# 3-5 生成AIの基礎と展望

東京大学 数理・情報教育研究センター 2024年6月1日

#### 概要

- 基盤モデルと生成AIについて理解する.
- 実世界で進む生成AIの応用と革新について理解する.
- 大規模言語モデルについて理解する.
- 拡散モデルについて理解する.
- 生成AIの留意事項について理解する.

## 本教材の目次

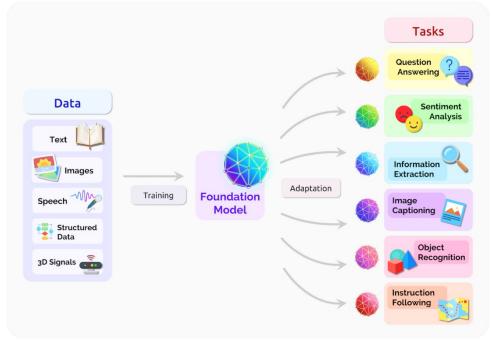
1. 基盤モデル	4
2. 生成AI	6
3. 実世界で進む生成AIの応用と革新	8
4. 大規模言語モデル	18
5. 拡散モデル	24
6. マルチモーダルモデル	25
7. 生成AIの留意事項	26
8. 自己教師あり学習	27
9. Transformer (トラスフォーマ)	29
10. Vision Transformer, CLIP	34
11. 生成モデル	36

#### 基盤モデル

大規模で多種多様なデータをもとに学習され、さまざまなタスクに適 応可能な機械学習のモデルを基盤モデルと呼びます。

"A foundation model is any model that is trained on broad data (generally using self-supervision at scale) that can be adapted (e.g., fine-tuned) to a wide range of downstream tasks." 1)

基盤モデルは主に深層学習<sup>2)</sup>の技術に立脚しています。



文献1)の図 より引用

- 1) R. Bommasani, et al. On the Opportunities and Risks of Foundation Models. arXiv:2108.07258, 2021.
- 2) 応用基礎レベル教材 3-4. 深層学習の基礎と展望 参照

## 基盤モデル

- 基盤モデルは、転移学習という機械学習の方法の考え方がもとになっています。
- 転移学習は、あるタスクから学習した「知識」を他のタスクの問題を解くことに適用する方法です。
- 基盤モデルの学習においては、補助的なタスクで「事前学習」(主に自己教師あり学習<sup>1)</sup>による)を行なったモデルを対象のタスクに適応するように学習(ファインチューニング)する、というように転移学習を行います。
- 基盤モデルはこのような転移学習を大規模に行うことで構築されますが、そのためには(1)GPU(Graphics Processing Unit)のようなハードウェア、(2)トランスフォーマ(Transformer) $^{2)}$ のように並列処理を適用可能なモデル、(3)大規模なデータが必要となります。

- 1) 本教材の項目「自己教師あり学習」参照
- 2) 本教材の項目「トランスフォーマ」参照

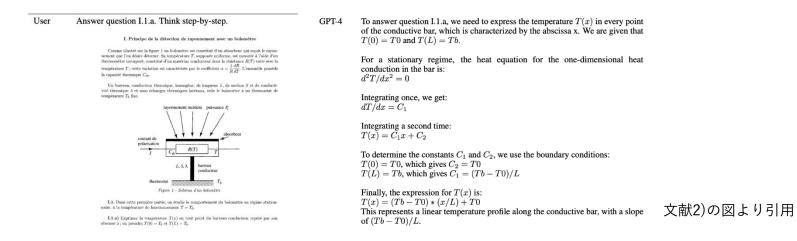
#### 生成AI

- 主に基盤モデルに基づく生成モデル<sup>1)</sup>を用いて、テキスト、画像、音声、動画、その他のコンテンツを生成する人工知能(AI)技術の総称を生成AIと呼びます。
- 生成AIの利用においては、ユーザが自然言語による指示(プロンプト)を入力することで、そのプロンプトに従ってコンテンツを生成 (より一般的にはタスクを実行)することが可能です。
- 例えば、2022年11月に公開されたChatGPT<sup>2)</sup>は、大規模なテキストデータをもとに事前学習された大規模言語モデルを基盤モデルとして、プロンプトに対してテキストなどのコンテンツを生成する生成AIのサービスです。

- 1) 本教材の項目「生成モデル」参照
- 2) <a href="https://chat.openai.com/">https://chat.openai.com/</a>

#### 生成AI

- 2023年3月に公開された大規模言語モデルであるGPT-4では、従来のモデルに比べて言語生成、知識活用、知識推論、記号推論、数学推論などの能力がさらに向上したことが報告されています1)。
- また,生物学,経済学,化学,物理学などの複数の専門的な試験においても一定程度の成績を収めることが報告されています<sup>2)</sup>。



- 数学、プログラミング、医学、法学などでは人間のレベルのパフォーマンスを達成という予備的な実験結果の報告もなされています<sup>3)</sup>。
  - 1) W. Zhao, et al. A Survey of Large Language Models. arXiv:2303.18223, 2023.
  - 2) OpenAI, et al. GPT-4 Technical Report. arXiv:2303.08774, 2023.
  - 3) S. Bubeck, et al. Sparks of Artificial General Intelligence: Early experiments with GPT-4. arXiv:2303.12712, 2023.

- 生成AIは、テキスト、画像、音声、動画などさまざまなモダリティを 入出力としてコンテンツの生成を可能にします.
- 以下のように、入力に対してさまざまなコンテンツを生成可能な生成 モデルの技術とそれらのサービスが現在研究開発されています。
  - テキストの生成 (Text-to-Text, Image-to-Text, Speech-to-Text など)
  - コードの生成 (Text-to-Code など)
  - 画像の生成 (Text-to-Image, Image-to-Image など)
  - 動画の生成 (Text-to-Video など)
  - 音楽の生成 (Text-to-Audio, Text-to-Music)

そのほか、3Dモデルの生成、ロボットやシミュレーションに応用可能な行動や動作系列の生成、などを行う生成AI技術に関する研究開発も進められています。

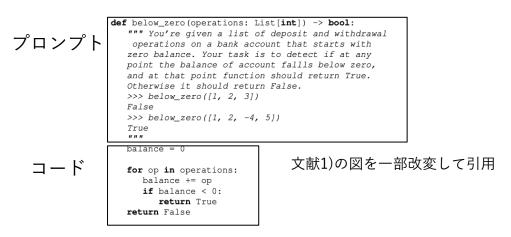
#### テキストの生成

- 大規模なテキストデータをもとに事前学習された大規模言語モデル<sup>1)</sup>により、入力となる指示(プロンプト)に対してテキストを生成.
- テキスト作成の支援、テキストの翻訳・要約・執筆・校正、対話応答な どへ応用することができます。
- 代表的な大規模言語モデルとそのサービス
  - サービス:<u>ChatGPT</u>(OpenAI社)
    - モデルはGPT (Generative Pre-trained Transformer)の系統
  - モデルおよびサービス: Gemini (Google社)
  - モデルおよびサービス: Claude (Anthropic社)
  - サービス:<u>Meta AI</u>(Meta社)
    - モデルはオープンソースの<u>LlaMA</u> (Large Language Model Meta AI)
  - モデル:Mistral, Mixtral (Mistral AI社)
- 日本語大規模言語モデルのまとめ(Ilm-jp, 国立情報学研究所)

1) 本教材の項目「大規模言語モデル」参照

#### コードの生成

- テキストを生成する大規模言語モデルの技術は、プログラムのソースコードを生成するText-to-Codeモデルにも応用されています。
- Text-to-Codeモデルはプログラミングやソフトウェア開発の支援へ応用することができます.
  - Text-to-CodeモデルのCodex 1)でのコード生成の例

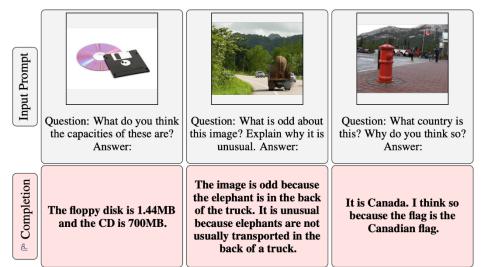


- Text-to-Codeモデルのサービスの例
  - <u>GitHub Copilot</u> (GitHub社)

1) M. Chen, et al. Evaluating Large Language Models Trained on Code. arXiv:2107.03374, 2021.

## 画像からテキストの生成

- 大規模なテキストと画像のデータをもとに事前学習されたマルチモー ダルな基盤モデルである画像-言語(Vision-Language)モデルにより、 テキストや画像を入力としてテキストを生成。
- キャプションの生成,画像を用いた質問応答や対話などへ応用することができます.
  - 画像言語モデルFlamingo<sup>1)</sup>での画像質問応答の例



文献1)の図より引用

1) J.-B. Alayrac, et al. Flamingo: a visual language model for few-shot learning. arXiv:2204.14198, 2022.

## 音声からテキストの生成

- 音声入力に対してテキストを生成する音声認識技術の開発も行われています.
  - 従来の音声認識技術に対して、基盤モデルに基づく技術が近年導入されています。
  - 例えば、ChatGPTのサービスを提供しているOpenAI社は、音声認識と自動翻訳のためのモデルである $\underline{\text{Whisper}}^{1)}$ を公開しています。
- 音声からテキストの生成は、対話システムや自動応答システム、文字起 こしや翻訳の自動化などへ応用することができます。
- テキストから音声を生成するText-to-Speech (TTS) の技術も発展しており、人のような自然な発話を生成することも可能になっています.
- 音声認識やTTSにより、生成AIを対話的に利用可能にするサービスの開発も行われています。
  - 例えば2024年5月に発表されたGPT-4oでは、人と会話するのと同様にほぼ遅延なく音声で生成AIと対話することも可能になっています<sup>2)</sup>.
    - 1) A. Radford, el al. Robust Speech Recognition via Large-Scale Weak Supervision. arXiv:2204.14198, 2022.
    - 2) <a href="https://openai.com/index/hello-gpt-4o/">https://openai.com/index/hello-gpt-4o/</a> (2024-06-01 閲覧)

#### 画像の生成

- 大規模なテキストと画像のデータをもとに事前学習されたマルチモーダルな基盤モデルであるText-to-Imageモデルにより、入力(プロンプト)に対して画像を生成。
  - 生成モデルや拡散モデルなどの技術が活用されています¹)。
- 画像作成や編集の支援に応用することができます.
  - Text-to-ImageモデルのDALL-E2<sup>2)</sup>での画像生成の例



vibrant portrait painting of Salvador Dalí with a robotic half face



a shiba inu wearing a beret and black turtleneck



a close up of a handpalm with leaves growing from it

文献2)の図より引用

13

- 1) 本教材の項目「拡散モデル」,項目「生成モデル」参照
- 2) A. Ramesh, et al. Hierarchical Text-Conditional Image Generation with CLIP Latents. arXiv:2204.06125, 2022.

## 画像の生成

- 代表的な画像生成のサービス
  - DALL-E (ChatGPTで利用可能) (OpenAI社)





文献1)の図より引用

Stable Diffusion (Stability AI社)





文献2)の図より引用

- <u>Midjourney</u> (Midjourney社)
  - 1) J. Betker, et al. Improving Image Generation with Better Captions. <a href="https://cdn.openai.com/papers/dall-e-3.pdf">https://cdn.openai.com/papers/dall-e-3.pdf</a>, 2023.
  - 2) D. Podell, et al. SDXL: Improving Latent Diffusion Models for High-Resolution Image Synthesis. arXiv:2307.01952, 2023.

#### 動画の生成

- Text-to-Imageモデルを発展させ、入力(プロンプト)に対して動画を生成するText-to-Videoモデルの開発も行われています。
- Text-to-Videoモデルは現実の物理世界をシミュレートする世界シミュレーターへの応用も期待できます。
  - Text-to-VideoモデルのSORA<sup>1)</sup>での1分程度の動画生成の例<sup>2)</sup>

#### **Text Prompt**

A stylish woman walks down a Tokyo street filled with warm glowing neon andanimated city signage. She wears a black leather jacket, a long red dress, and black boots, and carries a black purse. She wears sunglasses and red lipstick. She walks confidently and casually. The street is damp and reflective, creating a mirror effect of the colorful lights. Many pedestrians walk about.

#### Generated video





文献2)の図より引用

- 1) Video generation models as world simulators, https://openai.com/index/video-generation-models-as-world-simulators (2024-05-01 閲覧)
- 2) Y. Liu, et al. Sora: A Review on Background, Technology, Limitations, and Opportunities of Large Vision Models. arXiv:2402.17177v3, 2024.

#### 音楽の生成

- 入力(プロンプト)に対してオーディオ(音声やサウンドエフェクトなど)を生成する Text-to-Audioモデルや音楽を生成する Text-to-Musicモデルの開発も行われています。
  - オーディオ生成の例
    - Audiobox<sup>1)</sup>
    - AudioLDM<sup>2)</sup>
  - 音楽生成の例
    - MusicGen<sup>3)</sup>
    - MusicLM<sup>4)</sup>

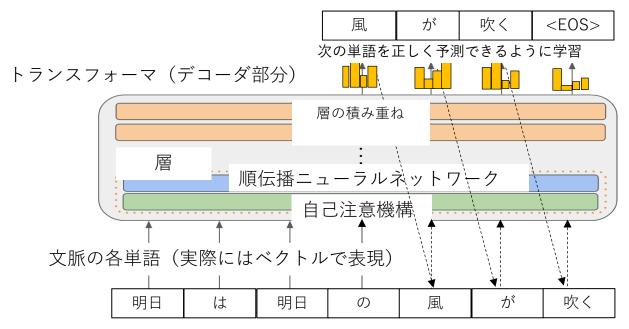
- 1) A. Vyas, et al. Audiobox: Unified Audio Generation with Natural Language Prompts. arXiv:2312.15821v1, 2023.
- 2) H. Liu, et al. AudioLDM: Text-to-Audio Generation with Latent Diffusion Models. arXiv:2301.12503v3, 2023.
- 3) J. Copet, et al. Simple and Controllable Music Generation. arXiv:2306.05284v3, 2024.
- 4) A. Agostinelli, et al. MusicLM: Generating Music From Text. arXiv2301.11325, 2023.

#### 産業応用

生成AI技術のさまざまな産業応用が現在検討・試行されています.

- ビジネス
  - 文書作成支援・要約・翻訳、マーケティング支援、カスタマーサービス、ブレインストーミングなど
- 教育
  - 教材作成、個別指導、自動評価、外国語学習など
- 医療
  - 診断・治療支援, 創薬 など
- メディア・広告
  - コンテンツ・シナリオ・デザイン作成や編集 など
- 芸術
  - 絵画や音楽作成 など
- IT
  - プログラミング・ソフトウェア開発・データ分析支援 など

- 大規模なテキストデータをもとに事前学習された基盤モデルは、大規模言語モデル(large language model, LLM)と呼ばれています。
  - LLMは主にトランスフォーマ $^{1)}$ と呼ばれるニューラルネットワークのアーキテクチャに基づいています。
  - 特にテキストの生成に用いられるLLMの事前学習では、与えられた文脈 (単語(実際にはトークン)の系列)に続く次の単語(実際には確率分 布)を予測するという自己教師あり学習<sup>2)</sup>を行います。



- 1) 本教材の項目「トランスフォーマ」参照
- 2) 本教材の項目「自己教師あり学習」参照

## 事前学習

- LLMの事前学習は、与えられた単語(トークン)の系列に続く次の単語の出現確率をモデルが正しく予測できるようになるように、トランスフォーマを構成するニューラルネットワークの重み(パラメータ)を学習することで行われます。
- この学習には、大規模なテキストデータが用いられます。
  - 例えば、ウェブ、本、コード、学術文献などのテキストデータを含みます.
  - 2020年に開発されたLLMであるGPT-3<sup>1)</sup>では、およそ3000億トークン(数百万冊の本に相当)のテキストデータをもとに事前学習が行われ、モデルのパラメータの総数は1750億に及びます。
- LLMの学習には膨大なコストが必要となります.
  - 例えば、オープンソースのLLMであるLLaMA(2023年に公開)は、650 億パラメータのモデルの学習のため、2048個のGPUで21日間の期間を要した(数百万ドルのコストに相当)、と報告されています $^{2)}$ .
    - 1) T. B. Brown, et al. Language Models are Few-Shot Learners. arXiv:2005.14165, 2020.
    - 2) H. Touvron, et al. LLaMA: Open and Efficient Foundation Language Models. arXiv:2302.13971, 2023.

#### 指示チューニングとアライメント

- LLMの学習では、事前学習に続いて指示チューニング $^{1)}$ およびアライメント $^{2)}$ によりモデルのファインチューニング $^{3)}$ を行います。これにより、LLMは入力された指示に対して適切な応答を返すように学習されます。
  - 指示チューニング (Instruction tuning)
    - 推論,要約,質問応答などのさまざまな問題(タスク)について,入力とそれに対する適切な応答のペアをもとにモデルをファインチューニングします.
    - 多様なタスクで指示チューニングすることで、モデルが未知のタスク にも適応できるようにその能力が引き出されると考えられています。
  - アライメント (Alignment)
    - 入力に対して倫理的・社会的に適切な応答を返すように実際の人のフィードバックをもとにモデルをファインチューニングします。
    - 人のフィードバックによる強化学習(reinforcement learning from human feedback, RLHF)や直接選好最適化(direct preference optimization, DPO)などの手法が用いられます。
      - 1) J. Wei, et al. Finetuned Language Models Are Zero-Shot Learners. arXiv:2109.01652, 2021.
      - 2) L. Ouyang, et al. Training language models to follow instructions with human feedback. arXiv:2203.02155, 2022.
      - 3) 本教材の項目「基盤モデル」参照

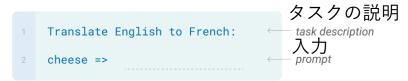
## 文脈内学習

- 事前学習と指示チューニングにより学習されたLLMは、対象のタスクに関する入力と応答の例をいくつか与えるだけで、そのタスクに適応させることができることがわかってきています。
- このように例の教示によって未知のタスクに適応させる学習を文脈内 学習(In-Context Learning)と呼びます.

ゼロショット (タスクの説明のみを与える)

#### Zero-shot

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

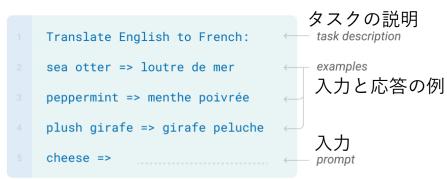


文献1)の図より一部改変して引用

フューショット(タスクの説明に加えて, 入力と応答の例を与える)

#### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.



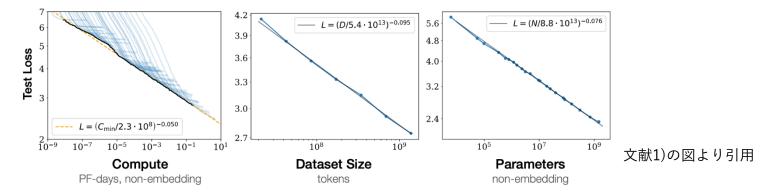
1) T. B. Brown, et al. Language Models are Few-Shot Learners. arXiv:2005.14165, 2020.

## プロンプトエンジニアリング

- 文脈内学習において、入力となる指示(プロンプト)の与え方を工夫することで、LLMの能力を引き出すことができることがわかってきています。
- タスクに応じてLLMが適切な応答を返すように、さまざまなプロンプトの作り方(プロンプトエンジニアリング)が提案されています。
  - 例えば、「ロールプレイング」(Role-Playing)という方法では、LLMに特定の役割や状況をあらかじめ想定させ、そのもとでタスクを解かせます。
  - 「思考の連鎖」(Chain-of-Thought, CoT)という方法では、文脈内学習のフューショットにおいて入力から応答の例を教示する際に、その中間の過程を含めて教示することで、推論を要する複雑なタスクを解かせます。また、ゼロショットにおいても「ステップ・バイ・ステップで考えてみましょう」と指示することでLLMの適切な推論の能力を引き出すことがわかっています。
  - プロンプトエンジアリングの実践例
    - Prompt Engineering Guide (DAIR.AI)
    - <u>文章生成AI利活用ガイドライン</u>(東京都デジタルサービス局)
    - <u>中小企業のための「生成AI」活用入門ガイド</u>(東京商工会議所)

#### スケーリング則

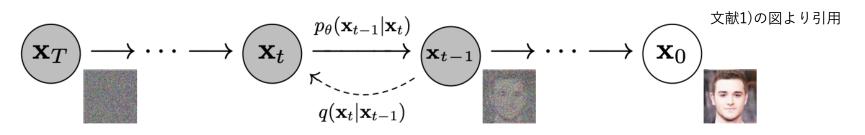
LLMのサイズ(パラメータの数)やデータセットのサイズ(トークンの数)や学習に用いる計算資源の量が大きくなると、スケーリング則(べき乗則)の関係に沿ってモデルの性能も向上(テストデータに対する損失が減少)することが報告されています¹).



- また、LLMのモデルのサイズが大きくなるとさまざまなタスクにおいて文脈内学習の性能が向上するというLLMの創発(Emerging Abilities)が起こるという議論もなされています<sup>2)</sup>。
- 現在,巨大テック企業を中心にLLMの開発競争が起こっています.
  - 1) J. Kaplan, et al. Scaling Laws for Neural Language Models. arXiv:2001.08361, 2020.
  - 2) J. Wei, et al. Emergent Abilities of Large Language Models. arXiv:2206.07682, 2022.

#### 拡散モデル

- 拡散モデルは、画像の生成に活用できる生成モデルです。
  - 拡散モデル<sup>1)</sup>では、データにノイズを付加することを繰り返す(下図の右から左の方向)拡散過程により、データをノイズへ変換します。
  - 次に,この変換を逆の方向(下図の左から右の方向)へたどり,ノイズを除去することを繰り返すことで,ノイズからデータを生成する逆拡散過程を学習します.
  - これにより、拡散モデルは高品質で多様な画像生成を可能にします.
  - またテキストで条件付けを行った画像を生成することも可能にします.



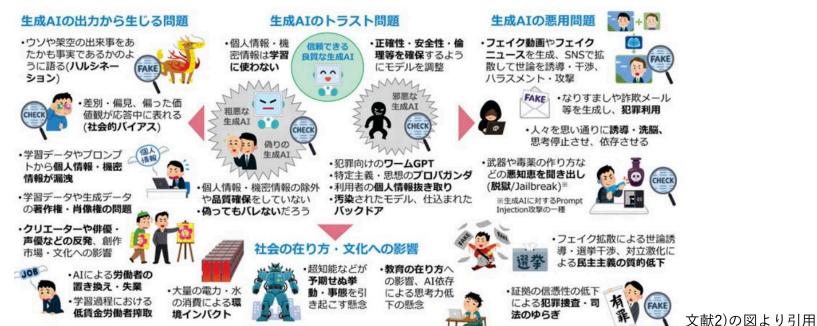
- 拡散モデルを確率微分方程式で一般化した方法も提案されています2).
  - 1) J. Ho, et al. Denoising Diffusion Probabilistic Models. arXiv:2006.11239, 2020.
  - 2) Y. Song, et al. Score-Based Generative Modeling through Stochastic Differential Equations, arXiv:2011.13456, 2020.

#### マルチモーダルモデル

- 画像-言語(Vision-Language)モデルのように異なるモダリティを横断的に扱うことを可能にした基盤モデルの技術 $^{1)}$ により,テキスト,画像,音声,動画などさまざまなモダリティに対応したマルチモーダルモデルとそのサービスが開発されています。
  - マルチモーダルモデルの例
    - <u>LLaVA</u> (Large Language and Vision Assistant) (University of Wisconsin–Madison, Microsoft Research)
    - <u>GPT-4V</u> (OpenAI社)
    - <u>Gemini Pro Vision</u> (Google社)
    - <u>Qwen-VL</u> (Alibaba社)
- 2024年5月時点では、テキスト、画像、音声などを入出力可能なマルチモーダルモデルを元にしたサービスとして<u>GPT-4o</u>(OpenAI社)や <u>Gemini 1.5 Pro Vision</u>(Google社)などが公開されています。
  - 1) 本教材の項目「CLIP」参照

#### 生成AI技術の留意事項

- 生成AIを含む人工知能(AI)技術の急速な発展に伴い、AIのリスクと その対応に関する議論も行われています。
  - 内閣府のAI戦略会議<sup>1)</sup>では、個人情報保護、AIと知的財産権との関係、 偽・誤情報等、雇用への影響、ガイドラインと履行確保等、をAIに関する リスクの論点として整理しています。
  - JST・CRDS(研究開発戦略センター)がまとめた次世代AIモデルに関する戦略プロポーザル $^{2)}$ では、生成AIがもたらすリスクを挙げています。



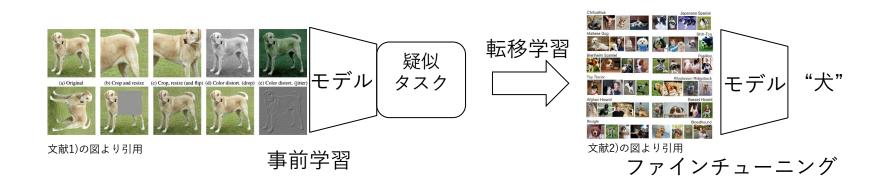
- 1) 内閣府 AI戦略会議, https://www8.cao.go.jp/cstp/ai/ai\_senryaku/1kai/shiryo2.pdf (2024-06-01 閲覧)
- 2) JST・CRDS 戦略プロポーザル 次世代AIモデルの研究開発, CRDS-FY2023-SP-03, 2024.

#### 自己教師あり学習

- 教師あり学習(特に深層学習に基づく学習)では,一般にモデルの学習のために大規模なラベル付き教師データ(人手でラベル付けされたアノテーションデータ)が必要となります.
  - 例えば、画像認識のモデルの学習使用されるデータセットである ImageNetでは1000万以上の画像データについて、各画像にどのような物体が写っているかを示すラベルが2万以上のカテゴリから選ばれ付与されています.
- 自己教師あり学習では、ラベル付き教師データを必要としない補助的なタスク(疑似タスク)によりモデルの事前学習を行います。
  - 例えば画像であれば
    - 一部を欠落させた画像の残りの部分から欠落部分を復元するタスク
    - 様々な変換を施した対の画像についてそれらが元は同じ画像か異なる 画像かを予測するタスク
  - 例えば言語であれば
    - 単語が一部欠落した文脈からその欠落単語を予測するタスク
    - 与えられた文脈に続く単語を予測するタスク

#### 自己教師あり学習

- 転移学習の考え方に従うと、自己教師あり学習において疑似タスクによりモデルの事前学習を行うことは、他のタスクの問題を解くために有用な「知識」をデータから学習していることになります。
- 自己教師あり学習で事前学習したモデルを他のタスクに適用する時、 適用先のタスクを下流タスクと呼ぶことがあります。
- 下流タスクの小規模なラベル付き教師データ(例えば画像であれば画像分類のためのデータセット)があれば、事前学習済みのモデルをそれらの教師データで学習する(これをファインチューニングと呼びます)ことで、モデルを下流タスクに適応させることができます。



<sup>1)</sup> T. Chen, et al. A Simple Framework for Contrastive Learning of Visual Representations. arXiv:2002.05709, 2020.

28

<sup>2)</sup> A. Khosla, et al. Novel dataset for fine-grained image categorization. Workshop on Fine-Grained Visual Categorization, 2011.

#### トランスフォーマ

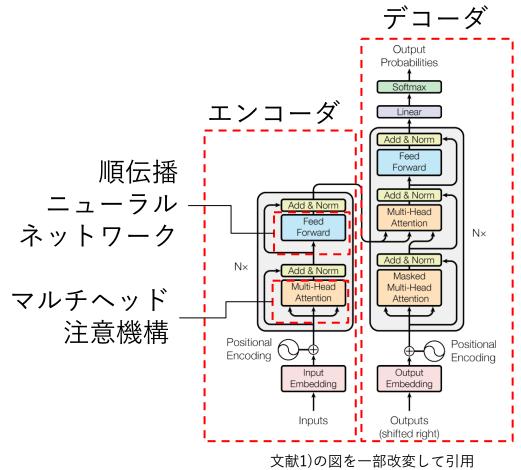
- トランスフォーマ( $\mathsf{Transformer}$ )  $^{1)}$ は, 2017年に提案されたニューラルネットワークのアーキテクチャーです.
- 可変長の系列の情報(例えばテキストの単語の系列)を処理可能なトランスフォーマは、自己注意機構(self attention mechanism)により系列内の離れた情報間の依存関係(例えばテキスト中の離れた単語間の関係)を扱うことに可能にしました。
- トランスフォーマは入力の系列を内部表現に変換するエンコーダと内 部表現を出力の系列に変換するデコーダからなります.
- トランスフォーマのエンコーダやデコーダは基盤モデルのアーキテクチャとしても用いられています。
  - 例えば、後述するBERT (Bidirectional Encoder Representations from Transformers)とGPT (Generative Pre-Trained Transformer)<sup>2)</sup>は、それぞれトランスフォーマのエンコーダとデコーダのアーキテクチャをもとにした言語モデルです。
    - 1) A. Vaswani, et al. Attention Is All You Need. arXiv:1706.03762, 2017.
    - 2) A. Radfold, et al. Improving language understanding by generative pre-training. Technical Report, 2018.

#### トランスフォーマ

- トランスフォーマのエンコーダ・デコーダの主な構造は, (マルチ ヘッド)注意機構と順伝播ニューラルネットワークから構成される層 を多層に積み上げたものになっています(次ページ図参照).
  - 系列内の離れた情報間の依存関係を扱うことを可能にした注意機構は、タスクを解くため系列のどの情報に注目すべきかを適切に選択する役割を果たしています.
  - 順伝播ニューラルネットワークは入力層・中間層・出力層からなる多層 パーセプトロンであり、情報を保存しそれらを想起するための記憶の役割 を果たしています。
  - トランスフォーマでは、これらの処理を多段に繰り返すことで、必要な情報の選択と過去の情報との関連付けを行い、それらの情報をもとにタスクを解くということを行なっています。
- トランスフォーマで行われる処理は並列化可能であるため, GPUのようなハードウェアを用いて効率的にその処理を行うことができます.

#### トランスフォーマ

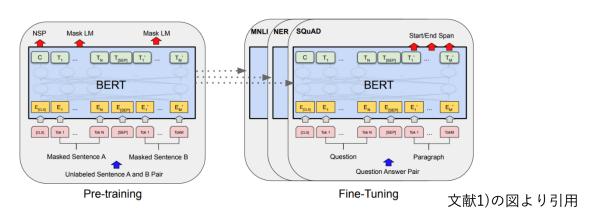
トランスフォーマの構造



1) A. Vaswani, et al. Attention Is All You Need. arXiv:1706.03762, 2017.

## トランスフォーマと大規模言語モデル

- 2018年に提案されたBERT (Bidirectional Encoder Representations from Transformers) <sup>1)</sup>は、トランスフォーマのエンコーダ部分のアーキテクチャをもとにした言語モデルです。
- BERTでは、テキスト中の単語の穴埋め問題(masked language model)および与えられた対の文が元々連続していたかという次文予測の問題を解くという自己教師あり学習で事前学習を行います。
- BERTの事前学習済みモデルを, さまざまな自然言語処理のタスクでファインチューニングすることで, モデルをそれらの下流タスクに適応させることができるようになりました.



1) J. Devlin, et al. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. arXiv:1810.04805, 2018.

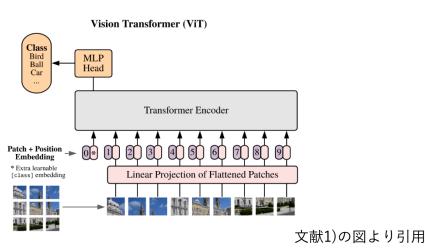
#### トランスフォーマと大規模言語モデル

- 2018年に提案されたGPT(Generative Pre-Trained Transformer)<sup>1)</sup>は、トランスフォーマのデコーダ部分のアーキテクチャーをもとにした言語モデルです。
- GPTでは、与えられた文脈に続く単語を予測するという問題を解くという自己教師あり学習で事前学習を行います.
  - このように事前学習されたモデルは自己回帰(autoregressive)言語モデルと呼ばれます。
- BERT同様, GPTの事前学習済みモデルをさまざまな自然言語処理の 下流タスクに適応させることができるようになりました.
- GPTはその後,事前学習に用いるデータやモデルの規模を大きくすることで現在の大規模言語モデル(LLM)の基盤となるGPT-2(2019年),GPT-3(2020年)へと発展しました.
- 2022年には、GPT-3.5のLLMをもとにしたChatGPTのサービスが OpenAI社から公開されました。

<sup>1)</sup> A. Radfold, et al. Improving language understanding by generative pre-training. Technical Report, 2018.

## Vision Transformer (ViT)

- Vision Transformer (ViT)¹¹)では、画像をパッチに分割してそれらを トークンとしてトランスフォーマのエンコーダに入力し、画像分類タ スクを解く事前学習を行います。
- ViTは画像データでの自己教師あり学習(例えば画像の欠落を復元するなど)による事前学習に用いることもできます.
- ViTの事前学習済みモデルをさまざまな画像認識のタスクでファイン チューニングすることで、モデルをそれらの下流タスクに適応させる ができます。

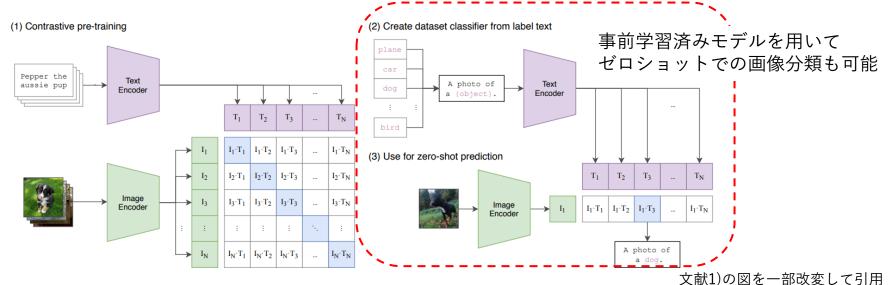


1) A. Dosovitskiy, et al. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. arXiv:2010.11929, 2020.

#### **CLIP**

- CLIP (Contrastive Language-Image Pretraining)<sup>1)</sup> は、大規模な画像とテキストのペアをもとに学習されたマルチモーダル(画像-言語)の基盤モデルです。
- テキスト

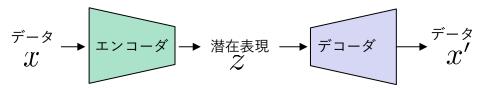
  →画像の生成へ応用されています。
  - テキストと画像それぞれをエンコーダ(例えばトランフォーマとViT)で 内部表現(ベクトル)に変換します。
  - 対となる画像とテキスト間のベクトルの内積が大きくなるように事前学習を行います.



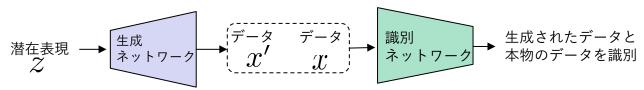
1) A. Radford, et al. Learning Transferable Visual Models From Natural Language Supervision. arXiv:2103.00020, 2021.

#### 生成モデル

- データの同時確率分布をモデル化する手法は、生成モデルと呼ばれます.
- 生成モデルにより新たなデータサンプルを生成することができます。
  - 特に深層学習の技術に基づく深層生成モデルでは、深層ニューラルネット ワークにより、データとその潜在表現と間の複雑な写像を学習します。
  - 深層生成モデルの技術は生成AIにも活用されてます.
    - 代表的な深層生成モデル
      - 変分オートエンコーダ (Variational Autoencoder, VAE)



• 敵対的生成ネットワーク(Generative Adversarial Networks, GAN)



• フローベースモデル(Flow-based Model)

