていました。そんな生活の中で色々見てきた結果、がんばってPIになりたいとは思えず（そもそもなれるのかという問題も大いにありましたが）、かと言ってPIにならずに生き残る道と

いうのも極めて狭いように見えたため、どうしようかな、と置いていたところ、なんか舟があったので「乗せて」と言ってみたら乗せてもらえて、後はなんとなく現在に至った次第です。乗ってみたら意外と舟を漕ぐのが大変で、もう少し漕ぎ賃が欲しいもんだと思ったりすることもあります。とりあえず、スパコンのユーザーから管理者に鞍替えし、今のところボチボチやっています。

### スパコン運用の効率化

現在の計算科学研究センターのスパコンは以下の図1のような全体構造になっています。この内、実際に計算を行うのは演算ノード領域になります。演算ノード領域（図2）には演算ノードがたくさんありますが、基本的には通常のパソコンをやたらと多数集めたようなものです。性能面や安定性やネットワークとかで色々差があります。スパコンの演算性能を効率的に

使うには、これら演算ノードができるだけサボらないように、可能な限りたくさんジョブ（計算）を割り当て、できるだけブラックな労働を課す必要があります。そのような理念の元、ユーザーの不利にならない範囲で混んでいる領域のジョブを自動的に空いている領域に振り分け、演算性能を搾取するような設定をセンターでは以前から構築していました。この振り分けですが、私の入所時点では優先度に問題があったり、一部手動で操作するなどしていた部分があったりもしました。私がこの運用管理に携わるようになってからはこの仕組みについて継続して修正や拡張を行い、現在は図3のような形で、自動的にほぼこちらが思い描いた通りの順序で振り分けが行えるようになっています。完全に公平な振り分けだということはできませんが、限られた資源をかなり効率的に運用できています。以上を改善を行うに

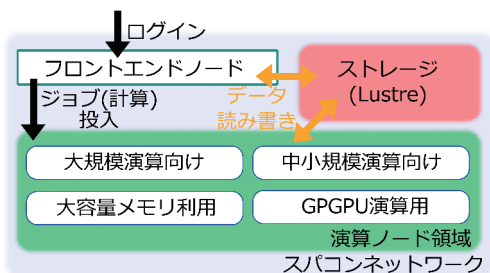


図1 計算センタースパコンの概要  
演算用領域については、計算時間などに応じて追加の内部区分があります。



図2 演算ノードのラック（一部）  
写っているラックには大規模演算向けのノード（合計約800台）が搭載されています。

当たって、システム側の設定変更に加えて、jsub というセンタースパコンのジョブ投入用ソフトウェアも重要な役割を果たしています。このような重要なソフトは継承という観点からも複数の眼で見ることが大切です。そのため、私の発案が採択されGitLabというソフトを導入しました。現在ではソフトのコードを他のメンバーにも見える形で管理しています (図4)。

## ストレージ問題

一方、現在、我々計算センターのユーザーの皆様から最も評判が悪いであろう部分の一つがストレージのアクセス速度だと思われます。体感的に「遅い」と感じているユーザーが多数でしょう。自分自身ですらそう思うことがあります。計算センターではLustreという分散ファイルシステムが用いられています。このLustreは研究室レベルで用いることはほとんどありませんが、巨大なファイルや大規模データの並列書き込みの扱いに優れているため、スパコンでは非常に一般的に採用されています。それくらいに優秀な仕組みではあるのですが、大量の小さい

ファイルの取り扱いはどうしても苦手です。そのように細かいファイルを大量に読むような状況は、ユーザーが自身のプログラムで解析データを読み込む際、あるいは科学計算では非常に著名なAnaconda環境を読み込む際、さらには大量のファイルがあるディレクトリでファイルリストを表示する際などに発生します。端的に言えばものすごく日常的に発生しています。その結果としてユーザーにディスクが遅いという認識が生まれやすい状況になっていると思われま。これまでセンターとベンダーが一緒になって色々としてきましたが、現時点ではうまい手が無い状況です。今後Data on MDTのような新しい機能が使えるような状況となれば改善の可能性があるのですが……。というわけで、ユーザーの皆様にはもうしばらくご不便おかけすることになりますが、何卒お許し下さい。

## ユーザーの声

さて、スパコンの利用者から管理者になるに当たって、要望というものの、例えばこういうことがしたいから

何とかしろ、こういうソフト使いたいから入れてくれ、といったものもつとあるのだと想像していました。しかし、実際こっち側に来て分かったことはこのような要望は思ったよりも来ないということです。無理難題のような要望も少なくはないという話も聞いていたので、ちょっと意外でした。なお、新規利用者向けのガイド (<https://ccportal.ims.ac.jp/quickstartguide>) や英語情報の提供、自分のジョブが実行されるまでどの程度待つ必要があるのかを概算でも良いので知りたい (waitest コマンドで実装)、という頻出の要望については完璧とは絶対言えませんが、ガイドや新規コマンドの作成である程度対応しています。

そんな感じに要望が来ないことは仕事量的な意味で楽ではあるのですが、致命的ではないけど実は割と不便を感じている、というような微妙な問題が顕在化しないということでもあるので、大歓迎とは言えません。ユーザー側では不満を感じているけど管理側では気づいてなくて、言ってさえくれさえすれば実は対応できるケースも存在しているような気はしています。一例と

## 大規模計算向け領域(ノード単位で割り当て)

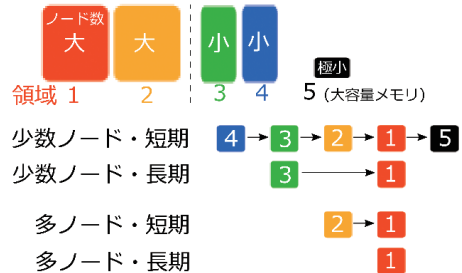


図3 ジョブの割り当て優先度。「大規模演算向け」の一部ジョブは「大容量メモリ利用」領域に割り当てられることがあります。

```

431     filter( lambda k:
432         ( ( self.jobs[k].get( "job_state", None ) == "Q" or \
433           self.jobs[k].get( "job_state", None ) == "H" ) ) and \
434           not uh.search( self.jobs[k].get( "Hold_Types", "n" ) ),
435           self.jobs.keys() )
436
437     # priority-sorted jobs
438     self.q_jobsids = sorted( q_jobsids,
439                             key=lambda x:
440                               ( -int(self.jobs[x].get( "Priority", None ))
441                                 self.jobs[x].get( "qtime", None ) )
442                             )
443
444     def sortRunningJobs( self ):
445         r_jobsids = filter( lambda k:
446             self.jobs[k].get( "job_state", None ) == "R",
447             self.jobs.keys() )
448         self.r_jobsids = sorted( r_jobsids,
449                                 key = lambda k: self.jobs[k].get( "remaining_t
450                                     reverse = False )
451
452     def checkResourcesStat( self, q, u, g, ncpu, ngpus ):
453         cur_cpu_u = ncpu
454         cur_cpu_g = ncpu
455         cur_gpu_u = ngpus
456         cur_gpu_g = ngpus
457         if u in self.stat_u:
458             cur_cpu_u += self.stat_u[u][0] # cpu - user
  
```

図4 GitLabで管理されているコード。



して、最近対応した複数のジョブを順番に実行する簡単な仕組みについての要望があります。ジョブを順番に実行する仕組みそのものは元々あったのですが、だいぶ煩雑でパツとは使えず、使おうとするどうしても一手間かかる面倒くさいものだったので、それをもっと簡単にできる方法はないかという要望です。これについては幸運にも需要を私の方でも理解できていて、仕組み自体はちょっと時間があれば実装できそうだったので、実装しました（現在

の jsub コマンドで利用可能な --step 及び --stepany オプションです）。現在、この機能は要望を出した人以外にもそれなりに使われているようです。この機能には元々潜在的な需要はあったのですが、管理側でそのような潜在的な需要を読み切ることが一般的に極めて困難です。かといって、現状の受付体制のままで要望をどしどし御寄せ下さい、とアピールしても期待通りに要望が集まるわけでもなさそうです。もうちょっと気楽に要望を投げつけられ

るような仕組み、例えば登録ユーザーしか使えないけどほぼ匿名（本当に完全に匿名にすると別の問題が発生しそうなので、IPアドレスくらいはどっかに記録されるでしょうが）で要望を出せる仕組みでしょうか。できるものをできる範囲で実現して行きたいです。ただ、何か新しいモノを作っていくにはやはりちょっと漕ぎ手が足りてない、というのが実際のところですよ。

## 覧古考新25 | 2003年

最近、例の大学共同利用機関の法人化問題で振り回されている。この改革によって素晴らしい体制と組織が出来そうと言うことであれば良いが、どうもそういう明るい光が見えてこない現実の中で、それでもこれを進めて行かねばならず、しかもそれに荷担しなくてはならないのは誠に心外であり健康に良くない。誰かが、「山本五十六の心境ですね」と冷やかしてくれたが、残念ながら「言い得て妙」である。多大な時間と労力を費やしているが、日本の基礎学術は100年の計で本当に大丈夫なのだろうかと心底心配になる。国に資金が十分でない時の良くない改革は最悪である。

.....

よく、日本の組織改革は「着せ替え人形だ」と言われる。体裁だけが違って、中身が変わらないのである。今回は、「着せ替えで中身が腐らない様に気を付けなくてはならない」と言う皮肉な状況にあるのではないだろうか!! 少しでも良い方向に向かう様に（多大な!）努力をしなくてはならない。

何はともあれ、基礎科学、基礎学術の意義とその重要性について我々自身が改めて深く考え直さなくてはならない。社会に対する責任を自覚すると共に、意義と重要性を強く訴えていく努力をする必要がある。「基礎科学は子供の様なものだから大事にしなくてはいけない」という表現がある。「将来有用な大人に育って行く者がいるからである」と言うこともあるが、むしろ、「子供それ自身に存在意義がある」からなのである。基礎科学、基礎学術もそれ自身に存在意義があるのである。ただ、我々が殻に閉じこもって奢っているだけでは許されない。説明責任は果たさなくてはならない。しかし、何と言っても、最大の問題は、科学者自身が気概を失い、プラグマティズムと応用研究重視の風潮に流されているのではないかと危惧されることである。

分子研レターズ No.48 「レターズ：最近思うこと」(2003年)

中村 宏樹（分子科学研究所教授）