

# AIのCO<sub>2</sub>排出量削減のための モデル軽量化技術

Light-weight Model Technologies to Reduce  
CO<sub>2</sub> Emissions from AI



株式会社NTTデータ

〒135-6033 東京都江東区豊洲3-3-3豊洲センタービル  
Tel: 03-5546-8051 Fax: 03-5546-2405  
<https://www.nttdata.com/jp/ja/>



## 目次

CHAPTER.1

はじめに

CHAPTER.2

AIと消費電力

CHAPTER.3

AIの消費電力を削減するアプローチ

CHAPTER.4

AIモデル軽量化による消費電力の削減

CHAPTER.5

実験結果

CHAPTER.6

今後の展望







## CHAPTER.1

### はじめに

近年、AIの技術進化による精度向上が進む一方で、AIサービスの開発・活用の際に必要な電力量の増加に起因するCO<sub>2</sub>排出量が増加しています。

自然言語処理で革新をもたらしたBERTの登場以降、AIモデルの大規模化により、自然言語処理、画像処理、音声処理等の多くの分野で劇的な精度向上が実現し、AIは様々なビジネス変革に応用されています。しかし、大規模AIモデルでは、アルゴリズムの訓練や実際の活用のフェーズで多大な電力が必要となります。

今後、AIに起因する消費電力はますます増加するとみられており、消費電力削減に向けた取り組みが必要不可欠です。

本稿では、AIの消費電力削減に向けた世の中の動向と、ソフトウェアによる消費電力削減のアプローチの詳細について説明します。



## CHAPTER.2

### AIと消費電力

#### 2.1 AIが地球環境に与える影響

Applied Materials社のCEO Gary Dickerson氏は、現在のAIの採用率が進み、ハードウェアやソフトウェアに技術革新がない場合、データセンターにおける全世界の消費電力の割合は、2019年の2%から2025年には10%まで増加する、と推測しています※1。また、低炭素社会戦略センターによると、2018年から2030年にかけてのデータセンターの消費電力は、日本で14 TWhから90 TWh、全世界では190 TWhから3,000 TWhまで拡大する見通しです※2。このように、AIに起因する消費電力は今後急上昇することが予想され、消費電力削減に向けた取り組みが求められます。

#### 2.2 AIによる消費電力の内訳

ディープラーニングを用いてAIモデルを開発し、活用するまでには、大きく「学習」と「推論」の2つのフェーズがあります。

##### 学習フェーズ

AIモデルの開発者が大量の学習データをモデルに読み込ませることで、目的となるAIタスク(機械翻訳や画像分類等)を実行できるようにモデルに教えるプロセスです。

学習フェーズは、汎用的なモデルを開発する「事前学習」と、個別のタスクに最適化する「ファインチューニング」の2つに大別されます。大規模言語モデル開発を含む自然言語処理の学習フェーズの例を説明します。

事前学習では、大規模なデータセットを用いて汎用的な学習済みモデルを開発します。例えばGoogleが提案したBERTは3.4億、OpenAIが提案したGPT-3は1750億ものパラメーターを備えたモデルとなっており、この膨大なパラメーターを学習するには数十GBの大規模なデータが必要となります。

そのため、開発に必要な消費電力も膨大で、GPT-3の学習に要する電力は1,287 MWh、CO<sub>2</sub>排出量552トンに相当すると言われてしています※3。ただし、事前学習モデルの開発は一部の企業や研究機関に限られ、一般のAI開発者はこれらのモデルを活用可能です。

ファインチューニングでは、比較的少ないデータを用いて、開発された事前学習モデルを個別のタスクに最適化します。質問応答、文書要約、文書分類等の具体的なタスクを設定した上で、解きたいタスクのデータを数件～数万件程度用意して学習することで、モデルの精度を向上させます。事前学習と比べて消費電力は少なくなりますが、個別のAI開発プロジェクトごとに行う必要があり、実行される回数は多くなります。

※1) <https://observer.com/2019/08/artificial-intelligence-bitcoin-cloud-computing-climate-change/>

※2) 情報化社会の進展がエネルギー消費に与える影響(Vol.1-IT 機器の消費電力の現状と将来予測-

<https://dl.ndl.go.jp/view/prepareDownload?itemId=info%3Andljp%2Fpid%2F11546567&contentNo=1>

※3) David Patterson et al., Carbon Emissions and Large Neural Network Training



## 推論フェーズ

利用者がトレーニング済のモデルに新しいデータを入力して予測結果の出力を得るプロセスです。

## 学習フェーズと推論フェーズの比較

一般に、学習フェーズと推論フェーズの1回あたりの消費電力は、学習フェーズの方が大きいです。モデルの精度を向上するためにデータ量やパラメーター数を増やすほど、学習フェーズの消費電力は大きくなっていきます。

一方で、学習フェーズはモデルの開発者が1回または複数回行えば開発が完了するのに対し、推論フェーズは予測のたびに実行されるため、用途によっては数万回から数百兆回以上も実行されることになります。

では、AIのライフサイクル全体では学習フェーズと推論フェーズではどちらの消費電力が大きいのでしょうか。

- Amazon Web Servicesを利用する顧客からは、機械学習サービスにかかるコストのうち90%は推論が占めるとの声があがっています。<sup>※4</sup>
- NVIDIA社のCEO Jensen Huang氏は、2019年に機械学習コストの80-90%は推論フェーズによるものであると推定しています。<sup>※5</sup>
- 自社で大規模なモデル開発を行うGoogleでも、2022年の報告で、自社で機械学習のために消費する電力のうち60%が推論、40%が学習に起因すると述べています。<sup>※6</sup>
- Facebookは、自社のサーバー経由のニューラルネットワークを使い、全世界のスマートフォン上で1日あたり200兆回の推論を実行していると報告しています。<sup>※7</sup>

以上のように、推論では1回あたりの消費電力は学習より小さいものの、実行回数が「学習」と比べて桁違いに大きくなるため、トータルの消費電力では、推論が学習よりも大きくなるのがわかります(図1)。

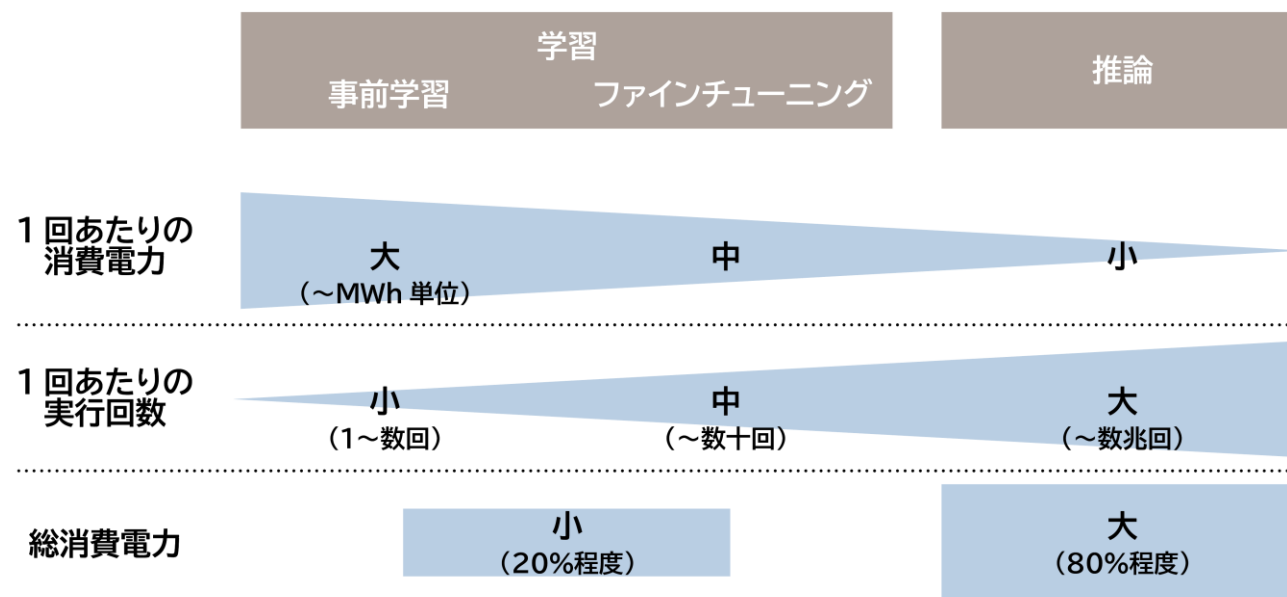


図1: AIフェーズごとの実行回数と消費電力の関係

※4) <https://aws.amazon.com/jp/blogs/aws/amazon-ec2-update-inf1-instances-with-aws-inferentia-chips-for-high-performance-cost-effective-inferencing/>

※5) <https://www.hpcwire.com/2019/03/19/aws-upgrades-its-gpu-backed-ai-inference-platform/>

※6) <https://ai.googleblog.com/2022/02/good-news-about-carbon-footprint-of.html>

※7) <https://engineering.fb.com/2018/05/02/ai-research/announcing-pytorch-1-0-for-both-research-and-production/>





## 2.3 AIの消費電力に対する動向

### 学習における計算量と消費電力の推移

2012年に発表されたAlexNetは、画像からの物体認識のための畳み込みニューラルネットワーク(CNN)です。当時開催された画像分類に関するコンテスト(ILSVRC)にて、次点より10pt以上低いエラー率15.3%で優勝し、ディープラーニングが注目を集めるきっかけとなりました。

当時は画期的な精度を誇ったAlexNetですが、ディープラーニングのアルゴリズムに関する技術革新はすさまじく、AlexNetと同等の精度を実現するために必要な学習の計算量は、2012年から2019年までの7年間で44分の1にまで減少しています。これは、約16ヶ月ごとに必要な計算量が半減していることを意味します。<sup>※8</sup>

一方で、ディープラーニングでは学習に用いるパラメーターが増加するほど精度が高くなることから、学習の計算量も爆発的に増加してきました。最高の精度のモデルを得るための計算量は、2012年から2018年までの6年間で33万倍に増加しました。これは約3.4ヶ月ごとに必要な計算量が倍増していることを意味します。<sup>※9</sup>このように、計算量の増加はアルゴリズムの改善による計算量の低減をはるかに上回るペースで進んでおり、消費電力やCO<sub>2</sub>排出量の増加に繋がっています。

### 推論における計算量と消費電力の推移

推論に要する計算量も、増加傾向にあります。

Canzianiら(2017)<sup>※10</sup>は、ディープラーニングの登場以降、AIモデルの目標が実際の推論時間に関係なく最高の精度を実現することに偏ってきたことを問題視しており、モデルの精度向上の実現が、消費電力や計算リソースの使用率に影響を与えていると指摘しています。

Desislavovら(2021)<sup>※11</sup>は、コンピュータビジョン(CV)や自然言語処理(NLP)の分野での推論に要する計算量や消費電力の推移を報告しています。

報告によると、その時点で最高の精度を得るために必要な推論時の計算量は、

- CVでは、1.42 GFLOPs<sup>※12</sup>(2012年)から5,270 GFLOPs(2021年)へ約3,700倍に増加しました。
- NLPでは、54 GFLOPs(2017年)から740,000 GFLOPs(2020年)へ約13,700倍に増加しました。

計算量の増加に伴い、計算に用いられるGPUのエネルギー効率の改善も進んできましたが、それでも1回の推論に要する消費電力はCVで最大300倍以上、NLPで最大500倍以上増加したことになります。

※8) <https://arxiv.org/abs/2005.04305>

※9) <https://openai.com/blog/ai-and-compute/>

※10) <https://arxiv.org/abs/1605.07678>

※11) <https://arxiv.org/abs/2109.05472>

※12) GFLOPs: モデルの計算コストを表現するためによく使われる指標。浮動小数点の演算回数を表すFLOPsを10億倍したものの。



# AIの消費電力を削減するアプローチ

AIの消費電力を削減するアプローチは以下の3つに大別されます。

### 3.1 クラウド/データセンターにおける削減

学習や推論における計算リソースが実際に消費される設備の消費電力の削減のことです。

例えば、データセンターの電気代の内訳は、50%をサーバーが、30%を電源と冷却系が、20%をその他ストレージ等が占めると言われています。これら機器の省エネを実現することで、消費電力の削減に繋がります。

- 例えば、NTTデータで実証実験中のデータセンターの「液浸冷却システム」では、冷却に必要な消費電力を、従来型のデータセンターと比べて最大97%削減できることを確認しています。<sup>※13</sup>
- また、データセンターの室内環境の可視化システムによりサーバーラームの過冷却を抑制することで、冷却エネルギーを約35%削減することにも成功しています。<sup>※14</sup>
- Googleでは、自社のデータセンターのサーバー冷却に使用する電力に機械学習を用いて最適化したことで、消費電力を40%削減できたと発表しています。<sup>※15</sup>

### 3.2 ハードウェアにおける削減

個別の機器・部品の電力利用最適化による削減のことです。

具体的には、クラウドやデータセンターにおける個々のサーバー、推論を実行する多数のエッジデバイスや、それらの内部のGPUやCPUが対象となります。

AIモデルの特性や推論が実行されるデバイスの制約を踏まえ、最適なハードウェアを選択することで、省エネを実現できる可能性があります。

また、各ハードウェアのメーカーも、ディープラーニングに特化したデバイスの開発に取り組んでおり、同一の計算を実行するための消費電力は年々削減されています。

- 先述したDesislavov らの報告によると、GPUが1Wで処理できる計算量は、2010年に7 GFLOPs/W未満でしたが、2020年時点では100 GFLOPs/Wを超え、CNN や NLP に特化した GPU では300 GFLOPs/Wを超えるものも登場しています。
- ディープラーニングに特化したプロセッサの開発も進んでいます。例えば、プリファード・ネットワークス社は、ディープラーニングの特徴である「行列演算」に最適化した専用チップを開発し、コストを抑えながら実行性能を向上。このチップを搭載したスーパーコンピュータにより、省電力性能ランキングGreen500で2020年、2021年に世界1位を獲得しています。

### 3.3 ソフトウェアにおける削減

AI開発者自身でのアルゴリズムの工夫等による消費電力の削減のことです。

例えば、AIモデル軽量化、エネルギー効率の良いプログラミング言語の選択、データ分析に適したファイル形式の選択が挙げられます。

- AIモデル軽量化: モデルの計算量を小さくすることで推論時間を短縮、メモリ使用量を削減し、推論時の消費電力を削減することが可能です。具体的な手法は次章で説明します。
- プログラミング言語: Python等のインタプリタ型言語は、コードを逐次解釈しながらプログラムを実行するため、処理速度が相対的に遅くなります。CやC++等の事前コンパイル言語を選択することで処理速度が改善し、消費電力削減に繋がる可能性があります。
- ファイル形式: モデル構築では大量のデータをキャプチャして前処理する必要があるため、効率的なデータの処理やストレージによって消費電力を削減することが可能です。例えば、Parquetはデータ保存と検索のために設計された列指向ファイル形式であり、CSVファイルより効率的にデータ処理することが可能です。

※13) <https://www.nttdata.com/jp/ja/news/release/2022/060601/>

※14) <https://www.nttdata.com/jp/ja/news/release/2022/072901/>

※15) <https://www.deepmind.com/blog/deepmind-ai-reduces-google-data-centre-cooling-bill-by-40>





# AIモデル軽量化による消費電力の削減

ソフトウェアによる推論時の消費電力を削減するためのモデル軽量化の代表的な実現方法には、次の3つがあります。

- (1) Pruning(枝刈り)
- (2) Distillation(蒸留)
- (3) Quantization(量子化)

どのモデル軽量化の方法も、電力削減とモデルの予測精度はトレードオフの関係にあります。モデルを軽量化すると推論時の計算量が少なくなり消費電力は減りますが、モデルが覚えている単語の数や文法規則の理解度を落とすことに相当するので、精度が悪くなってしまいます。

NTTデータでは、できるだけ精度を落とさずに電力削減効果を最大化するため、これらの手法を組み合わせることでアルゴリズムやパラメーターを最適化する手法を開発しています。

## 4.1 Pruning (枝刈り)

一般に、Pruningは、ディープラーニングのモデルから重要度の低いパラメーター(通常、0に最も近いノードまたは重み)を削除することを意味します。

モデル軽量化のアプローチとして、重み行列の行と列を削除するMatrix Pruning(重み行列の枝刈り)と、レイヤーそのものをモデルから削除するLayer Removal(レイヤー削除)の2通りを説明します(図 2)。

- Matrix Pruning  
推論時の消費電力の削減を考える場合、重み行列からいくつかの行と列を削除することでモデル全体の計算量を小さくする方法が効果的です。
- Layer Removal  
多くのレイヤーから構成される大規模な言語モデルの場合、モデルからいくつかのレイヤーを削除するアプローチが可能です。レイヤー削除は、適度な量であれば最終的な精度への影響は少ないことが示されています。\*16

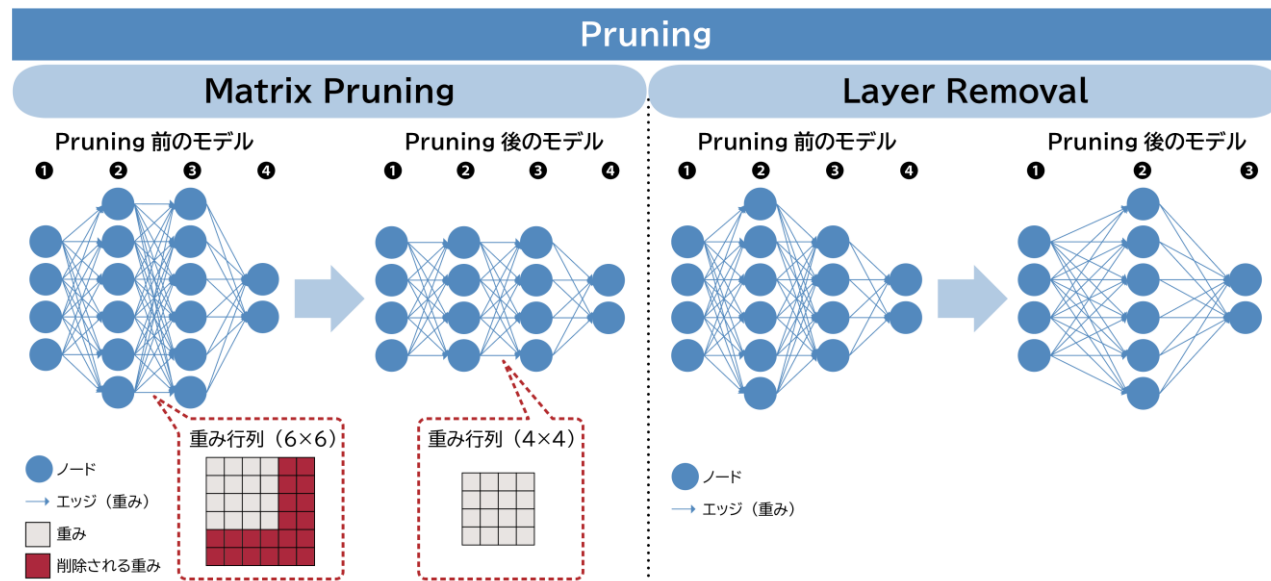


図 2: Pruning手法のイメージ

\*16) <https://arxiv.org/pdf/2004.03844.pdf>

## 4.2 Distillation (蒸留)

Distillationは、モデル圧縮の最も一般的な手法の1つです。学習済の大きい「教師」モデルの出力結果を小さい「生徒」モデルに学習させて、できるだけ「教師」モデルと同じ出力を出すように訓練することを指します。この手法を用いると生徒モデルの精度向上に繋がりますが、十分な精度改善効果を得るためには、大量のデータでの学習が必要です。

代表的な適用例には、小規模な言語モデルとして広く使用されているDistillBERTや、エッジデバイス・モバイル機器で使用できるほど小さいTinyBERT・MobileBERTがあります。

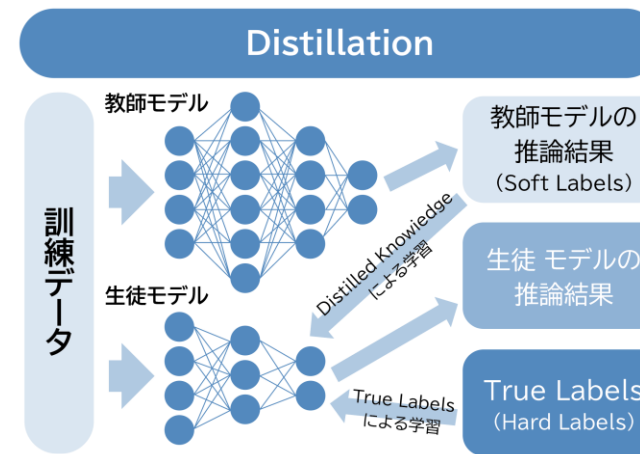
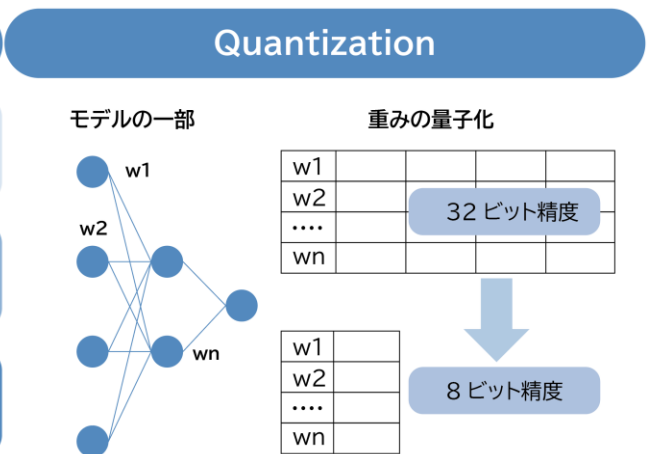


図 3: DistillationとQuantizationのイメージ

## 4.3 Quantization(量子化)

Quantization(量子化)は、前の2つの手法と異なり、モデルの構造を変更するのではなく、重みのビット数を減らすことで計算を近似する手法です。モデルの精度はほとんど変わりませんが、モデルに必要なメモリは大幅に削減されます。

しかし、これらの計算はほとんどのGPUでサポートされていないため、主にCPUを用いるエッジデバイスやモバイル機器上で推論を実施する場合に適用することが可能です。



## 4.4 モデル圧縮技術の組み合わせ

それぞれの手法には利点・欠点がありますが、適切に組み合わせることで、様々な要件に対応可能で持続可能・グリーンなAIモデルを実現できる可能性があります。一方で、近年のAIに関する技術進化のスピードは早く、BERT以降、事前学習された大規模AIモデルが続々と登場しているため、その消費電力削減を実現する技術も特定のモデルに特化するのではなく、様々なAIモデルにタイムリーに対応できる汎用性が求められます。

NTTデータでは現在、精度・消費電力(CO<sub>2</sub>排出量)・GPUでの利用可否といった観点から最適な手法の組み合わせを検証することで、特定のAIモデルに限定せず、任意のAIモデルで推論時の消費電力削減を可能とする汎用的な軽量化手法の確立を目指しています。

次の章で、NTTデータが検証中の手法について、詳細を見ていきます。





## CHAPTER.5

# モデル軽量化技術の検証

### 5.1 検証条件

前章で紹介したモデル軽量化の3つの手法は一般的なもののですが、これらの組み合わせによる最適化はAI開発者個人のノウハウに委ねられています。

NTTデータでは、これまでに自然言語処理の分野で培ってきた知見を活用し、最適な軽量化手法の組み合わせを実現することを目指しています。このNTTデータが検証している手法では、任意の言語モデルに対し、1回限りの調整で軽量化モデルを構築できることが利点になります。

この手法の検証にあたり、モデル圧縮の度合いを調整するため、様々な条件で組合せたLayer Removal、Matrix Pruningと、Distillationを適用しました。

- 具体的には、各言語モデルに対し、軽量化を行わない「ファインチューニング」で得られるモデルを基準とし、軽量化後のモデルサイズが大きい順に、「2/3レイヤー」「1/3レイヤー」「1/3レイヤー+ Matrix Pruning」のように3通りの軽量化を適用しました。
- なお、すべての軽量化モデルに対し、精度向上のためにDistillationを適用しています。

この検証で、Layer RemovalとMatrix Pruningは、Layer Removalを先に適用することが最適であることが分かりました。最初にLayer Removalを適用すると、精度劣化を抑えながらモデル圧縮をすることが可能です。次に、ファインチューニングを実施しながらMatrix Pruningを適用すると、失った情報を再学習するための重みを維持しながら重み行列を削減していくことになるためです。

その他の検証の条件は下記の通りです。

- 言語モデルとして、BERT、GPT-2、T5、OPTを対象としました。
- データセットとして、自然言語処理の精度検証に広く用いられるIMDbやGLUEを対象としました。<sup>※17-18</sup>
- 推論に要する消費電力からCO<sub>2</sub>の相対排出量を換算して評価を行いました。  
軽量化前のモデルでの精度・CO<sub>2</sub>排出量を1とした場合の相対精度(%accuracy)・相対排出量(%emission)を比較しています。
- 推論はすべてNVIDIA A40 GPU上で行いました。

推論にCPUデバイスを用いる場合、Quantizationを適用することも可能です。次節の検証結果はすべてGPUによる検証のため直接比較はできないものの、別途行った検証結果では、量子化モデルは非量子化モデルに対し、99%の相対精度を維持しながら、推論時間とCO<sub>2</sub>排出量をともに半減できることを確認しています。(ただし、CPUはGPUより推論時間や排出量が大きくなるため、量子化を適用すべきケースは、エッジデバイスやモバイル端末などGPUを適用できないケースに限られます。)

次節で、NTTデータが検証中の手法で実際に推論をした精度と、その推論に要するCO<sub>2</sub>排出量の結果を示します。



※17) IMDb: 自然言語処理の検証に広く用いられるデータセットであり、映画のレビューコメントがポジティブかネガティブかの2値分類に用いられます。データ数は学習用、テスト用ともに25,000です。

※18) GLUE: 自然言語処理用のデータセットであり、IMDbと比べてやや難易度の高い9つのタスクから構成されています。データ数はタスクによって異なります。





## 5.2 検証結果

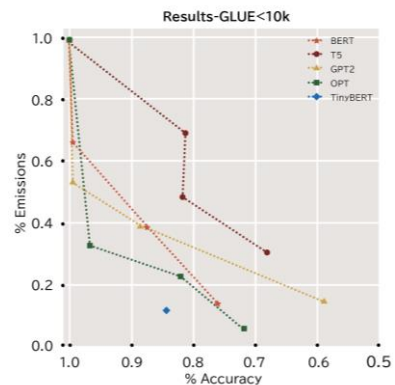


図4: GLUE(データ数10,000未満)の検証結果

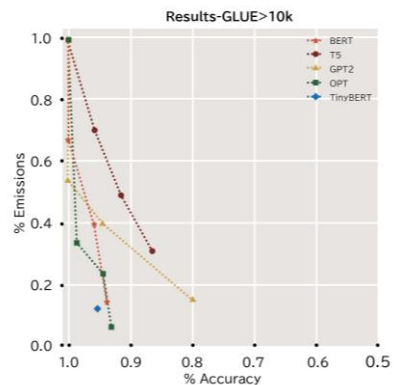


図5: GLUE(データ数10,000以上)の検証結果

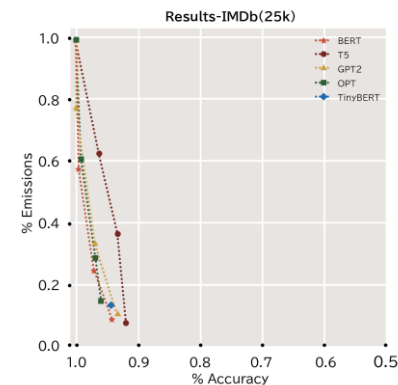


図6: IMDb(データ数25,000)の検証結果

BERT、T5、GPT-2、OPTの4種類のモデルについて検証した結果を示します。

なお、タスクによってデータ数が異なるGLUEでは、データ数によって異なる傾向の結果が得られたので、データ数 10,000未満 (<10K) と 10,000以上 (>10K)のタスクに分けて評価を行いました。

図4はGLUE (<10K)、図5はGLUE (>10K)、図6はIMDbの結果を表します。

また、BERTの検証結果との比較対象として、一般的な軽量化モデルであるTinyBERTの結果も示しています。全体の傾向として、すべての言語モデルについて、相対精度を可能な限り維持しながら、CO<sub>2</sub>排出量を削減できることを確認しました。

### タスクごとの傾向

- GLUE (<10K) では、モデル軽量化が少ない範囲では相対精度95%を維持できるものの、モデルの軽量化が進むと、Pruningによって失った情報を効果的に再学習することができなくなり、精度劣化が進んでいます。
- GLUE (>10K) では、相対的にPruningを行う余地が大きくなるため、精度劣化が抑制されています。
- データ数が25,000と大きいIMDbでは、すべての言語モデルについて、相対的に優れたパフォーマンスを発揮しました。最も軽量化した場合で、相対精度90%を維持しながら、相対排出量は10%程度まで削減できています。

### 言語モデルごとの傾向

- BERTモデルの結果についてさらに詳細にみみると、データセットが少ないGLUE (<10K)の場合は、Pruningが進むに連れて精度が劣化しています。TinyBERTでも精度劣化しているものの、その影響は小さく留まっています。

- 一方で、データセットが大きいIMDbやGLUE(>10K)の場合は、われわれの手法でBERTのモデル圧縮を進めた場合も、TinyBERTと同様の精度を維持することができています。全体としての傾向は、言語モデルとしての構成がBERTと類似しているOPTでも同様です。
- T5では、4つの言語モデルの中で全体的にやや精度劣化した結果となりました。T5はseq2seqと呼ばれる生成モデルであり、他のモデルと比べると、文の生成を行うために高いレベルでの言語理解を要するモデルとなっています。そのため、特にGLUE (<10K)のようにデータ数が少ない条件では、レイヤー削除が少ない場合でも精度への悪影響が大きくなることがわかります。
- GPT-2では、全体的な傾向はBERTやOPTと同様ですが、グラフの右端にあたるMatrix Pruningを適用した結果のみ傾向が異なりました。GPT-2は、重み行列が相互に結合する珍しい構成をとるため、行と列の順序が非常に重要になります。そのため、Matrix Pruningを適用すると、この学習した順序が崩れてしまい、モデルの精度は著しく劣化します。

以上の結果から、モデルとそのアーキテクチャに応じて、異なるPruningの手法を用いる必要があることがわかります。例えば、

- 使用可能なデータが少ない場合は、T5など生成系のモデルにPruningを適用することは推奨されません。それ以外の言語モデルを用いる場合は、Layer Removalを数層の削減に留めることで一定の精度が維持できます。
- GPT-2のような特徴的なアーキテクチャのモデルでは、Matrix Pruningは行わず、Layer Removalのみを行うことが推奨されます。

このように、モデル軽量化が適切に実行されれば、元のモデルの数分の1の計算量で推論を実行することができ、CO<sub>2</sub>排出量を10分の1程度まで削減できるAIモデルを作成することが可能です。

また、この手法の副次的なメリットとして、表 1に示す通り、モデルサイズや推論時間についてもCO<sub>2</sub>排出量と同様の傾向で削減することができます。表 1ではBERTの結果を示していますが、他の言語モデルについても同様の結果が得られています。

このように、異なるモデル軽量化技術を組み合わせることで、言語モデルの種類によらず、精度・CO<sub>2</sub>排出量・モデルサイズ・推論時間といった様々な条件に合った軽量化モデルが得られる可能性があります。

Model	Model Type	% Accuracy	% Size	% Time
BERT	Finetuned	1	1	1
BERT	8 layers (2/3)	99.7%	74.2%	86.9%
BERT	4 layers (1/3)	91.9%	48.2%	61.4%
BERT	4 layers (1/3) + Matrix pruning	85.6%	21.5%	26.9%

表 1: BERTの検証結果抜粋



## CHAPTER.6

### 今後の展望

以上のように、精度とCO<sub>2</sub>排出量のトレードオフはモデルの種類とデータセットに依存しますが、同一の手法で、任意のモデルに対して、精度の大部分を維持したまま排出量削減ができる可能性を確認できています。この手法を実際のAI開発に適用すれば、ユーザが要求する相対精度やCO<sub>2</sub>排出量といった目標を可能な限り満たすモデル探索の最適化が実現可能です。

また、AIのモデル軽量化手法について、さらなる研究が進められています。将来的には学習の際に時間をかけることなく、精度劣化を抑えてCO<sub>2</sub>排出量を削減できる新しい技術が登場することが期待されます。例えば、

- QuantizationをサポートするGPUの実用化
- 行列計算の高速化に繋がる「疎行列」のディープラーニングへの適用

今後、当社では、将来登場する技術を取り入れながら、CO<sub>2</sub>排出量と精度のバランスをより最適化する軽量化モデルの開発を継続していきます。

さらに、言語モデルに加えて、画像・音声・マルチモーダルモデルへの適用、効果検証を進めていきます。

このホワイトペーパーで説明したように、ソフトウェアの観点からのCO<sub>2</sub>削減アプローチは様々あります。

こうした手法をまとめ標準化する団体として非営利団体のGSF (Green Software Foundation)があり、NTTデータも加盟し、脱炭素に向けた取組を強化していきます。

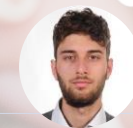
私たちは環境に優しいAIを開発することで、今後もAIによるイノベーションとサステナブルな社会の実現を推進していきます。



#### 著者情報



**野村 雄司**  
技術革新統括本部 システム技術本部  
データ&インテリジェンス技術部 課長



**Paolo Tirota**  
技術革新統括本部 システム技術本部  
データ&インテリジェンス技術部 勤務



**中野 篤**  
技術革新統括本部 システム技術本部  
データ&インテリジェンス技術部 主任

2023年3月