

# GPU による汎用計算について

## Getting Started with GPGPU

菅 準 一

# GPU による汎用計算について

菅 準 一

CPU (Central Processing Unit) の処理能力は飛躍的に伸びたが、解きたい問題によっては、その処理時間のさらなる短縮を望む事態が生じる。場合によっては、プログラムの一部に並列処理を用いることで処理時間を短縮できる。しかしクロック数が高くてコア数とスレッド数の多い CPU の導入には、多額の費用がかかり実現が困難な状況がある。そこで、数値計算の用途に BTO (Build To Order) したものではないが、マルチコア CPU のパソコンに、CPU と比較すると安価な (GPGPU 専用の Tesla ではない) GPU (Graphics Processing Units) を導入することで、処理時間の短縮を図る GPGPU (General-Purpose computing on Graphics Processing Units) つまり GPU による汎用計算について考えてみる。

キーワード：GPGPU, CUDA, Python, Ubuntu, Docker

## 目次

1. はじめに
2. Python で GPGPU
3. OpenSSH で Windows ノートから GPGPU
4. Ubuntu の Docker と LXD を利用する
5. Docker とはどんなもの
6. おわりに

## 1 はじめに

最近まで CPU の処理能力は飛躍的に伸びてきたが、解きたい問題によっては、その処理時間のさらなる短縮を望む事態が生じている。場合によっては、プログラムの一部に並列処理を用いることで処理時間を短縮できる。しかしクロック数が高くてコア数とスレッド数の多い CPU の導入には、多額の費用がかかり実現が困難な状況がある。そこで、数値計算の

用途に BTO したものではないが、マルチコア CPU のパソコンに、CPU と比較すると安価な (GPGPU 専用の Tesla ではない) GPU を導入することで、処理時間の短縮を図る GPGPU つまり GPU による汎用計算について考えてみる。

GPGPU を実現するためには、GPGPU に対応する GPU が必要である。GPU を積んだグラフィックカードの 2 大メーカーである AMD 社と NVIDIA 社の製品であれば対応している。ここでは比較的安価なグラフィックカードである NVIDIA 社の GeForce 1060 (3G) を使って Ubuntu 16.10 上で GPGPU を体験してみる。今回用いた PC は Windows 10 Home がプレインストールされ販売されているもので、メインメモリを 64G に増やしている。以前メモリが 16G の Windows 10 パソコンで Python を使って数値計算しているとき、ユーザーメモリを使い切ってしまう、パソコンがカクカクとしか反応しなくなったのでメモリを増強している。また Ubuntu 16.10 を後からインストールして、Windows10 と Ubuntu16.10 をデュアルブートできるようにしている。Ubuntu 16.10 をインストールする際に、Windows 10 のファストブート機能は無効化して、起動時に Del キーの連打で UEFI を呼び出せるようにしている。ただ、UEFI のセキュアブートは残したままであった。

また、なぜ長期的にサポートされる Ubuntu 16.04 LTS でないかという点、Ubuntu 16.04 LTS 日本語 Remix の iso イメージのグラフィックドライバが古くて GeForce 1060 に対応していないこと、マザーボードのチップセットが Intel X99 で CPU ソケットが LGA2011-3 であり、グラフィック機能のない CPU が載っていること。ディスプレイへの表示も外付けの GeForce 1060 ボードによるので、Ubuntu 16.04 LTS 日本語版 Remix iso イメージから作った Live USB で起動しても、ディスプレイに何も表示されない。Ubuntu 16.10 は、サポート期間が 9 か月と短い点、Ubuntu 16.10 日本語版 Remix の iso イメージから作成した Live USB から起動すると、ディスプレイに表示できたので、Ubuntu 16.10 をインストールして用いている<sup>1)</sup>。

---

1) ただデュアルブートにしなくても、Windows 10 HOME 64 上でも Ubuntu を動かすことができる。一つは VirtualBox など Virtual 環境に Ubuntu をインストールして使う。Network インストールできる通信環境が整い、また Virtual マシンに割ける PC のメモリやグラフィックメモリ、ディスク容量の増加に伴い快適に動くようになった。もう一つは Windows10 Anniversary Update に含まれる形で提供されるようになった WSL (Windows Subsystem for Linux) だ。WSL を有効にして、セキュリティレベルを開発者モードに設定することで、Ubuntu がダウンロードされて、WSL として Ubuntu を動かすことができる。最初は CLI で面食らうと思うが、X Windows System や軽量デスクトップをインストールすると GUI 環境が構築できて、FireFox を起動できる。ただ日本語が入力できないなど面白いことが起こり元々は英語圏で開発されたものであることを痛感することになる。上記の Virtual 環境や WSL のいずれにおいても、GPGPU 環境を構築して良いパフォーマンスを生み出せないと感じた。そこで、Mac ユーザーや Linux ユーザーにとっては馴染みの Windows とのデュアルブート環境の構築である。昨年 Microsoft 社が Linux Foundation の法人会員として加入したことによって、Linux の勢力圏に大きな変化が起こるかもしれない。

GPGPU するために NVIDIA のドライバと CUDA をインストールする必要がある<sup>2)</sup>。Ubuntu を通常インストールすると使える NVIDIA の GPU 向けオープン・ソース・ソフトウェアのグラフィックス・ドライバ(nouveau)は、ディスプレイに表示させる本来の機能とは別の機能である GPGPU に対応するようには作られていない。そこで、NVIDIA が Linux 向けに提供しているグラフィックス・ドライバと GPU に汎用計算を行わせる開発環境「CUDA」をインストールする必要がある。UEFI のセキュアブートをオフにしておき、標準の nouveau ドライバを無効化して、DKMS パッケージをインストールして、Ubuntu の公式リポジトリから GeForce 1060 対応ドライバをインストールする<sup>3)</sup>。ドライバと CUDA Toolkit がインストールできたら、GPU が認識されているか確認のためにターミナルで

```
nvdiia-smi
```

コマンドを実行して、図1のような表示が出ると良い。GPU の ID は 0, 名前は GeForce GTX 106..., 温度は 27C などが表示されている。PCI Bus ID が 0000:03:00.0 メモリが 3010MiB (メビバイト)

```
kt@kt-X99-504A:~$ nvidia-smi
Fri Apr 14 11:45:56 2017

+-----+
| NVIDIA-SMI 375.39                Driver Version: 375.39          |
+-----+-----+
| GPU Name      Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf    Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Compute M. |
+-----+-----+-----+-----+
|   0  GeForce GTX 106...    Off   | 0000:03:00.0  On    |           N/A       |
| 38%   27C   P8      6W / 120W | 206MiB / 3010MiB |      0%      Default |
+-----+-----+-----+-----+

+-----+
| Processes:                        GPU Memory |
| GPU       PID    Type   Process name      Usage   |
+-----+-----+-----+-----+
|   0       1506    G     /usr/lib/xorg/Xorg  162MiB |
|   0       3256    G     /usr/bin/compliz   39MiB  |
|   0       4669    G     /usr/lib/firefox/firefox  1MiB   |
+-----+-----+-----+-----+
```

さらに、CUDA Toolkit がうまくインストールできているか確認するため、CUDA のサンプルが動くかどうか確認する。標準的なインストールだと CUDA のサンプルは /usr/local/cuda-8.0/NVIDIA CUDA-8.0 Samples ディレクトリにある。1 Utilities/deviceQuery ディレクトリに移動して

- 
- 2) CUDA ではなくて OpenCL を使うという選択肢もある。その他にも新しい技術が登場しているので、それらを使うのもいいと思う。今回は CUDA を使ってみる。
  - 3) うかつにも、このような手順を踏まずにインストールしようとして、インストールが途中で止まっていることを確認せずに電源を落としてしまって、再起動したら GUI でログインできなくなりました。うっかりして、そんなことになったら、Ctrl+Alt+F1 で GUI を抜けて CLI モードに入り、ログインして NVIDIA 関連のファイルを purge してやり直すとい。私の場合は Windows 10 とデュアルブートなので、再起動して Grub の画面で Windows10 を起動し、必要な情報を手に入れたが、ノートパソコンやスマホを手元に置いておくとい。また、Windows10 も起動しない最悪の場合を考えて、LiveUSB も手元に用意しておくとい。

sudo make

コマンドで deviceQuery のバイナリを作り

./deviceQuery

を実行すると、以下のような表示がされる<sup>4)</sup>。

```

kt@kt-X99-504A: /usr/local/cuda-8.0/NVIDIA_CUDA-8.0_Samples/1_Utils/1/deviceQuery$ ./deviceQuery
./deviceQuery Starting...

  CUDA Device Query (Runtime API) version (CUDA static linking)

Detected 1 CUDA Capable device(s)

Device 0: "GeForce GTX 1060 3GB"
  CUDA Driver Version / Runtime Version      8.0 / 8.0
  CUDA Capability Major/Minor version number: 6.1
  Total amount of global memory:             3011 MBytes (3157131264 bytes)
  ( 9) Multiprocessors, (128) CUDA Cores/MP: 1152 CUDA Cores
  GPU Max Clock rate:                        1709 MHz (1.71 GHz)
  Memory Clock rate:                         4004 Mhz
  Memory Bus Width:                           192-bit
  L2 Cache Size:                              1572864 bytes
  Maximum Texture Dimension Size (x,y,z)     1D=(131072), 2D=(131072, 65536), 3D=(16384, 16384, 16384)
  Maximum Layered 1D Texture Size, (num) layers 1D=(32768), 2048 Layers
  Maximum Layered 2D Texture Size, (num) layers 2D=(32768, 32768), 2048 Layers
  Total amount of constant memory:           65536 bytes
  Total amount of shared memory per block:   49152 bytes
  Total number of registers available per block: 65536
  Warp size:                                  32
  Maximum number of threads per multiprocessor: 2048
  Maximum number of threads per block:       1024
  Max dimension size of a thread block (x,y,z): (1024, 1024, 64)
  Max dimension size of a grid size (x,y,z): (2147483647, 65535, 65535)
  Maximum memory pitch:                       2147483647 bytes
  Texture alignment:                           512 bytes
  Concurrent copy and kernel execution:       Yes with 2 copy engine(s)
  Run time limit on kernels:                   Yes
  Integrated GPU sharing Host Memory:         No
  Support host page-locked memory mapping:    Yes
  Alignment requirement for Surfaces:         Yes
  Device has ECC support:                       Disabled
  Device supports Unified Addressing (UVA):    Yes
  Device PCI Domain ID / Bus ID / location ID: 0 / 3 / 0
  Compute Mode:
    < Default (multiple host threads can use ::cudaSetDevice()) with device simultaneously >

deviceQuery, CUDA Driver = CUDART, CUDA Driver Version = 8.0, CUDA Runtime Version = 8.0, NumDevs = 1,
Device0 = GeForce GTX 1060 3GB
Result = PASS

```

## 2 Python で GPGPU

CUDA を使って GPGPU のためのプログラムを直接書いてもよいが、Python で GPGPU をすることを考える。そこで、Python の環境を整えるために、Continuum Analytics の Anaconda をインストールする<sup>5)</sup>。Python で CUDA のプログラムを書けるものとして、まず頭に浮かぶの

- 4) ただ、残念なことにその他のサンプルでは、「gcc5 以上には対応していない」というエラーがでる。とりあえずインストールされている gcc6 を downgrade するために gcc-4.9 と g++-4.9 をダウンロード・インストールして、それぞれを /usr/local/cuda/bin/gcc と /usr/local/cuda/bin/g++ にリンクさせると make できる。
- 5) Anaconda のダウンロードページの中ほどにある Download for Linux タブをクリックして左側に書かれた 3 つの指示に従ってインストールするとよい。右側にある Python 3.6 version の 64-BIT INSTALLER (緑色) をクリックするとダウンロードが始まる。ダウンロードしたファイルの MD5 or SHA-256 でファイルが壊れたり改ざんされていないかチェックする。terminal で bash Anaconda3-4.3.1-Linux-x86\_64.sh をタイプして指示に従う。ファイル名が長いので Tab キーによる補完機能を使うと間違いがない。terminal で python と打ち込んで Python と Anaconda のバージョンが表示されればインストールが成功したことになる。Anaconda をインストールすると、package manager と environment manager の役割を果たす conda も一緒にインストールされているので、まずは conda 自体を conda update conda で update する。

はPyCUDAだろう。しかしながら、ネットで少し調べてみると、CuPyというものがあることが分かった。これはNumpyのGPU版という位置づけのもので、ChainerというDeep Learning用のPythonパッケージでGPUバックエンドとして動いている。Chainerをインストールすると自動的にインストールされる。そこで、まずChainerがインストールされているかどうか確かめる。

`python -c "import chainer"` をターミナル上で実行するとエラーが出ることから、インストールされていないことが分かる。そこで、condaでインストールできるか確かめる。conda search chainer を実行しても表示されないの、condaではインストールできないことが分かる。condaコマンドでインストールできない場合はpipを使ってインストールすることになる。

```
pip install chainer
```

インストールがうまくいったので、MNIST exampleを試す。

```
python train_mnist.py -g=0
```

でgpuを使う設定で計算させると、53秒ほどで計算が終了した。CPUだけの場合は225秒ほどなので、CPUだけでも待てないほどの処理時間ではないけれど、GPUを使うと4倍くらい計算が早くなっている。GPUを高性能なものに替えると処理をさらに高速化可能であると思われる<sup>6)</sup>。

### 3 OpenSSHでWindowsノートからGPGPU

GPGPUのメリットはわかるが、普段使っているWindowsノートにはGPUは入っていない。普段持ち歩いているWindowsノートでGPGPUできないだろうか。一つの解決法としてGPUがついているデスクトップPCに、ユーザーとして登録してもらい、インターネットを通じてログインしてWindowsノートでデスクトップPCを操作することが考えられる。インターネットを通じて離れた場所からPCを使うには、データが盗み見られても大丈夫なようにしておく必要がある。そのためには、通信内容を暗号化することが考えられるが、一般に普及しているのはOpenSSHを用いた暗号通信である。そこでGPUのついたPC(ホスト)にUbuntu用のSSHサーバーをインストールし、WindowsノートにWindows用OpenSSHをインストールしてSSHクライアントとしてホストにログインできるようにする。Ubuntu PCをSSHサーバにして、SSHクライアントのWindowsノートで作成した公開鍵を

---

6) AnacondaでインストールしたPythonのバージョンは3.6でなぜかNVIDIAのcuDNN(Deep Neural Network library)がenableにならないエラーが出たので、再インストールしたがうまくいかなかった。そこで一度頭を切り替えて、condaのenvironment management機能を使って、Chainer用のPython 3.5.1の環境を作成してインストールすると、`import cupy.cudnn`のエラーが消えてcuDNNのenableエラーも表示されなくなった。しかし残念なことに処理速度に変わりはない。cuDNNをダウンロードするにはNVIDIAのAccelerated Computing Programのメンバーになることが必要である。

受け取り登録する簡単な手順は以下の通りである。

1. Ubuntu PC に SSH をインストール

```
sudo_apt-get_install_ssh
```

さらに、この PC の ip アドレスを確認する

```
ifconfig
```

2. Windows ノートに OpenSSH をインストールする。

OpenSSH のダウンロード・ページからダウンロードした 64bit 用の OpenSSH-Win64.zip ファイルを適当な場所 (ディレクトリ) に展開する。ディレクトリに Path を通す。

3. 暗号化に使う鍵を Windows ノートで作成する

```
ssh-keygen コマンド (. \ssh-keygen.exe) を使って
```

```
public/private rsa key を作る。
```

作ったキーを保管するファイルを尋ねられるが、デフォルトのままでもよいので、そのまま Enter キーを押す、更に passphrase を入力するよう求められるので、入力する。この passphrase が login に必要なので書き留めておく。すると、デフォルトの場所に id\_rsa (秘密鍵) と id\_rsa.pub (公開鍵) のファイルが作られる。

4. 公開鍵ファイル id\_rsa.pub をホストに送って認可鍵として保管してもらう

この公開鍵をサーバーに登録すれば暗号化通信ができる。Windows ノート (クライアント) から Ubuntu PC (ホスト) に向けて公開鍵を scp コマンドで送る

```
scp_id_rsa.pub_server-ip-address:~/
```

5. ホストの Ubuntu PC では、送られてきた公開鍵ファイル id\_rsa.pub を認可鍵ファイル authorized.keys に追加する

鍵を保管するディレクトリを作る (sudo で一時的に root になる)

```
sudo_mkdir_._ssh
```

```
sudo_cat_id_rsa.pub_>>_._ssh/authorized.keys
```

これで Windows ノート (クライアント) から SSH を用いた暗号化通信ができて、ホストの GPGPU を体験できる。

6. Windows ノートのコマンドプロンプトで

```
ssh_user@host-ip-adress
```

鍵の作成時に passphrase を設定したので、そのとき設定した passphrase を入力する

user として Ubuntu PC のキーボードから直接ログインするときの user password ではないので注意すること。また Ubuntu PC の文字コードが UTF-8 なので、Windows のコマンドプロンプトからログインすると文字化けする場合がある。文字化けしたら、一度 exit で接続を切って、

chcp 65001

を実行して、コマンドプロンプトの文字コードを UTF-8 に変更してから、再接続するとよい。

#### 4 Ubuntu の Docker と LXD を利用する

Mac を使っていると 2001 年の Mac OS X 発売以来毎年 OS の新しいバージョンが発売され、新しくなるたびに新しい OS を購入してインストールするというのが続いていて、ソフトウェアによると新しいバージョンの OS では動かないという事態が発生する。そこで USB の外付け HDD に、PC 本体とは異なるバージョンをインストールして、起動ディスクを使い分けることで対応してきた。いまも最新の macOS (Sierra) を外付けの SSD にインストールして、MAMP で moodle の最新版を使い STACK を利用する試験運用に使っている。名刺入れより小さい 256G の SSD で試験運用環境をポケットに入れて持ち歩くことができる。ただ、35G とサイズが大きいのが問題だ。MAMP に加えて moodle や STACK, Maxima, Gnuplot, Xterm など、必要なアプリをダウンロードして設定したものだけに絞って軽量化して、適当なバージョンの OS が起動しているサーバやデスクトップに簡単に受け渡し配備できないだろうか。

一台の物理マシン（ハードウェア）上に Virtual マシン（ファイル）を複数インストールして、macOS X や Ubuntu, Windows など異なる OS を使い分けたり、同じ OS でも違うバージョンを使ったり、64 ビットの物理マシン上に 32 ビットの仮想マシンを動かすことができるなど、仮想環境は一台の物理マシンを有効に使うために重要な役割を担ってきた。そのような仮想化技術は、ホスト OS 型とハイパーバイザー型、に分けることができる。ホスト OS 型は例えば Mac OS X（ホスト OS）が動いているマシンに VirtualBox などの仮想化ソフトをインストールして仮想マシンを立ち上げて、たとえば Windows や Ubuntu などの OS（ゲスト OS）を動かす。このときゲスト OS を動かすためには、物理マシンの上にホスト OS と仮想化ソフトの 2 つの層が必要である。他方ハイパーバイザー型の仮想化技術では、複数の OS を動かすためのハイパーバイザーと呼ばれる仮想化のための 1 層が物理マシンの上に作られる。最近注目を浴びているコンテナ型の仮想化は OS とコンテナ管理ソフトウェアの上にコンテナによる仮想環境を提供する。Docker 社の Docker が最近注目されている。また、LXD は LXC が提供する API を利用したコンテナ管理コマンドやサービスで、ホストとは異なる OS をインストールしたコンテナをいくつか作って、様々な運用検証環境を整備できる。

Ubuntu のプログラムをインストールして動かすためには、Ubuntu のカーネルのバージョンとプログラムのバージョンがうまく合わないと、インストールにすら失敗してしまう。う

かつにシステムをバージョンアップしないことに注意しないとイケないが、最新のデバイスを買ってしまって、最新のバージョンの Ubuntu をインストールしてしまうと様々な困難に遭遇することになる。とくに NVIDIA のデバイスドライバが関わると、一苦労する、上で記したように、NVIDIA のドライバと CUDA, cuDNN をインストールして、Chainer を動かすことに成功したが、さらに Tensorflow を、色々な依存関係をクリアするように調整して、インストールして試してみる時間的余裕がない。そこで、Tensorflow が動く組み合わせをパッケージにしてくれたものがあれば、ありがたいと思ってネットを調べたら nvidia-docker なるものがあり、それを使って Tensorflow を簡単に利用できるようなので試してみた<sup>7)</sup>。簡単そうでもやはり手間が必要である。Github の nvidia-docker のリポジトリに行き Installation のページを見ると、次のような記述がある。

The list of prerequisites for running nvidia-docker is described below.

1. GNU/Linux x86 64 with kernel version > 3:10
2. Docker >= 1:9 (official docker-engine only)
3. NVIDIA GPU with Architecture > Fermi (2.1)<sup>8)</sup>
4. NVIDIA drivers >= 340:29 with binary nvidia-modprobe

Your driver version might limit your CUDA capabilities (see CUDA requirements)

さらに CUDA toolkit version 8.0 を利用する場合には、Driver version が 367.48 以上でないといけないと書いてあった。

kernel のバージョンは `uname -rv` コマンドで表示できる。バージョンは 4.8 なので十分である。

NVIDIA GPU は GeForce 1060 (3G) で Pascal Architecture なので十分である。NVIDIA ドライバーのバージョンは `nvidia-smi` コマンドで確認すると、375.39 であり十分である。残るは

---

7) nvidia-docker については、<https://github.com/NVIDIA/nvidia-docker/wiki/nvidia-docker> に説明がある。nvidia-docker is a thin wrapper on top of docker and act as a drop-in replacement for the docker command line interface. This binary is provided as a convenience to automatically detect and setup GPU containers leveraging NVIDIA hardware. Internally, nvidia-docker calls docker and relies on the NVIDIA Docker plugin to discover driver files and GPU devices. Note that nvidia-docker only modifies the behavior of the run and create Docker commands. All the other commands are just pass-through to the docker command line interface. As a result, you can't execute GPU code when building a Docker image. If the nvidia-docker-plugin is installed on your host and running locally, no additional step is needed. nvidia-docker will perform what is necessary by querying the plugin when containers using NVIDIA GPUs need to be launched. したがって、nvidia-docker は docker command line interface の完全互換であり、docker コマンドの run と create だけに修正を加える。また nvidia-docker から GPU の情報を得るための nvidia-smi コマンドを使うためには、ドライバファイルや GPU を探してくれる nvidia-docker-plugin が起動していないといけないことがわかる。

8) NVIDIA の architecture には著名な科学者の名前がつけられてきた。Tesla, Fermi, Kepler, Maxwell そして昨年市場に投入された Pascal, ようやくこの3月に姿をみせた次世代の Volta という具合だ。括弧の中の数字 2.1 Compute Capability のバージョンである。Pascal architecture の GeForce GTX 1060 の Compute Capability は 6.1

Docker のバージョンが 1.9 以上でしかも official docker-engine only ということで、条件を満たす official docker-engine をインストールする。

Ubuntu の公式リポジトリで提供している docker.io ではなくて、Docker 社のリポジトリの docker-engine でないといけないようだが、見当たらない。Docker 社のページにある無料の Docker CommunityEdition (Docker CE) をインストールすればよいことが分かったので、ダウンロードしてインストールする。次に、プラグインである nvidia-docker をインストールしようとして、Quick start の例にならって `dpkg -i` を試してみたが、依存するファイル `stvrc` と `file-rc` がないというエラーでインストールが止まってしまう。そこで、ソースファイルをダウンロードして `make install` するとエラーメッセージをはきながら、とにかくインストールができた。

```
sudo_docker_images
```

 で表示させると、`nvidia-docker` が build されている。

そこで、`nvidia-docker` が使えるか試しに、GPU の情報を `nvidia-smi` を使って表示させてみる。

```
sudo_nvidia-docker_run_nvidia/cuda_nvidia-smi
```

`plugin not found` というエラーで止まってしまう。そこで、手動でプラグインを起動することにする。

```
sudo_nvidia-docker-plugin
```

これでプラグインがアクティベートされる。

プロセスが動いているので、別のターミナルを立ち上げて

```
sudo_nvidia-docker_run_nvidia/cuda_nvidia-smi
```

と入力すると GPU のドライバのバージョンが 375.39 であることなどが表示され、ドライバのインストールと GPU の認識がうまくいっていることがわかる。

以上で、`docker-ce` と `nvidia-docker` がインストールできたので、つぎは `nvidia-docker` を使うとどれくらい容易く TensorFlow をインストールして運用できるか確かめてみる。

Tensorflow の github の Docker ディレクトリや Docker Hub から tensorflow のイメージ

`gcr.io/tensorflow/tensorflow:xxx-gpu` (TensorFlow with all dependencies and support for NVidia CUDA) をためす。ここで、`gpu` がついていることを確認する。`gpu` がないと、CPU だけを使うものであり GPU は使えない。`xxx` はバージョンである。

```
sudo_nvidia-docker_run_--itd_--name=tensorflow-xxx-gpu_-p_8888:8888_
-p_6006:6006_gcr.io/tensorflow/tensorflow:xxx-gpu
```

うまくいくと、`sudo_docker_images` や `sudo_docker_ps_-a` で一覧表示させる (`-a` オプションをつけないと停止中の docker は表示されない) と、`gcr.io/tensorflow/tensorflow:xxx-gpu` のイメージファイルをもつ `tensorflow-xxx-gpu` という名前の docker コンテナが作られている。

そこで、この docker コンテナに入って作業をする。

```
sudo_nvidia-docker_exec_-it_tensorflow-xxx-gpu_bash
```

dockerの中に入ってbashでインタラクティブ(-it)に作業ができ、プロンプトの表示が次のように変化する(ホストの時の\$から#に代わる)<sup>9)</sup>。

```
root@docker-ID:/notebooks#
```

Githubからtensorflowをクローンする。このとき、gitが無いというエラーが出たらgitをインストールするとよい。

```
git_clone_-b_r0.10_--single-branch_--recurse-submodules_https://
git.com/tensorflow/tensorflow.git
```

うまくクローンが作れたら

カレントディレクトリの下にtensorflowというディレクトリが作られているので、lsコマンドで確認するとよい。

cdコマンドでtensorflow/tensorflow/models/image/cifar10/まで移動する。このとき、Tabキーの入力補完機能を使うと間違いが少ない。

```
python_cifar10_multi_gpu_train.py_--num_gpus=1_--Vmax_steps=1000
```

を実行する。

以下の図のように、cudaのライブラリが読み込まれていることがわかる。簡単すぎるので当惑する。

```
root@98e945d66ef6:/notebooks/tensorflow/tensorflow/models/image/cifar10
init .py cifar10 input.py cifar10 train.py
root@98e945d66ef6:/notebooks/tensorflow/tensorflow/models/image/cifar10# python cifar10_multi_gpu_train.py
--num_gpus=1 --max_steps=1000
I tensorflow/stream_executor/dso_loader.cc:108] successfully opened CUDA library libcublas.so locally
I tensorflow/stream_executor/dso_loader.cc:108] successfully opened CUDA library libcudnn.so locally
I tensorflow/stream_executor/dso_loader.cc:108] successfully opened CUDA library libcufft.so locally
I tensorflow/stream_executor/dso_loader.cc:108] successfully opened CUDA library libcuda.so.1 locally
I tensorflow/stream_executor/dso_loader.cc:108] successfully opened CUDA library libcurand.so locally
Filling queue with 20000 CIFAR images before starting to train. This will take a few minutes.
I tensorflow/core/common_runtime/gpu/gpu_init.cc:102] Found device 0 with properties:
name: GeForce GTX 1060 3GB
major: 6 minor: 1 MemoryClockRate (GHz) 1.7085
pciBusID 0000:03:00.0
total memory: 2.94GiB
free memory: 2.71GiB
I tensorflow/core/common_runtime/gpu/gpu_init.cc:126] DMA: 0
I tensorflow/core/common_runtime/gpu/gpu_init.cc:136] 0: Y
I tensorflow/core/common_runtime/gpu/gpu_device.cc:838] Creating TensorFlow device (/gpu:0) -> (device: 0,
name: GeForce GTX 1060 3GB, pci bus id: 0000:03:00.0)
2017-05-15 09:19:24.945228: step 0, loss = 4.68 (4.3 examples/sec; 29.670 sec/batch)
2017-05-15 09:19:26.109095: step 10, loss = 4.60 (1928.0 examples/sec; 0.066 sec/batch)
2017-05-15 09:19:26.797880: step 20, loss = 4.55 (1829.3 examples/sec; 0.070 sec/batch)
2017-05-15 09:19:27.478645: step 30, loss = 4.42 (1914.0 examples/sec; 0.067 sec/batch)
2017-05-15 09:19:28.165126: step 40, loss = 4.29 (1840.6 examples/sec; 0.070 sec/batch)
```

9) 一度シャットダウンして、翌日に再現しようとするときは、sudo nvidia-docker-pluginでプラグインを立ち上げておいて、別のターミナルでsudo docker start tensorflow-xxx-gpuでdockerをスタートさせてからでないこと、sudo nvidia-docker execコマンドでdockerの中には入れない症状が現れる場合がある。そこで、ホームディレクトリに上記のコマンドを列挙したシェルスクリプトファイルを作る。1行目にシバンと呼ばれるおまじないを書いておく。#!/bin/bash 次の行にはsudo nvidia-docker-plugin & 行の最後の&は、起動したプロセスをバックグラウンドで実行させるためのもので、次の行のコマンドsudo docker start tensorflow-xxx-gpuに進み、最後の行でsudo docker exec -it tensorflow-xxx-gpu bashを実行させる。ファイルを保存したら、ターミナルでbash nvidia-docker.shと実行して、tensorflow-xxx-gpu docker内のnotebooksディレクトリに入っていることを確認するとよい。

このとき、プログラム実行中の GPU の状況を `nvidia-smi` コマンドで確認したものが、以下の図である。

```

kt@kt-X99-02A: ~
kt@kt-X99-02A:~$ nvidia-smi
Mon May 15 18:20:29 2017

+-----+
| NVIDIA-SMI 375.39                Driver Version: 375.39          |
+-----+-----+
| GPU Name      Persistence-M| Bus-Id        Disp.A | Volatile Uncorr. ECC |
| Fan  Temp  Perf  Pwr:Usage/Cap|      Memory-Usage | GPU-Util  Compute M. |
+-----+-----+-----+-----+
| 0   GeForce GTX 106...   Off | 0000:03:00.0  On    |      47%      N/A   |
| 38%   44C   P2     37W / 120W | 2792MiB / 3010MiB |           | Default |
+-----+-----+-----+-----+

Processes:
+-----+-----+-----+-----+
| GPU    PID  Type  Process name      GPU Memory |
| Usage |
+-----+-----+-----+-----+
| 0      2762  G    /usr/lib/xorg/Xorg  139MiB |
| 0      4910  G    compiz             37MiB  |
| 0      7402  C    python             2611MiB|
+-----+-----+-----+-----+

kt@kt-X99-02A:~$

```

GPU の温度は 44 度で、消費電力は 37W である。GPU のメモリ 3010MiB のうち 2792MiB を使っていて、python の Process が 2611MiB 利用して計算等を行っている。Volatile GPU-Util が 47% と表示されており、メモリがボトルネックになっているのかもしれない。GeForce GTX 1060 (3G) が GeForce GTX 1060 (6G) の廉価版であることが何らかの影響を与えているのだろうか。この点は調査を要する。

ここで、`tensorflow-xxx-gpu` を立ち上げたときに、コンテナの中で `jupyter notebook` サーバが起動して 8888 ポートを使っている。このコンテナの 8888 ポートとホストの 8888 ポートをマップしている。ホストからコンテナに Firefox でアクセスするために、Firefox を立ち上げて `http://localhost:8888` と入力すると `jupyter notebook` の Home が開く。1 `hello tensorflow.ipynb` や 3 `mnist from scratch` などを開いてセルを実行させると、エラーもなく実行できる。

#### 4.1 TensorBoard を使う

コンテナの `tensorflow` を用いて簡単な計算をして、その結果を TensorBoard で表示できる。

TensorBoard は TensorFlow とともにインストールされて、パスが通っているので、`docker` の中で `tensorboard` コマンドで起動すると、

```
Starting TensorBoard 41 on port 6006
```

```
(You can navigate to http://172.17.0.2:6006)
```

と表示される。

`docker` の外に出て、ホスト上で `firefox` などのブラウザを起動して、URL に `172.17.0.2:6006`

と入力すると TensorBoard の画面が表示される。

## 5 Docker とはどんなもの

Docker についてももう少し知りたいと Docker 社のページをのぞいてみると、ソフトウェアのコンテナのプラットフォームだという。「僕の PC では動くけど」(“works on my machine”)と言われても困るケースがある。プログラムを明示的な合意に基づいて共同開発する場合以外にも、コンピュータプログラムを用いた先行研究を検証したり、先行研究のプログラムを利用しようとするとき、同じ環境を構築するのに時間がかかり、結局は動かなかったということがしばしば起きる。コンテナを使うとソフトウェアを動かすのに必要なすべてをコンテナにパッケージでき、どのような環境に配備しても動くようになるらしい。VM はマシーンレベルの仮想化技術であり、VM 上に OS 全部をインストール必要とするが、コンテナはソフトウェアが動くためのライブラリと設定だけで十分なのである。インストールや設定に時間を費やさなくて済み。複雑な依存関係は Docker イメージから引っ張てくると良い。

必要に迫られて nvidia-docker を使い Tensorflow 等の GPU を利用する Python のパッケージを動くようにしたかったので、Docker に手を出したから、Docker の全体像を語るほどの知識はいまはない。ただ、初心者にとって必要最小限の情報は伝えることができると思う。

Docker を支えるものは、イメージ (image)、コンテナ (container)、サービス (service)、スタック (stack) とスワム (swarm) である。イメージからコンテナを作って、コンテナで作業をする。イメージは既製品と手作りの二つがある。既製品は Docker Hub をはじめとする様々なレジストリ (Registry) に保存されている色々なレジストリの Image ファイルとして提供されている。手作りは Dockerfile という名のファイルに必要事項を書き、必要なものを同じディレクトリに集めておいて作る。したがって、作業するためのディレクトリをまず mkdir コマンドで作って、そこで作業する。手作りの際も既製品をもとに作ることが多いように思える。

Docker Hub にある既製品の Image は、username/repository:tag で指定して使うことができる。docker run コマンドで元になる Image を指定してコンテナをつくる。このとき -name= でコンテナに名前をつけることができ、-p でコンテナのポートとホストのポートを対応付けること (mapping) ができ、コンテナ内で提供される Web ベースのサービスにホストの Web Browser からこのポートを通じてアクセスできるようになる。

Get started with Docker Part 1:Orientation を見ながら Docker について理解を深めてみるとよいと思う。ただ、読み進めるためには、次のものについて予備知識があったほうがよい。

1. IP addresses and Ports
2. Virtual Machines

### 3. Editing configuration files

#### 4. Basic familiarity with the ideas of code dependencies and building

#### 5. Machine resource usage terms, like CPU percentages, RAM use in bytes, etc.

これに加えて、terminalでの基本的なLinuxコマンド操作 cd, ls, cp, mkdir, cat, curl. やファイル作成・編集に用いる vi などのエディターに慣れている方がよい。

Get started Part 2: Containers を読み進んで friendlyhello イメージを作って、イメージをもとにコンテナを作り所望の結果が出た後で、Docker アカウントを作って Docker の public registry にイメージをアップするように書かれている。その部分を、ローカルマシン上にプライベートな registry を作ってそこに保存するようにしてみる。

```
$sudo_docker_run_-d_-p_5000:5000_-v_/var/opt:/var/lib/registry_registry:2.3.0
```

プライベートな registry を作成運用するためのイメージもたくさんあるが、public registry にあるイメージ registry のバージョン 2.3.0 を使ってプライベートレジストリのためのコンテナを作成する。-d でコンテナをバックグラウンドで走らせ、-p でコンテナのポート 5000 とホストのポート 5000 をマップする。-v でホストの /var/opt ディレクトリをコンテナの /var/lib/registry ディレクトリにマウントしてコンテナに加えた変更をホストの /var/opt に保存する。そうすることで、コンテナを停止すると基本的には加えた変更は消去されるのだが、ホスト上にその内容を保存することで消去から復元することができる。

先に作ったイメージ friendlyhello にローカル registry の場所である localhost:5000 を指定し、さらに repository の名前をつける。

```
$sudo_docker_tag_friendlyhello_localhost:5000/repository名/friendlyhello
```

\$sudo\_docker\_images で localhost:5000/repository 名 /friendlyhello—ができていることを確認して、このイメージを push する。

```
$sudo_docker_push_localhost:5000/repository名/friendlyhello
```

そうすると、マウントしたホストのディレクトリ /var/opt に docker/registry/v2/repositories/repository 名 /friendlyhello—ディレクトリが作られイメージが保存されている。

このプライベート registry のイメージは

```
$sudo_docker_pull_localhost:5000/repository名/friendlyhello
```

で呼び出すことができる。

Get Started, Part 4: Swarms を試してみると、複数のコンテナを使ったアプリの利用やコンテナにかかる負荷の調整などについて理解が深まる。

## 5.1 Docker for Windows について

Docker for Windows もあり今度は圧倒的シェアをもつ Windows 10 環境で試してみたくな

る。Microsoft の 64bit windows 10 pro でハイパーバイザベースの x64 向け仮想化システムである Hyper-V が動くことが必要である。VirtualBox は使えなくなるという。手元のノートパソコンで確認すると、「Windows の機能の有効化または無効化」の中に Hyper-V があって有効化して、再起動してみた。「Windows 管理ツール」に Hyper-V マネージャーが追加されている。VM 型の仮想化とは排他的なので再起動後確かに VirtualBox の Ubuntu16.04LTS が起動できなくなっている。Docker 社の Docker for Windows の Windows インストーラーパッケージ InstallDocker.msi をダウンロードしてインストールする。すぐに cmd で Docker コマンドが使える。docker version で docker のバージョンを確認すると、17.03.1-ce である。OS/Arch は Client は windows/amd64 であり、Server は linux/amd64 である。また Hyper-V マネージャーを起動すると、MobyLinuxVM という仮想マシンが稼働し 2048MB のメモリが割り当てられている。Docker for Windows をインストールした後 cmd で docker run hello-world を実行したら Hello from Docker! This message shows that your installation appears to be working correctly. と表示されて docker がうまく動いている。Windows ユーザーも docker が使えるようになった。

## 6 おわりに

スーパーコンピュータを使うほどでもないが、パソコンだと少し時間がかかるような問題を、パソコンを使って何度も試行錯誤するには、プログラムを書く時間が短く、計算時間が短い方がいい。プログラムを書く時間を短縮するには、Python のような簡易言語を使うのも一つの解決策だと思う。計算時間を短縮するための選択肢の一つにオンプレミス（自前）で GPGPU がある。ただ、高速で安全は通信が可能な環境がある場合は、オンプレミス（自前）なのかクラウドなのかという選択肢がある。今回はオンプレミスで GPGPU する環境を構築することを考えてみた。Linux の一つである Ubuntu を OS として選んで NVIDIA の GPU と CUDA, Python という組み合わせで考えた。Linux の Docker を使うとプログラムの Deployment にかかる時間を節約できることがわかった。

ただ、GeForce GTX 1060 (3G) では、すぐに GPU メモリを使い切ってしまうと計算エラーを起こすので、予算に余裕があれば、メモリの大きい GPU を 2 枚以上使えるように、マザーボードや CPU が複数の GPU に対応するように考えないといけない。CPU と GPU を効率的に冷却することも必要だろう。そうするとシステム電力にも配慮が必要であるし、交流を直流に変換するときの発熱によるロスが少ない電源ユニットを選ぶ必要が出てくる。NVIDIA のページを見ると、GeForce GTX 1060 (3G) の場合消費電力は 120W で、最小限必要な（システム）電力は 400W である。GeForce GTX 1080 Ti (11G) の場合は消費電力 250W で、最小限必要な（システム）電力は 600W ということだ。GPU を 2 枚利用する場合についても考えてみたい。ただ本稿を書くために使った PC の電源は 850W に増強しているが、

GPUを2枚使うのは慎重にしないとイケないだろう。

2016年のWindows 10 kernel with the Anniversary Updateでは、WSL (Windows Subsystem for Linux)としてUbuntuを動かすことができるようになった。Hyper-Vを有効にしてLinuxベースのDockerも利用できるようになった。さらに、もう一つWindows containers (LinuxではなくてWindowsネイティブのコンテナ) が使えるようにもなった。2017年のWindows 10 kernel with the Creators Updateでは、さらにWSLのUbuntuのバージョンが16.04になり、Linuxの主要な機能が動くようになっている。またHyper-Vを更新している。WindowsがLinuxの成果を着実に取り入れ、開発環境として使いやすくなってきている。ただ、本稿を執筆している時点では、Hyper-VはWindows 10 Pro, Enterprise, Educationに対応していて、Windows 10 Homeには対応していないので、Proにアップグレードする必要がある。Windows 10の最新の機能を使ってみたい学生は、そのような最新機能が早く反映されるWindows 10 ProがインストールされたPCを買っておくべきだろう。また本稿を書くにあたっては、まとまった書籍がなくネットの検索に頼るしかなく、無償で提供された情報に感謝します。とくにメーカーやソフトウェアの開発元のホームページを丹念に読むことを学生にもお勧めします。