フィッシャー情報行列のクロネッカー因子分解を用いた 深層ニューラルネットワークの分散学習

大友 広幸^{†1,a)} 大沢 和樹^{†1,b)} 横田 理央^{†1,c)}

概要:深層ニューラルネットワークの画像認識モデルの学習では、訓練データに対する誤差関数を最小に することを目的とする.パラメータを反復的に探索するために、ミニバッチを用いた確率的勾配降下法 (SGD)による学習が広く用いられている.近年では、大規模計算資源を活用した学習高速化のアプローチ として、巨大なミニバッチを分割して計算を分散化させる大規模分散学習が注目されているが、SGDによ る巨大なミニバッチでの学習においては、モデルの認識精度が劣化することが報告されている.一方で、 フィッシャー情報行列で定まるリーマン計量を用いて誤差関数の自然勾配を求める自然勾配法(NGD)で は、SGDより大幅に少ない反復数でモデルが学習される上に巨大なミニバッチでの学習においても、モデ ルの認識精度の劣化を抑えられるという特長を持っている.膨大な数のパラメータを持つモデルの NGD による学習では巨大なフィッシャー情報行列の逆行列計算が学習全体のボトルネックとなるため、この計 算を近似し、NGDを擬似的に実現する手法が提案されてきている.本研究ではクロネッカー因子分解を 用いてフィッシャー情報行列の逆行列を近似的に計算する学習手法(K-FAC)に注目し、K-FACによる大 規模分散学習における学習の高速化と画像認識モデルの認識精度について調査した.

Ootomo Hiroyuki $^{\dagger 1,a)}$ Oosawa Kazuki $^{\dagger 1,b)}$ Yokota Rio $^{\dagger 1,c)}$

1. はじめに

近年,画像認識をはじめとした多くの分野で深層学習を 用いた手法が広く用いられている.年々深層ニューラル ネットワークの構造は複雑さを増しており,強化学習を用 いて深層ニューラルネットワークを設計する研究 [1] も行 われている.

一方,深層ニューラルネットワークの学習においてはパ ラメータをユークリッド空間上の点と仮定した確率的勾 配降下法を用いた最適パラメータ探索が主流となってい る.これに対して Amari が考案した自然勾配法 [3] ではパ ラメータの集合がリーマン多様体上にあるとし,リーマン 計量であるフィッシャー情報行列の逆行列を用いて最適 パラメータを探索する.フィッシャー情報行列の大きさは ニューラルネットワークの大きさに依存し,ニューラル

情報処理学会

1

- Presently with Tokyo Institute of Technology
- ^{a)} ootomo.h@rio.gsic.titech.ac.jp
- $^{\rm b)} \quad {\rm oosawak} @rio.gsic.titech.ac.jp \\$
- $^{\rm c)} \quad {\rm rioyokota@gsic.titech.ac.jp}$

ネットワークの大きさが大きくなっている昨今,フィッ シャー情報行列を精度良く求めることはメモリ容量や計算 量の観点から困難である.そのためフィッシャー情報行列 の逆行列を近似的に求める様々な手法が提案されている.

本研究では Martens らが提案したクロネッカー因子分解 を用いた近似手法 K-FAC[4] を用いた深層ニューラルネッ トワークの分散学習において、ミニバッチサイズを変化さ せた場合や分散数を変化させた場合、フィッシャー情報行 列の逆行列計算を行うノードの数を変化させた場合での精 度と収束の傾向を調査した.

2. 関連研究

2.1 確率的勾配降下法

ニューラルネットワークでは入力 $x \in \mathbb{R}^{d}$ に対し出力 $y = f(x) \in \mathbb{R}^{d'}$ を計算する.例えば3層ニューラルネット ワークであれば,入力層と中間層1層と出力層からなり, 入力層と中間層及び中間層と出力層の間で計算が行われ る.l+1層ニューラルネットワークを考える.i層目では 入力 $a_{i-1} \in \mathbb{R}^{d_{i-1}}$ に対し出力 $a_i \in \mathbb{R}^{d_i}$ を計算する.この 計算は重み行列 $W'_i \in \mathbb{R}^{d_i \times d_{i-1}}$,バイアス $b_i \in \mathbb{R}^{d_i}$,非線 形な活性化関数 ϕ_i を用いて

IPSJ, Chiyoda, Tokyo 101-0062, Japan ^{†1} 現在,東京工業大学

$$s_i = W_i' a_{i-1} + b_i \tag{1}$$

$$a_i = \phi_i(s_i) \tag{2}$$

と書ける. ただし, $a_0 = x, a_l = y$ である. 本論文では便 宜上 $\bar{a}_i = \begin{bmatrix} 1 & a_i^T \end{bmatrix}^T$, $W_i = \begin{bmatrix} b_i & W'_i \end{bmatrix}$ として式 1 を

$$s_i = W_i \overline{a}_{i-1} \tag{3}$$

と書く.

ニューラルネットワークの学習では出力 y と期待される 出力 $z \in \mathbb{R}^{d'}$ との不一致度を測る損失関数 L(y,z) を定義 し,その期待値を目的関数とする最小化問題を解くために 反復的にニューラルネットワークのパラメータ W を更新 する.損失関数は、例えば訓練サンプルを用いて学習する 画像認識の分野では交差エントロピーが広く用いられて いる.Wをユークリッド空間上の点と仮定すると、勾配 $\nabla L = \frac{\partial L}{\partial W}$ は L が増加する方向を向いている.したがって 反復数 t でのパラメータ $w^{(t)}$ を勾配の逆方向、すなわち $-\nabla L$ 方向へ更新すれば L を減少させることができる.し たがってパラメータの更新式は

$$w^{(t+1)} = w^{(t)} - \epsilon \nabla L \tag{4}$$

と表される. *ϵ*は学習率と呼ばれ,更新の幅を決める係数 である.学習率はミニバッチサイズなどを考慮して決め, 反復数などに応じて途中変更なパラメータである.

確率的勾配降下法 (SGD;Stochastic Gradient Descent) ではミニバッチを用いて ∇L を計算し,反復的にパラメー タを更新する. ∇L の計算方法としては,出力 y と期待さ れる出力 z の誤差を出力層から入力層にかけて伝播させる ことで ∇L を計算する誤差逆伝播法 [2] が主流である.

2.2 自然勾配法

SGDではパラメータWをユークリッド空間上の点と仮定 し勾配を計算した.これに対し自然勾配法 (NGD;Narutal Gradient Descent) ではパラメータWの集合をリーマン計 量がフィッシャー情報行列Fで定まるリーマン多様体と見 なし勾配を計算する.フィッシャー情報行列は正定値対称 行列であるため逆行列が存在し,NGDのパラメータの更 新式は

$$w^{(t+1)} = w^{(t)} - \epsilon F^{-1} \nabla L \tag{5}$$

で表される.

SGD を用いたニューラルネットワークの学習では損失 関数 L の期待値を最小化するパラメータを探索するが, 鞍 点に陥り学習が停滞することがある. NGD は SGD と比 べて鞍点に陥りにくいことや反復数に対して学習が速いこ と,大きなミニバッチで学習を行っても精度が落ちにくい という特徴がある [5].一方,昨今のニューラルネットワー クではパラメータ数が増加しており,F 及び F⁻¹ を正確 に計算することはメモリ消費量や計算量の観点から困難 である.そこで F⁻¹をスパース行列近似 [6],低ランク近 似 [7],クロネッカー因子分解を用いた近似 [4],などを用 いて計算する手法が考案されている.

2.3 K-FAC

K-FAC(Kronecker-Factored Approximate Curvature) はクロネッカー因子分解を用いた近似手法として Martens らにより考案された [4]. K-FAC では 2 段階の近似を行う ことでメモリ消費量と計算量を抑えている.

2.3.1 近似理論

行列 *M* に対し $[M]_{i,j}$ は *M* の (i,j) 成分を,ベクトル *v* に対し $[v]_i$ は *v* の *i* 番目の成分を指すとする.vec を $M \in \mathbb{R}^{m \times n}$ を $v \in \mathbb{R}^{mn}$ に $[v]_{i+nj} = [M]_{i,j}$ を満たしなが ら変換するベクトル化関数と定義する.ニューラルネット ワークの全層でのパラメータを並べたベクトル θ を

$$\theta = \left[\operatorname{vec}(W_1)^T \operatorname{vec}(W_2)^T \cdots \operatorname{vec}(W_l)^T \right]^T$$
(6)

とする. *p*(*y*|*x*, *θ*) を入力 *x* に対する出力 *y* の確率密度関数 とし, 行列 *u* に対し

$$Du = -\frac{d\log p(y|x,\theta)}{du} \tag{7}$$

を定義すると、フィッシャー情報行列の定義から

$$F = E\left[\frac{d\log p(y|x,\theta)}{d\theta}\frac{d\log p(y|x,\theta)}{d\theta}^{T}\right]$$
(8)

$$= E \left[D\theta D\theta^T \right] \tag{9}$$

である. 式6より

$$D\theta = \left[\operatorname{vec}(DW_1)^T \operatorname{vec}(DW_2)^T \cdots \operatorname{vec}(DW_l)^T\right]^T (10)$$

であるから $1 \leq i,j \leq l$ として $F_{i,j} \in \mathbb{R}^{d_{i-1}d_i \times d_{j-1}d_j}$ を

$$F_{i,j} = E[\operatorname{vec}(DW_i)\operatorname{vec}(DW_j)^T]$$
(11)

と定義すると式9は

$$F = \begin{bmatrix} F_{1,1} & F_{1,2} & \cdots & F_{1,l} \\ F_{2,1} & F_{2,2} & \cdots & F_{2,l} \\ \vdots & \vdots & \ddots & \vdots \\ F_{l,1} & F_{l,2} & \cdots & F_{l,l} \end{bmatrix}$$
(12)

と書ける. $g_i = Ds_i$ とすると, 誤差逆伝播法の過程より vec $(DW_i) = g_i \overline{a_{i-1}}^T$ である. クロネッカー積 \otimes (付録 A.1) を用いると, vec $(uv^T) = v \otimes u$ であることから式 11 は

$$F_{i,j} = E\left[\operatorname{vec}(DW_i)\operatorname{vec}(DW_j)^T\right]$$

= $E\left[(\overline{a}_{i-1} \otimes g_i)(\overline{a}_{j-1} \otimes g_j)^T\right]$
= $E\left[(\overline{a}_{i-1} \otimes g_i)(\overline{a}_{j-1}^T \otimes g_j^T)\right]$
= $E\left[\overline{a}_{i-1}\overline{a}_{j-1}^T \otimes g_ig_j^T\right]$

と書ける. ここで $F_{i,j}$ に対して 1 段階目の近似を行い

$$F_{i,j} = E\left[\overline{a}_{i-1}\overline{a}_{j-1}^T \otimes g_i g_j^T\right]$$
$$\approx E\left[\overline{a}_{i-1}\overline{a}_{j-1}^T\right] \otimes E\left[g_i g_j^T\right]$$
$$= \overline{A}_{i-1,j-1} \otimes G_{i,j} \equiv \tilde{F}_{i,j}$$

を得る.ただし $\overline{A}_{i,j} = E\left[\overline{a}_i\overline{a}_j^T\right] \in \mathbb{R}^{d_i \times d_j}, G_{i,j} = E\left[g_ig_j^T\right] \in \mathbb{R}^{d_i \times d_j}$ であり、クロネッカー因子と呼ぶ. これにより F は

$$\tilde{F} = \begin{bmatrix} \tilde{F}_{1,1} & \tilde{F}_{1,2} & \cdots & \tilde{F}_{1,l} \\ \tilde{F}_{2,1} & \tilde{F}_{2,2} & \cdots & \tilde{F}_{2,l} \\ \vdots & \vdots & \ddots & \vdots \\ \tilde{F}_{l,1} & \tilde{F}_{l,2} & \cdots & \tilde{F}_{l,l} \end{bmatrix}$$
(13)

と近似される.

2 段階目の近似として \tilde{F} の非対角ブロックを零行列とする. これにより \tilde{F} は

$$\breve{F} = \operatorname{diag}(\tilde{F}_{1,1} \ \tilde{F}_{2,2} \ \cdots \ \tilde{F}_{l,l}) \tag{14}$$

と近似される. ここでクロネッカー積の性質より

$$(A \otimes B)^{-1} = A^{-1} \otimes B^{-1} \tag{15}$$

が成り立つため,

$$\breve{F}^{-1} = \operatorname{diag}(\tilde{F}_{1,1}^{-1} \ \tilde{F}_{2,2}^{-1} \ \cdots \ \tilde{F}_{l,l}^{-1})$$

$$= \operatorname{diag}(\overline{A}_{0,0}^{-1} \otimes G_{1,1}^{-1} \ \cdots \ \overline{A}_{l-1,l-1}^{-1} \otimes G_{l,l}^{-1})$$
(16)
(17)

が従い, クロネッカー因子 $A_{i,i}^{-1}, G_{i,i}^{-1}$ を計算することで \check{F}^{-1} を計算できる.

2.3.2 数値計算

K-FACの計算機上での計算は次の3つに大分できる.

- 勾配計算 (grad):
 学習データを用いて各層における a_i, g_i 及び ∇L を計 算する.
- クロネッカー因子更新 (cov):
 クロネッカー因子 Ā_{i,i}, G_{i,i} は a_ia^T_i, g_ig^T_i の期待値で あるが,それらの真の分布は訓練データからではわか り得ない.そこで反復毎に移動平均をとることで経験 的に推定する.
- クロネッカー因子逆行列計算 (inv):
 cov で計算したクロネッカー因子 Ā_{i,i}, G_{i,i} の逆行列 Ā_{i,i}, G_{i,i} を計算する.

3. 実験

実験設定

実験には図1に示すニューラルネットワークを用い, CIFAR-10*1を学習データとしてを用いた.また,実験を 行うにあたって次の3種類のノードを用意し,図2に示す





図 2 同期型データ並列

同期型データ並列とした.

- パラメータサーバー [12] ノード (ps ノード)
 ニューラルネットワークのパラメータ及び K-FAC の
 保持及び共有のみを行う.
- 勾配計算, クロネッカー因子更新ノード (grad ノード)
- ps ノードからパラメータ W を取得し, grad 計算を 行う.
- うち1ノード (チーフ) は cov 計算を行い,各 grad ノードが計算した勾配 ∇L を集計し平均をとったも のを ∇L_{avg} とし, $\check{F}^{-1}\nabla L_{avg}$ を用いて ps ノードのパ ラメータを更新する.
- 逆行列計算ノード (inv ノード)
 ps ノードから Ā_{i,i}, G_{i,i} を取得し inv 計算を行う.計算
 した Ā⁻¹_{i,i}, G⁻¹_{i,i} で ps ノードのそれらを更新する.grad
 ノードとは独立しており,常時 inv 計算をし続ける.

実装には TensorFlow^{*2} 1.4.0 を用い, ノード間の通信は Google が開発している RPC(Remote Procedure Call) で ある gRPC^{*3}を用いた.

各実験におけるノード数と各 grad ノードが計算するミ ニバッチサイズなどの分散計算のための設定は表 2 に示す. 平均損失 *L* を

 $\overline{L} = \operatorname{avg}_{3-\operatorname{Experiments}} \left[\operatorname{avg}_{\operatorname{Mini-Batch}} \left[L(y, z) \right] \right]$ (18)

のようにミニバッチでの損失関数値の平均をとり,さらに 3回実験を行った際の平均と定義する.

^{*1} https://www.cs.toronto.edu/ kriz/cifar.html

^{*2} https://www.tensorflow.org

^{*3} https://grpc.io





反復数



図 4 ミニバッチサイズ変更実験 (SGD との収束の比較)

実験1:ミニバッチサイズ変更

大きいサイズのミニバッチでの損失関数値の減少傾向を 調査するため、ミニバッチサイズを128から16384まで2 倍ずつ変化させた際の平均損失の減少傾向を比較した(図 3). 横軸は反復数,縦軸は平均損失を表している. ラベル は kfac-b[ミニバッチサイズ]である.

ミニバッチサイズが 1024 以上のときは平均損失の減少傾向に差は見られなかった.ミニバッチサイズが 128,256,512のときはミニバッチサイズが小さいほど平均損失の減少が 遅くなった.以上より K-FAC を用いた場合でもミニバッ チサイズを大きくするとある程度までは平均損失の減少速 度は速くなるが,それ以上収束の傾向が変わらなくなるミ ニバッチサイズが存在すると考えられる.

また,同じネットワークをバッチサイズを 512,2048, 8192 として SGD と K-FAC で学習した場合の平均損失の 減少傾向を比較した (図 4). 横軸は反復数,縦軸は平均損 失を表している. ラベルは [kfac/sgd]-b[ミニバッチサイズ] である. SGD ではバッチサイズが 2048 で最も平均損失 の減少が速く,512 と 8192 では遅くなった. これに対し K-FAC では,バッチサイズを大きくしても平均損失の減 少傾向が劣化することはなかった. これにより SGD に対



図 5 inv ノード数変更実験





する K-FAC の優位性を確認できた.

実験 2:inv ノード数変更

分散計算で inv ノード数 1, 2, 3, 4 と変化させた際の平 均損失の減少傾向を比較した (図 5). 横軸は反復数, 縦軸 は平均損失を表している. ラベルは kfac-inv[inv ノード数] である.

inv ノード数を変化させても平均損失の減少傾向に差が ないことが分かる. これはこのニューラルネットワークで は inv ノードが 1 台で十分な頻度 (1 反復に 1 回以上) の inv 計算を行えているからであると考えられる.

より多くのパラメータ W をもつニューラルネットワー クの学習では、1回の grad 計算より inv 計算にかかる時間 の方が長くなり、十分な頻度の inv 計算を行うためには inv ノードを複数台用意する必要が出てくると考えられる.

実験 3:ミニバッチ分割数変更



図7 ミニバッチ分割数変更実験(計算時間比較)

大きさ 16384 のミニバッチを 4, 8, 16, 32, 64 分割し た際の平均損失の減少傾向を比較した (図 6, 7). 横軸は図 6 では反復数, 図 7 では計算時間であり, 縦軸は平均損失 である. ラベルは kfac-div[分割数] である.

反復数で比較した場合,分割数を変化させても平均損失 の減少傾向は変わらなかった.計算時間で比較した場合, 分割数を大きくすることで平均損失の減少速度は速くなっ た.一般に分割数を大きくすると1つのgradノードの計 算量は減少し,全体の通信量は増加する.分散学習をする 際はこれを考慮して分割数を決めるべきであり,今回用い たニューラルネットワークと学習データでは更に大きい分 割数でも高速化できる可能性がある.

4. おわりに

本研究ではパラメータ W の数が比較的少ないニューラ ルネットワークを用いて K-FAC を用いた分散学習の有効 性と分散計算のための設定を確かめることができた.

一方でさらなる K-FAC を用いた分散学習の調査として 以下のような課題がある.

より多くのパラメータ W をもつニューラルネットワーク での実験

AlexNet[8], VGG[10], GoogLeNet[9], ResNet[11] と言っ た今回実験に用いたニューラルネットワークより多くのパ ラメータ W をもつニューラルネットワークを K-FAC を 用いて学習した場合の精度と収束性を調査する.本研究よ り十分な頻度で inv 計算を行うために必要な inv ノードの 数はニューラルネットワークのパラメータ W の要素数に よって決めるものと考えられるため,これらのニューラル ネットワークでどれほどの inv ノードが必要になるかを検 証する.また,大きなニューラルネットワークでの分散学 習ではパラメータを共有する際の通信コストが大きくなる ので,通信時間と計算時間を総合的に考えた場合の最適な 分散計算のための設定を調査する必要がある.

より大きな学習セットでの実験

最適なミニバッチサイズは学習データのクラス分類数に も依存する.本研究では10クラス分類の CIFAR-10を用 いたが,ILSVRC12*⁴で用いられた1000クラス分類用の データセットなどより多くのクラス分類を持つ学習データ を用いた場合の最適なミニバッチサイズを調査し,それに 対する最適な分散設定を調査する必要がある.

計算機に合わせた実装の改良

本研究では全ノードが gRPC を用いて通信を行った.しかし1ノード内に複数 GPU を持つ計算機を用いた場合では NVIDIA NVLink^{*5}などの高速な通信手段を備えている場合があり,これを有効に活用することでさらなる高速化が期待できる.

K-FAC の分散学習を考えると,複数 GPU を持つノード では例えば grad ノードと inv ノードの掛け持ちをさせる ことができ,掛け持ちのさせ方によって通信時間に変化が 見られると考えられる.

謝辞 本研究は JST, CREST の支援を受けたもので ある.

参考文献

- Barret Zoph , Vijay Vasudevan, Jonathan Shlens, Quoc V. Le: Learning Transferable Architectures for Scalable Image Recognition, arXiv 1707.07012
- David E. Rumelhart, Geoffrey E. Hinton, Ronald J. Williams: Learning representations by back-propagating errors, Nature 323
- [3] Amari Shun-ichi : Natural Gradient Works Efficiently in Learning, Neural Computation 1998
- [4] James Martens ,Roger Grosse : Optimizing Neural Networks with Kronecker-factored Approximate Curvature, ICML 2015
- Jimmy Ba, Roger Grosse, James Martens : Distributed Second-Order Optimization using Kronecker-Factored Approximations, ICLR 2017
- [6] Roger Grosse, Ruslan Salakhudinov : Scaling up Natural Gradient by Sparsely Factorizing the Inverse Fisher Matrix, PMLR 2015
- [7] Nicolas Le Roux, Pierre-Antoine Manzagol, Yoshua Bengio : Topmoumoute online natural gradient algorithm, NIPS 2007
- [8] Alex Krizhevsky and Ilya Sutskever and Geoffrey E. Hinton: ImageNet Classification with Deep Convolutional Neural Networks, NIPS 2012
- [9] Christian Szegedy ,Wei Liu ,Yangqing Jia ,Pierre Sermanet ,Scott Reed ,Dragomir Anguelov ,Dumitru Erhan ,Vincent Vanhoucke ,Andrew Rabinovich : Going Deeper with Convolutions, CVPR 2015
- [10] Simonyan, K. and Zisserman, A.: Very Deep Convolutional Networks for Large-Scale Image Recognition. CoRR 2014
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, Jian Sun:

^{*4} http://www.image-net.org/challenges/LSVRC

^{*5} http://www.nvidia.co.jp/object/nvlink-jp.html

[12] Mu Li, David G Andersen, Jun Woo Park, Alexander J Smola, Amr Ahmed, Vanja Josifovski, James Long, Eugene J Shekita, Bor-Yiing Su : Scaling Distributed Machine Learning with the Parameter Server, OSDI'14

分散計算のための設定

今回実験に使用した計算機は表1の2種類である.

付 録

A.1 クロネッカー積

$$A \in \mathbb{R}^{m \times n}, B \in \mathbb{R}^{m' \times n'}$$
に対するクロネッカー積は

$$A \otimes B = \begin{bmatrix} A_{1,1}B & A_{1,2}B & \cdots & A_{1,n}B \\ A_{2,1}B & A_{2,2}B & \cdots & A_{2,n}B \\ \vdots & \vdots & \ddots & \vdots \\ A_{m,1}B & A_{m,2}B & \cdots & A_{m,n}B \end{bmatrix}$$
(A.1)

と定義され,その大きさは mm'×nn'である.

計算機	略号	項目	型番					
TSUBAME 3.0^{*4}	Т	CPU	Intel Xeon E5-2680 $\mathrm{v4^{*5}}$					
		GPU	NVIDIA Tesla P100 *6					
reki	R	CPU	Intel Xeon E3-1220 v3 \ast7					
		GPU	NVIDIA GeForce GTX 1080 \ast8					
表 1 使用計算機								

実験は次に示す設定で行った.

実験	ラベル	#p	#g	#i	#b	計算機
	kfac-b128	1	1	1	128	R
	kfac-b256	1	1	1	256	R
	kfac-b512	1	1	1	512	R
	kfac-b1024	1	2	1	512	R
1	kfac-b2048	1	4	1	512	Т
	kfac-b4096	2	8	1	512	Т
	kfac-b8192	2	16	1	512	Т
	kfac-b16384	4	32	1	512	Т
	sgd-b512	1	1	-	512	R
	sgd-b2048	1	4	-	512	R
	sgd-b8192	2	8	-	1024	R
2	kfac-inv1	1	3	1	2048	Т
	kfac-inv2	1	3	2	2048	Т
	kfac-inv3	1	3	3	2048	Т
	kfac-inv4	1	3	4	2048	Т
3	kfac-div4	1	4	1	4096	Т
	kfac-div8	2	8	1	2048	Т
	kfac-div16	2	16	1	1024	Т
	kfac-div32	4	32	1	512	Т
	kfac-div64	7	64	1	256	Т

表 2 #p, #g, #i, #b はそれぞれ ps ノード数, grad ノード数, inv ノード数, 1 ノードが計算するミニバッチサイズを表す.

*4 http://www.gsic.titech.ac.jp/tsubame3

- *5 https://ark.intel.com/ja/products/91754/Intel-Xeon-Processor-E5-2680-v4-35M-Cache-2_40-GHz
- *6 http://www.nvidia.co.jp/object/tesla-p100-jp.html
- *7 https://ark.intel.com/ja/products/75052/Intel-Xeon-Processor-E3-1220-v3-8M-Cache-3_10-GHz
- *8 http://www.nvidia.co.jp/graphicscards/geforce/pascal/jp/gtx-1080