

2023.08.21



対話型文章生成AIの構築技術

— 各構成技術の役割や導入背景などの解説 —



東北大学
データ駆動科学・AI教育研究センター
大学院情報科学研究科 人工知能基礎学講座
鈴木潤

自己紹介

- 名前：鈴木 潤

専門分野：人工知能，機械学習，自然言語処理

- 2001.04 - 2018.03 日本電信電話株式会社 コミュニケーション科学基礎研究所
- 2018.04 - 2020.06 東北大学 大学院情報科学研究科 乾・鈴木研究室 准教授
- 2020.07 - 現在 東北大学 データ駆動科学・AI教育研究センター 教授
大学院情報科学研究科 人工知能基礎学講座 (協力講座)
(研究室は2021.4から)
- 2020.04 - 2022.04 Google LLC (Visiting Researcher) 

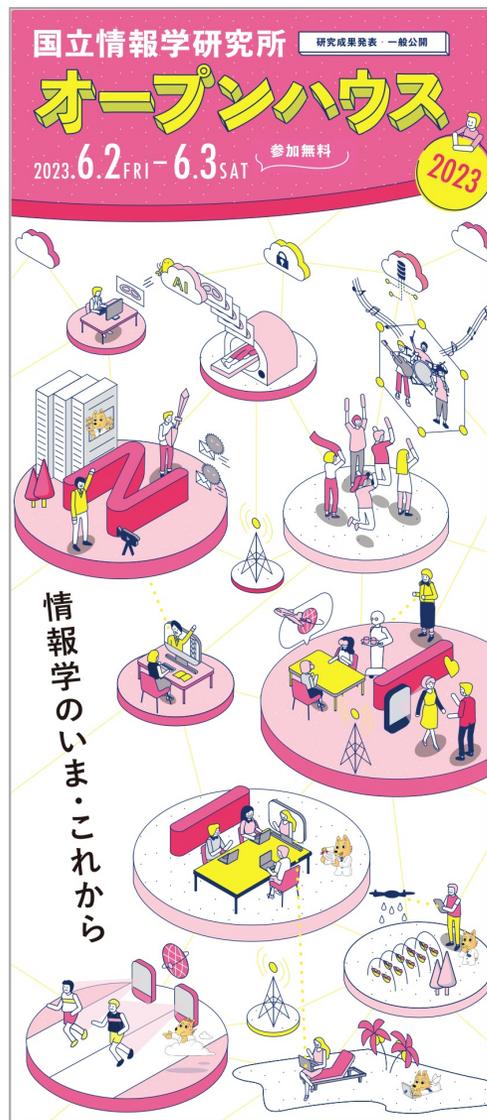
主に計算機が言語を効率的に学習する方法についての研究に従事

国立情報学研究所 研究成果発表 一般公開

オープンハウス

2023.6.2 FRI - 6.3 SAT 参加無料

2023



情報学のいま・これから

day1 6.2 FRI 13:00 - 18:30

13:00-15:00 一橋講堂 [※] (出席中継あり) for どなたでも
NII 活動報告
黒橋 祐夫 国立情報学研究所長

基調講演「ChatGPTを考える」
鈴木 潤 東北大学 データ駆動科学・AI 教育研究センター 教授
鈴木 久美 理研 AIP 言語情報アクセス技術チーム テクニカルアドバイザー
生員 直人 一橋大学大学院情報学研究所 システム専攻 教授

15:00-17:00 中会議場 for どなたでも
デモ・ポスターセッション
NII の研究者と議論可能なリアル会場メインのポスター発表

17:00-18:30 一橋講堂 [※] (出席中継あり) for 企業の方 大学の方
産官学連携セミナー
「フェイクメディア検出の最前線」

day2 6.3 SAT 10:30 - 17:30

13:00-14:30 一橋講堂 [※] (出席中継あり) for どなたでも
企画セッション
#情報研に聞きたい
研究者が語る情報学研究のいま！
Twitter ハッシュタグ #情報研に聞きたい で質問お待ちしております！

16:30-17:30 一橋講堂 [※] (出席中継あり) for 大学生 高専生 社会人
総合研究大学院大学 情報学コース
大学院説明会
NII で博士をとる！

10:30-16:30 中会議場ほか for 小学校 修一高学年
コンピュータサイエンスパーク
コンピュータを使わないで学ぶ遊び場！からだを使って
プログラミングを楽しもう！

10:30-16:30 中会議場 for どなたでも
デモ・ポスターセッション
研究者に研究内容を聞いてみよう！体験コーナーもあるよ！
デモ・ポスターセッションのコンテンツに限定すると
はのワークショップをブースで楽しんでほしい！

イベント参加方法 (参加無料)

詳細情報は
Webをチェック！

※ 学研総合センター (リアル会場) での開催です。一部プログラムは、オンラインでも開催します。参加にはイベントサイトへの登録が必要です。リアル会場での参加はチケット申し込み(有料) が必要。申込み詳細はイベントサイトをご覧ください。

会場：学研総合センター1階・2階 (東京都千代田区一ツ橋 2-1-2)
開催期間：6/2 (金) 13:00-18:30、6/3 (土) 10:30-17:30
開催方法：ハイブリッド開催 (リアル会場メイン・一部オンライン中継あり)
WEB サイト：https://www.nii.ac.jp/openhouse/
お問合せ・受付メールアドレス：oh@nii.ac.jp
TEL：03-4212-2131 後援：千代田区
〒101-8430 東京都千代田区一ツ橋 2-1-2 情報科学連携センター

NII OPEN HOUSE 2023

NII 国立情報学研究所
National Institute of Informatics

NIIオープンハウス 2023 基調講演
<https://www.nii.ac.jp/event/openhouse/2023/>

一般向け (非専門家向け)

言語処理学会第29回年次大会 (NLP2023)
 緊急パネル：ChatGPTで自然言語処理は終わ
 るのか？
https://www.anlp.jp/nlp2023/#special_panel

緊急パネル：ChatGPTで自然言語処理は終わるのか？

言語処理学会理事会主催によるパネルセッションを、3月14日(火) 13:10-13:50にH会場 (劇場ホール) で開催します。NLP2023に参加登録を行っている方が参加できます。

概要

大規模言語モデルの発展によって自然言語処理 (NLP) の方法論は大きく様変わりした。中でもOpen AIから発表されたChatGPTは言語モデルの利用方法を飛躍的に広げ、世間からも大きな関心を集めている。ChatGPTに代表される巨大言語モデルの出現でNLPはどう変わるか、はたまた終焉をむかえるのか？本パネルでは、NLP研究の一端で活躍されている5名のパネリストをお招きし、ChatGPTに関するフアクトの共有をはかるとともにNLP研究の今後についてご議論いただく。我々NLP/言語研究に携わる者自身が今起こっていることに対する理解を深め、今後なすべきことについて考えをめぐらす機会としたい。

動画・講演資料

- 言語処理学会YouTubeチャンネルにて本セッションの動画を公開しています

Top

お知らせ

開催案内

開催日時

会場

大会参加マニュアル

参加者限定Slackワークスペース

主催

スポンサー

参加登録

参加登録用ウェブサイト

参加費

【大会プログラム】

【大会発表要項】

招待講演

緊急パネル

チュートリアル

テーマセッション

ワークショップ

予稿集

表彰

表彰一覧

優秀賞・若手奨励賞

言語学賞

委員特別賞

スポンサー賞

優秀賞・若手奨励賞審査員

スポンサー展示

スポンサーイベント

共有

NLP2023 緊急パネル: ChatGPTで自然言語処理は...
言語処理学会理事会主催 緊急パネル
緊急パネル: ChatGPTで自然言語処理は終わるのか？

- ファシリテーター
 - 野澤太郎氏 (東北大)
- パネリスト
 - 黒橋 祐夫氏 (京大)
 - 相良 美穂氏 (ソバパ)
 - 佐藤 敏成氏 (LINE)
 - 鈴木 潤氏 (東北大)
 - 谷中 龍氏 (東大)
- 3月14日(火) 13:10-13:50, H会場 (劇場ホール)
 - Slack: #0314-1310-緊急パネル-h会場 (質問はこちらへ)
 - #ChatGPTで自然言語処理は終わるのか

見る YouTube

- 講演資料はこちら



PAKDD 2023

THE 27TH PACIFIC-ASIA CONFERENCE ON KNOWLEDGE DISCOVERY AND DATA MINING
25-28 May, 2023, Osaka, Japan
GRAND FRONT OSAKA

Organizations Calls Program Attend Participants Only Sponsors Awards

T2: A Gentle Introduction to Technologies Behind Language Models and Recent Achievement in ChatGPT (9:00–11:00)

Speakers:
Jun Suzuki (Tohoku University), Naoaki Okazaki (Tokyo Institute of Technology), Kyosuke Nishida (NTT Corporation)

Abstract:
Language models (LMs) have a long history in natural language processing (NLP) research. Their usage was mainly a text generation module in machine translation and speech recognition systems, used together with translation models or acoustic models. After the current neural era, LMs take a more essential role in the NLP field. In fact, LMs are integrated into any models/systems to tackle almost all the NLP tasks and provide state-of-the-art performance on conventional benchmarks. The usage of LMs is considered to be shifting to more like a world model of languages or a general-purpose feature generator of any language-related tasks. More recently, the public sometimes treats LMs like ChatGPT, and its successor GPT-4, as general-purpose AI after starting an online service in the public domain. This tutorial will first introduce some introductory topics we should know when discussing the recent advances in LMs like ChatGPT. We will then briefly introduce the technologies behind ChatGPT-like LMs. Additionally, we also provide ChatGPT's social impacts discussed recently in public.

Tutorial Notes:
[Part 1 & 2](#) [Part 3 & 4](#) [Part 5](#)

PAKDD 2023 Tutorial
<https://pakdd2023.org/tutorials/>

非NLP研究者向け

- ChatGPTに代表される対話型文章生成AIが、研究分野のみならず産業界や一般社会においても注目を集めている。対話型文章生成AIはニューラルネットによる言語モデルを基本（出発点）として幾つかの要素技術の集合体として構築されている。本講演では、対話型文章生成AIを構築する際に用いられる要素技術の集合をおおまかに整理して紹介する。また、**各技術の役割や導入背景などを説明**し、どのような過程を経て現在の対話型文章生成AIに到達したのかについて説明する

[注意事項] 事前のお断り

- 本資料には講演者独自の解釈が含まれています
 - ここに書かれていることは必ずしも一般的に普及している知見ではない場合があります
 - 議論や洞察を深めるためのきっかけとしてあえて書いている部分もあります

- 本資料の情報の正確性について
 - 対話型文章生成AI / ChatGPT 関連の情報は**頻繁に更新**
 - あくまで**2023.08.21**時点の情報

本資料の情報が時間経過と共に正しくなくなる可能性があることに注意

対話型文章生成AIとは？

一言で言うと...

[一般向け]

(人が人に話すような) **対話形式の指示**を受け付け
その指示に適した文章を生成する **文章生成器**

現在の代名詞：ChatGPT

[参考] 年単位のNLP分野の重要技術

- 
- 2013 Neural word embeddings / 単語分散表現
 - 2014 Neural encoder-decoder / 系列変換器
 - 2015 Attention mechanism / 注意機構
 - 2016 Subword / サブワード分割
 - 2017 Transformer / 自己注意型ネットワーク
 - 2018 Neural language models / ニューラル言語モデル
 - 2019 Masked neural LMs / マスク型言語モデル
 - 2020 Large LMs: GPT-3 / 巨大言語モデル
 - 2021 Prompt tuning, engineering / プロンプト開発
 - 2022 Instruction (Chat) tuning / 言語生成品質の向上
 - 2023 (???)

今日の話題

- 文章生成AI (言語モデル) がうまくいく第一要因は？
- 言語モデルに Transformer は たまたま？
- プロンプト (汎用/非特化型) のきっかけは？
- (LLM構築の取り組み)

- (昨今の)「言語モデル」が文章生成AIとしてうまく機能している第一要因は？

[参考] ChatGPTを構成する技術

- 基盤
 - 言語モデル
- 成功の要因
 - ① 多層ニューラルネットワーク (DNN) の利用
 - ② 大規模化 (パラメタ数/データ量)
 - ③ 指示文設計 (プロンプトエンジニアリング)
 - ④ 指示文 (+対話文) チューニング
 - ⑤ 人手点数付け結果の活用

[参考] ChatGPTを構成する技術

2階部分 (指示文の活用/学習)

指示文設計
指示文 (+対話文) 微調整学習
人手点数付けの利用



獲得するスキル?

指示文の理解

対話的なやり
とりの理解

不適切発言の
抑制

1階部分 (言語モデルの学習)

ニューラルネット
巨大言語モデル
大規模データ
(教師あり学習, 確率モデル)



文章の流暢さ
世界の知識



基盤：言語モデル

- 文章の出現確率を予測する確率モデル

語彙
(単語の集合)

に
。
これ
です
今日
1
…
仙台
良い
東京
は
…
天気
…

$P(Y = \$BOSS\$ \text{今日は良い天気です。}\$EOS\$)$

高

$P(Y = \$BOSS\$ \text{良いです。天気は今日}\$EOS\$)$

低

言語
モデル

確率は大量のデータ
から何かしらの形で
計算



言語モデルの例: n -gram 言語モデル

- Example: Google n -gram

BLOG >

All Our N-gram are Belong to You

THURSDAY, AUGUST 03, 2006

Posted by Alex Franz and Thorsten Brants, Google Machine Translation Team

Here at Google Research we have been using word **n -gram models** for a variety of R&D projects, such as **statistical machine translation**, speech recognition, **spelling correction**, entity detection, information extraction, and others. While such models have usually been estimated from training corpora containing at most a few billion words, we have been harnessing the vast power of Google's datacenters and distributed processing **infrastructure** to process larger and larger training corpora. We found that there's no data like more data, and scaled up the size of our data by one order of magnitude, and then another, and then one more - resulting in a training corpus of **one trillion words from public Web pages.**

We believe that the entire research community can benefit from access to such massive amounts of data. It will advance the state of the art, it will focus research in the promising direction of large-scale, data-driven approaches, and it will allow all research groups, no matter how large or small their computing resources, to play together. That's why we decided to share this enormous dataset with everyone. We processed **1,024,908,267,229 words** of running text and are publishing the counts for all **1,176,470,663 five-word sequences** that appear **at least 40 times**. There are **13,588,391 unique words**, after discarding words that appear **less than 200 times**.

<https://ai.googleblog.com/2006/08/all-our-n-gram-are-belong-to-you.html>

Context $Y_{<j}$	Target word y_j	Count	Probability
serve as the	incoming	92	0.00049
serve as the	independent	794	0.00423
serve as the	index	223	0.00119
serve as the	indication	72	0.00038
serve as the	indicator	120	0.00064
serve as the	indicators	45	0.00024
serve as the	indispensable	111	0.00059
serve as the	indispensible	40	0.00021
serve as the	individual	234	0.00125
serve as the	industrial	52	0.00028
serve as the	industry	607	0.00324
serve as the	info	42	0.00022
serve as the	informal	102	0.00054
serve as the	information	838	0.00447
serve as the	informational	41	0.00022
serve as the	initial	5331	0.02843
serve as the	initiating	125	0.00067
serve as the	initiation	63	0.00034
serve as the	initiator	81	0.00043
serve as the		

→ Total 187491

昔の言語モデルの使い方

- モデルが生成した候補文の尤もらしさを計測

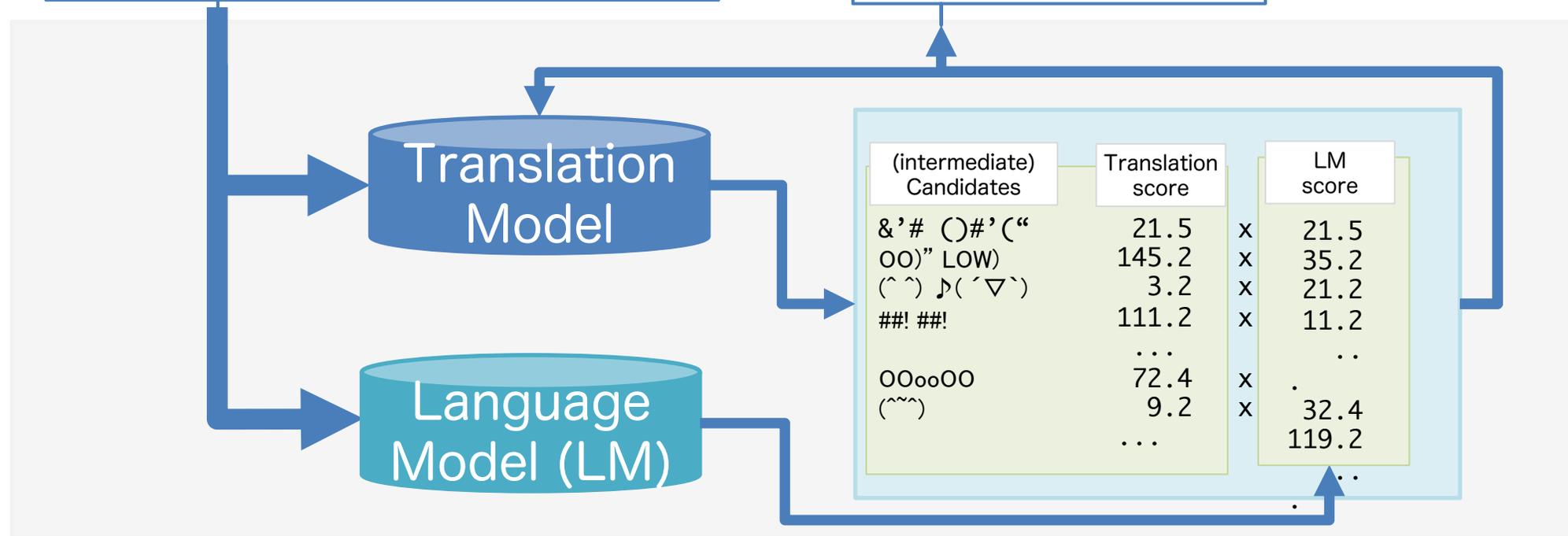
例：統計翻訳 (Statistical Machine Translation: SMT)

Input text

HI YWUY1 KKIUH WKUYN LO WUNI

Output text

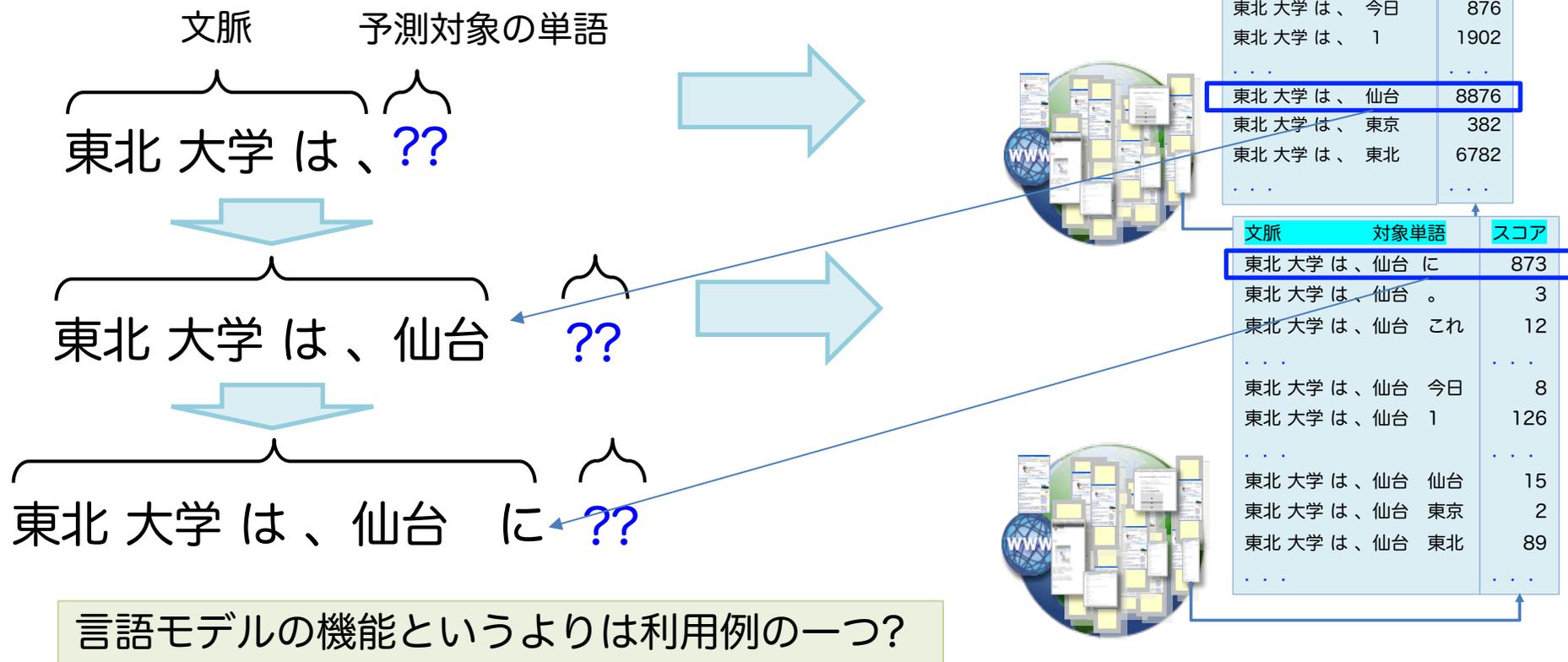
(^ ^) ♪('▽`) ##!



言語モデルによる文生成

- 文章の先頭から1単語ずつ次の単語を予測

=> 予測した単語は文脈として再利用
生成が終わるまで繰り返す

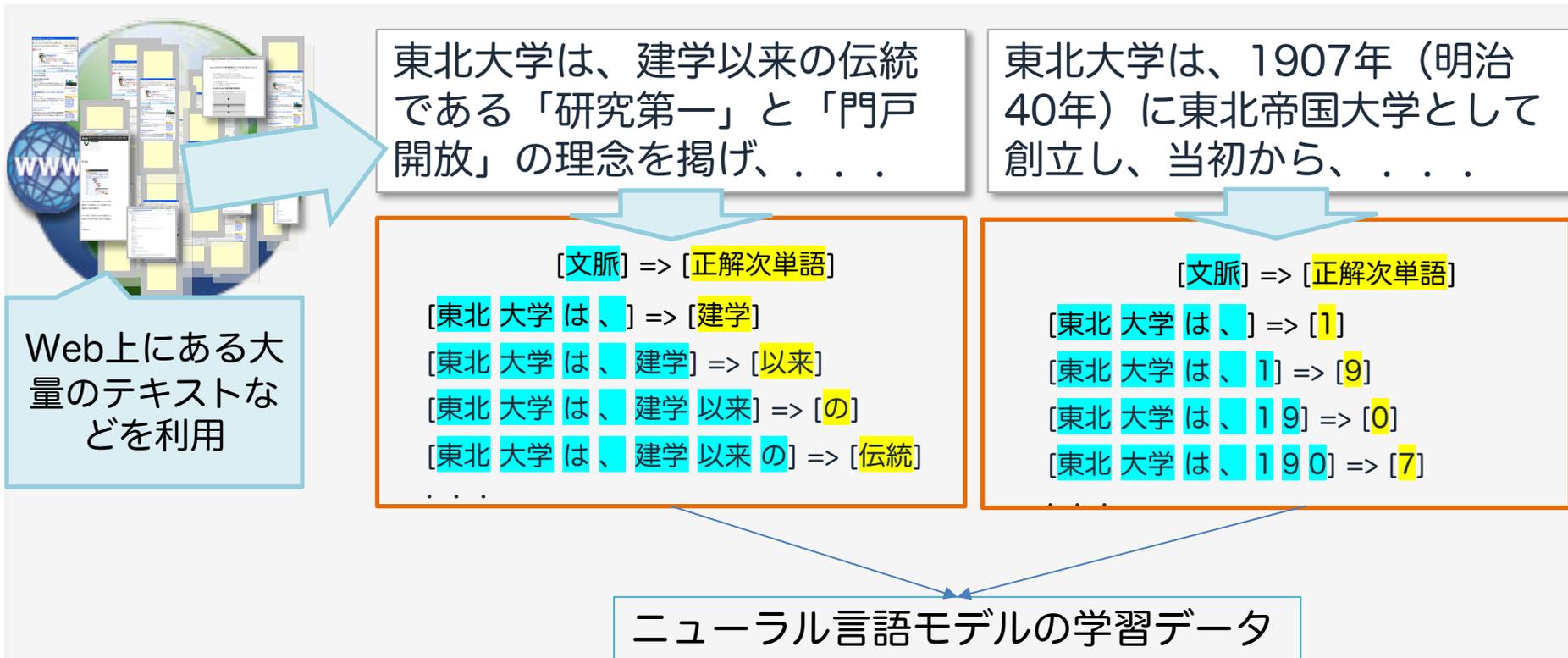


ニューラル言語モデル

DNN: deep neural network



- 言語モデル (確率モデル) をDNNで近似
 - [文脈] => [正解次単語] を分類問題として学習
 - 文章中の**全単語を対象**として学習



ニューラル言語モデル



TOHOKU UNIVERSITY

$$\operatorname{argmax}_{y_j} P_{\theta}(y_j | Y_{<j})$$

Vocabulary	Vocabulary	Vocabulary	Vocabulary
A	A	A	A
this	this	this	this
that	that	that	that
...
meet	meet	meet	meet
have	have	have	have
you	you	you	you
...
Nice	Nice	Nice	Nice
...
to	to	to	to
...
too	too	too	too
,	,	,	,
.	.	.	.

OUTPUT text~>

Nice to meet you , too ...

Neural LM (e.g., Generative pre-trained transformer: GPT)



文章生成AIがうまくいく第一要因は？

- 成功の要因は「言語モデル」の能力/機能？
 - 言語モデルはただの確率モデルのはず
 - n -gram LMをどんなに巨大化してもChatGPTにはならない

文章生成AIがうまくいく第一要因は？

- 成功の要因は「言語モデル」の能力/機能？
 - 言語モデルはただの確率モデルのはず
 - n -gram LMをどんなに巨大化してもChatGPTにはならない
- 何が実際の要因なのか？
 - 言語（単語）の意味や使われ方の観点で似ている/似ていないが判断できるようになったこと
=> 分散表現 (単語埋め込み) でも語られていた能力
 - 言語の「世界モデル」として覚えておくのに適しているモデルだったこと
 - Transformerだったのも運が良かった？

- 言語モデルに Transformer はたまたま？

[参考] Transformer

- 深層ニューラルネットワークの種類の一つ
 - 2017年：主に機械翻訳用として提案
翻訳関連で最初に流行った

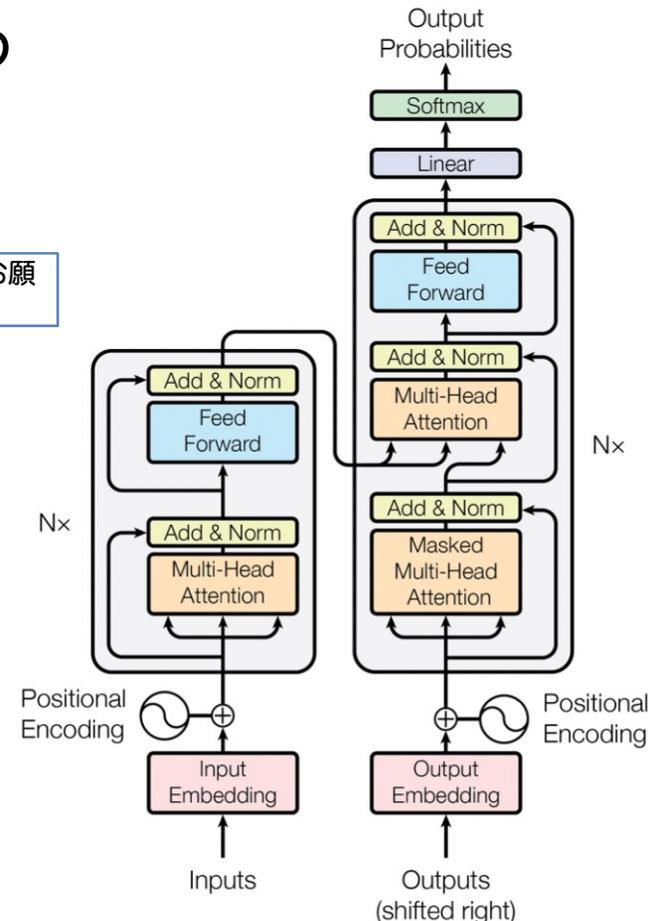
Attention Is All You Need

Ashish Vaswani* Google Brain avaswani@google.com	Noam Shazeer* Google Brain noam@google.com	Niki Parmar* Google Research nikip@google.com	Jakob Uszkoreit* Google Research usz@google.com
Llion Jones* Google Research llion@google.com	Aidan N. Gomez* † University of Toronto aidan@cs.toronto.edu	Łukasz Kaiser* Google Brain lukaszkaier@google.com	
Illia Polosukhin* † illia.polosukhin@gmail.com			



I'd like to have an aisle seat please.

通路側の席をお願いします



https://papers.nips.cc/paper_files/paper/2017/hash/3f5ee243547dee91fbd053c1c4a845aa-Abstract.html

- ニューラルネットに基づくAI技術において**最重要技術の一つ**



[参考] Transformer

- 2018: 言語モデルBERT (とGPT) のベースモデルとして採用

- => 言語の基盤モデルに

BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

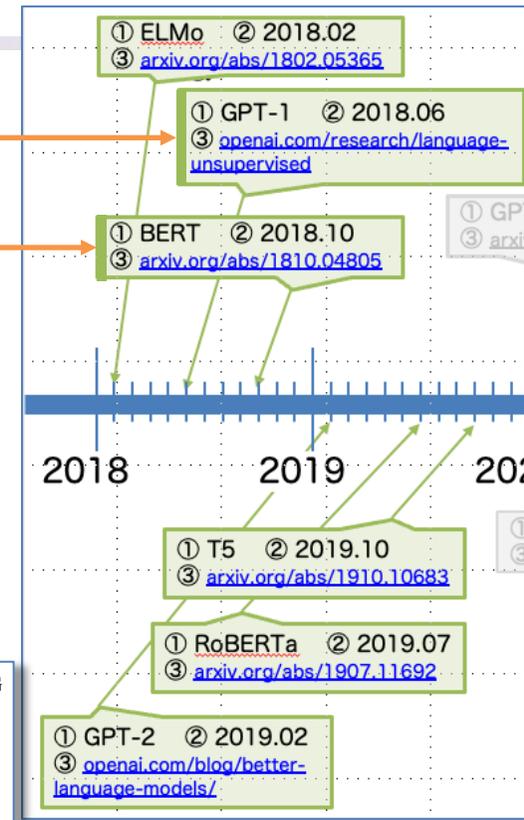
Jacob Devlin Ming-Wei Chang Kenton Lee Kristina Toutanova
Google AI Language

Improving Language Understanding by Generative Pre-Training

Alec Radford Karthik Narasimhan Tim Salimans Ilya Sutskever
OpenAI OpenAI OpenAI OpenAI
alec@openai.com karthikn@openai.com tim@openai.com ilyasu@openai.com

GENERATING WIKIPEDIA BY SUMMARIZING LONG SEQUENCES

Peter J. Liu*, Mohammad Saleh*, Etienne Pot†, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, Noam Shazeer
Google Brain
Mountain View, CA

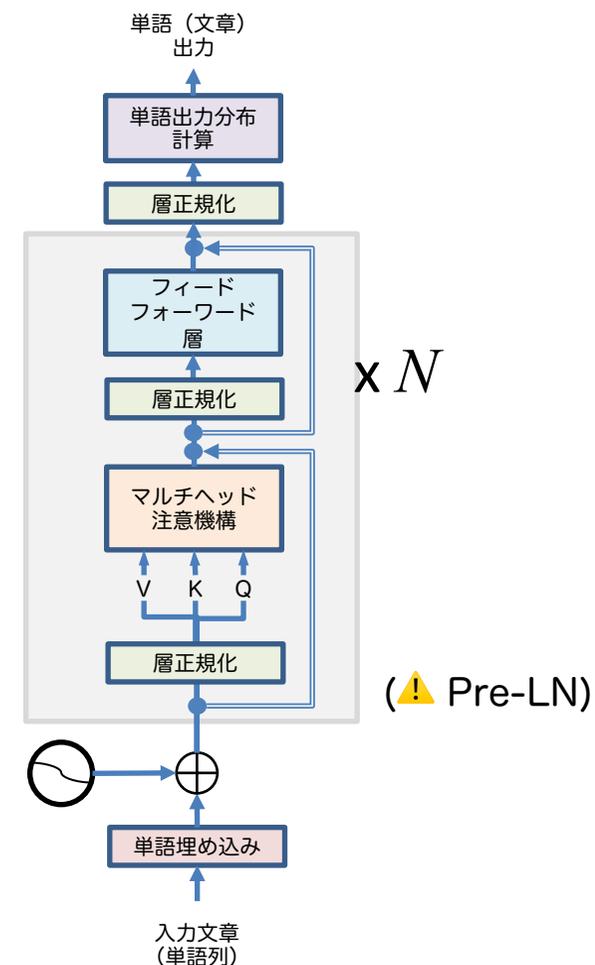


- 2023有名なニューラル言語モデルのほぼ全てで採用
 - 例：GPT-3 (ChatGPT/GPT-4), PaLM, LLaMA, OPT, BLOOM, ...
- 画像/音声/信号処理などでも広く普及

Transformer は言語モデルに適している？

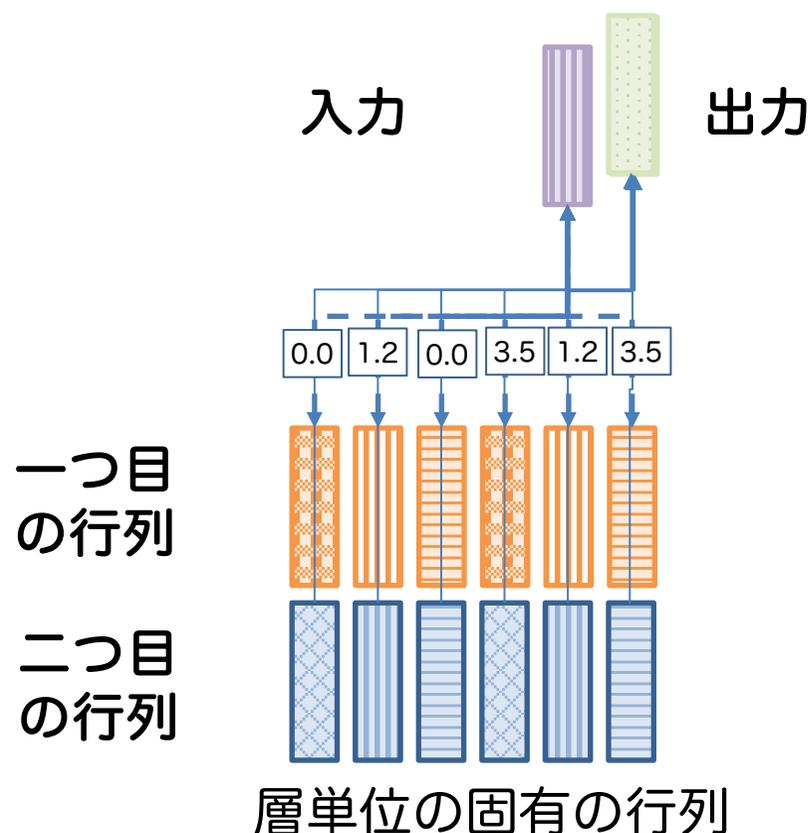
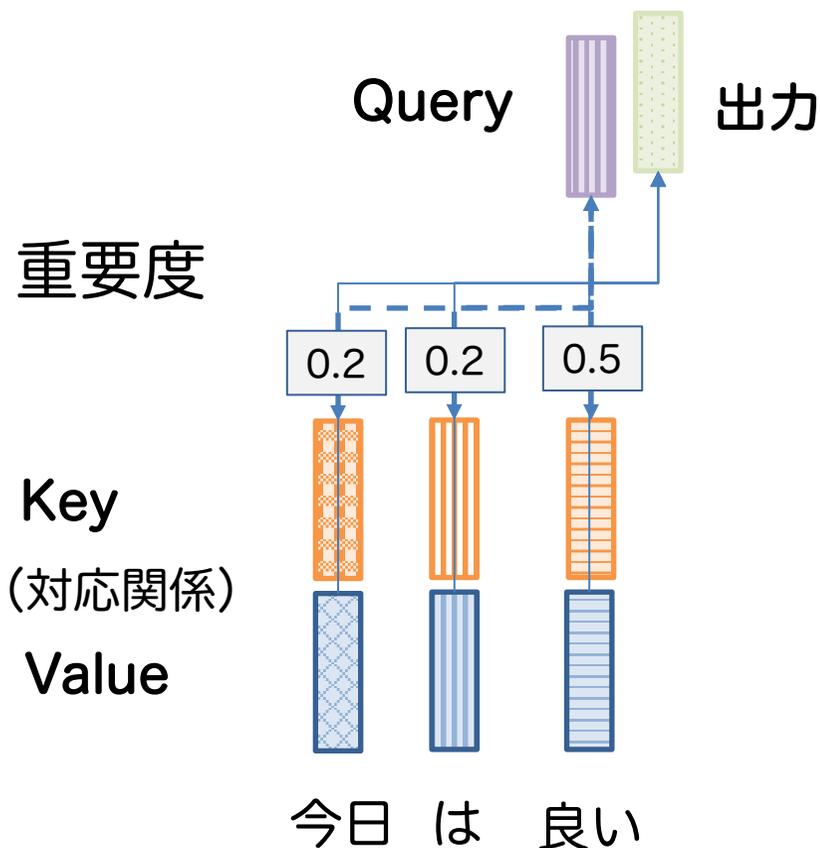
- 主にAttentionとFFNNで構成されている

GPT-3		12,288					
Encoder	in	out	bias	sub total	num	total	
embedding	50,257	12,288	0	617,558,016	1	617,558,016	
Position	12,288	4,096	0	50,331,648	1	50,331,648	
k	12,288	12,288	12,288	151,007,232	96	14,496,694,272	
v	12,288	12,288	12,288	151,007,232	96	14,496,694,272	
q	12,288	12,288	12,288	151,007,232	96	14,496,694,272	
out	12,288	12,288	12,288	151,007,232	96	14,496,694,272	
layernorm1	1	12,288	12,288	24,576	96	2,359,296	
ff1	12,288	49,152	49,152	604,028,928	96	57,986,777,088	
ff2	49,152	12,288	12,288	603,992,064	96	57,983,238,144	
layernorm2	1	12,288	12,288	24,576	96	2,359,296	
layernorm_last	1	12,288	12,288	24,576	1	24,576	
						174,629,425,152	



Transformer は言語モデルに適している？

- Attention / FFNN 共にメモリとその呼び出し処理と見做せる？



Transformer は言語モデルに適している？

- Attention / FFNN 共にメモリとその呼び出し処理と見做せる？
 - 学習データに過適応できる性質
 - 学習はなるべく手を抜く性質
 - Shortcut learning, superficial cues
 - (この辺りの性質が構造的な情報の学習をうまく実現する要因?)
- 「まる覚えする能力」が高いと想定できる
 - ただし、「文章」を丸覚えしているわけではなくベクトル列の情報を丸覚えしている（しようとしている）

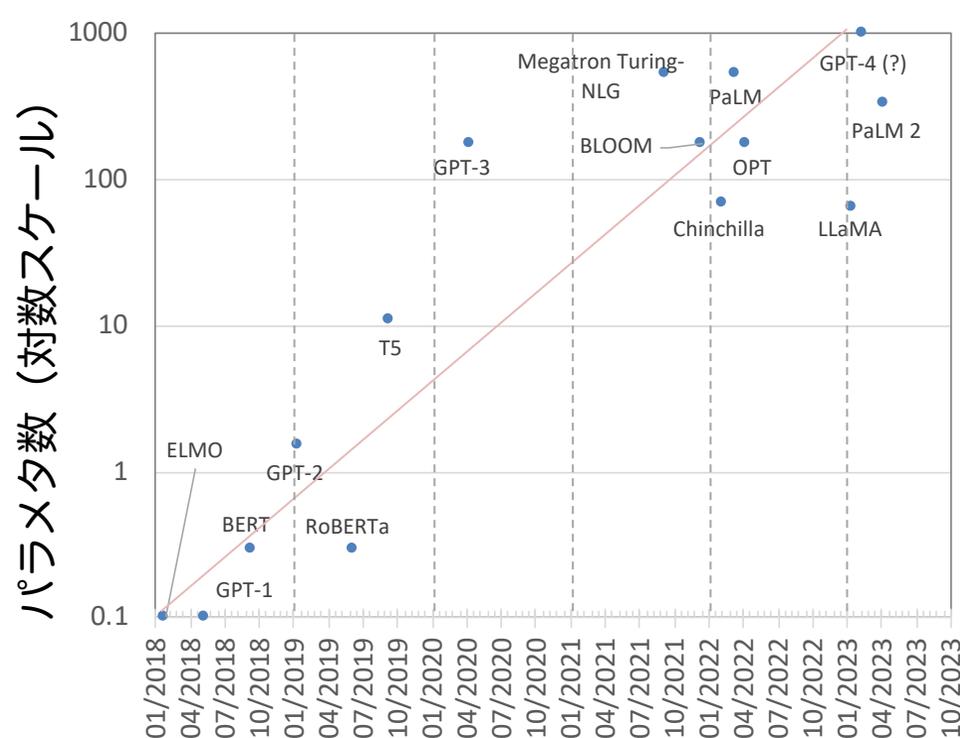
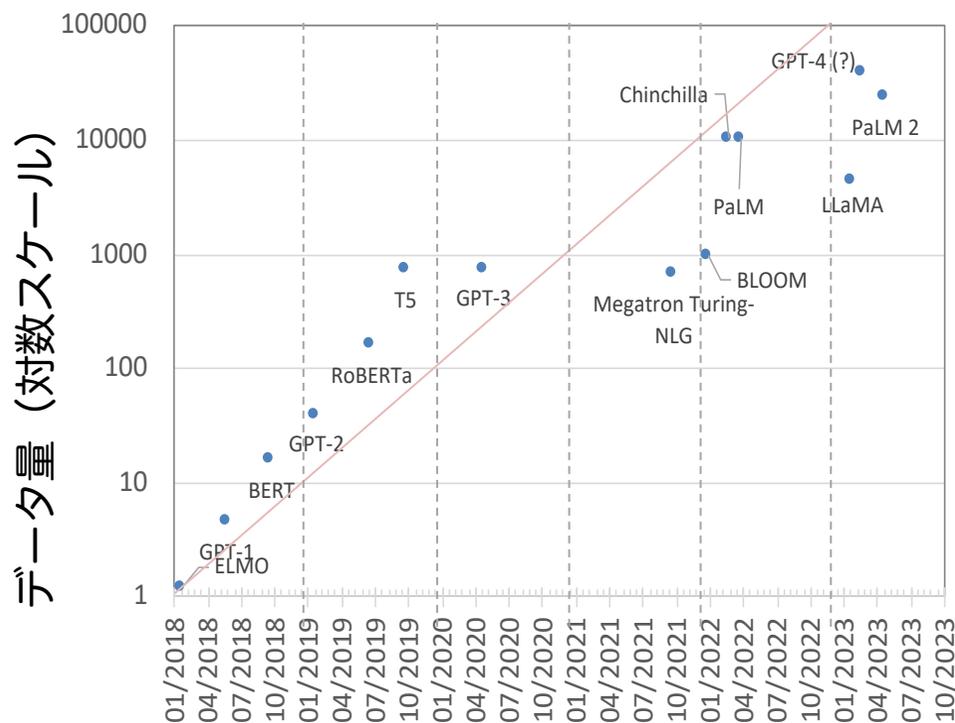
年毎の大規模化



● 学習データ量 / パラメタ数：年々順調に増加

	トークン数	ファイル容量
GPT-2:	0.01 B	0.04 TB
GPT-3:	0.5 B	0.75 TB
GPT-4:	13 T	20 TB (憶測)

	パラメタ数	必要メモリ量
GPT-2:	1.5 B	6 GB
GPT-3:	175 B	700 GB (ChatGPT)
GPT-4:	1,800 B	7,200 GB (憶測)

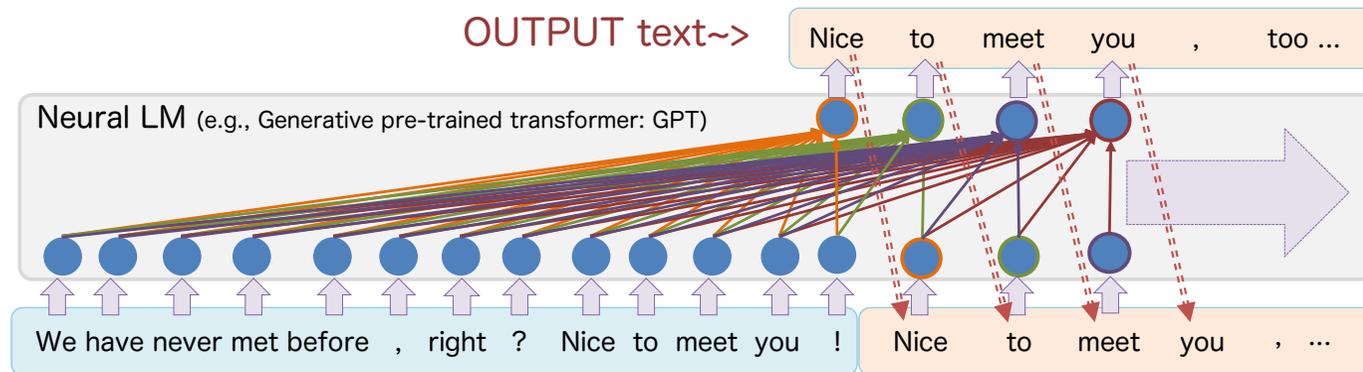


- プロンプト (汎用/非特化型) のきっかけは？

指示文 (プロンプト) とは？

文脈 = 指示文 (プロンプト) と考えてよい

- 言語モデル：与えられた**文脈**に基づいて文章の続きを生成 (文章の補完)

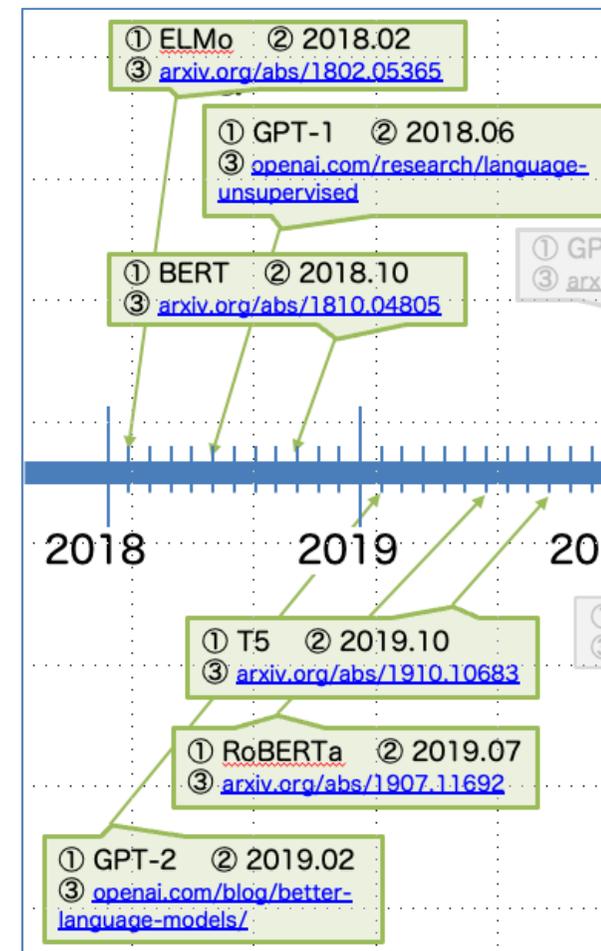
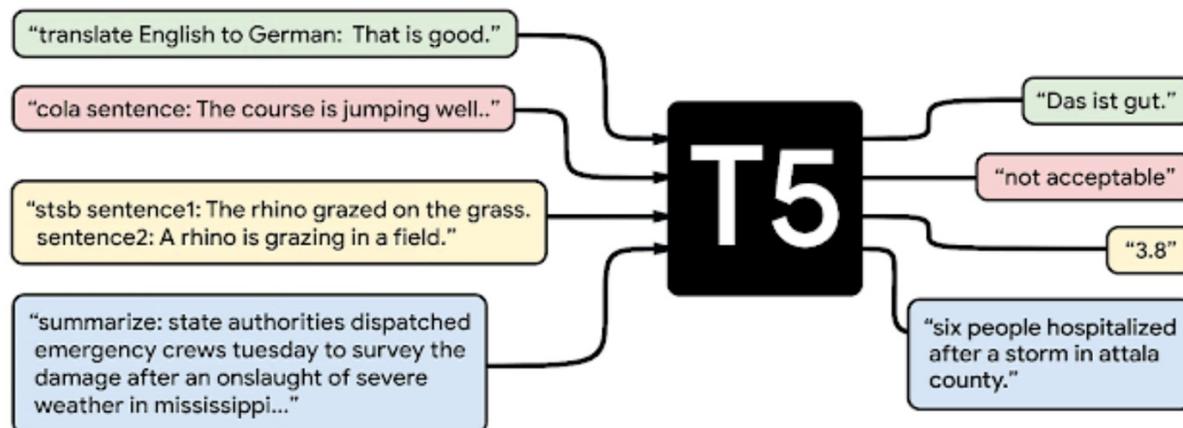


- **文脈**を工夫すると様々な文章を生成することが判明
- なぜ特別な用語を用いるのか
=> 単に「文脈」と呼ぶ以上の何か

[参考] 言語処理タスクの解き方

- 従来：それぞれのタスクに特化したモデル/方法の構築
 - 特定のタスクのみ効果的に解ける
- **一つのモデル**であらゆる言語処理タスクを解くことを目指す => **生成モデル**
 - 現在の対話型文章生成AIの先駆け (T5)

例：



Copied from <https://ai.googleblog.com/2020/02/exploring-transfer-learning-with-t5.html>

指示文設計 (プロンプトエンジニアリング)

- GPT-3: 指示文 (プロンプト) の考え方の転換点
 - Finetuningをしない / その代わりに事例を使う

Zero-shot

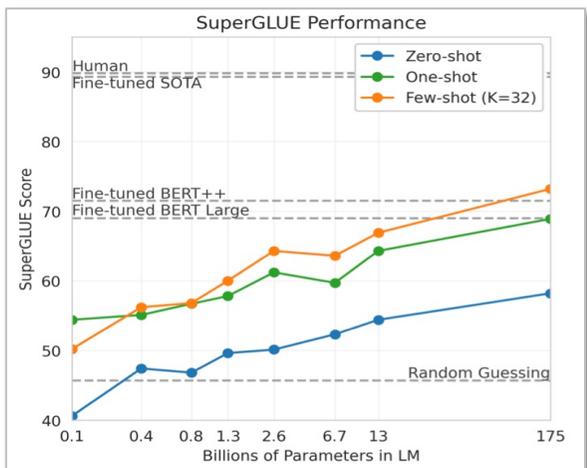
The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1 Translate English to French: ← task description
2 cheese => ..... ← prompt
```

Few-shot

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

```
1 Translate English to French: ← task description
2 sea otter => loutre de mer ← examples
3 peppermint => menthe poivrée ←
4 plush girafe => girafe peluche ←
5 cheese => ..... ← prompt
```



Language Models are Few-Shot Learners

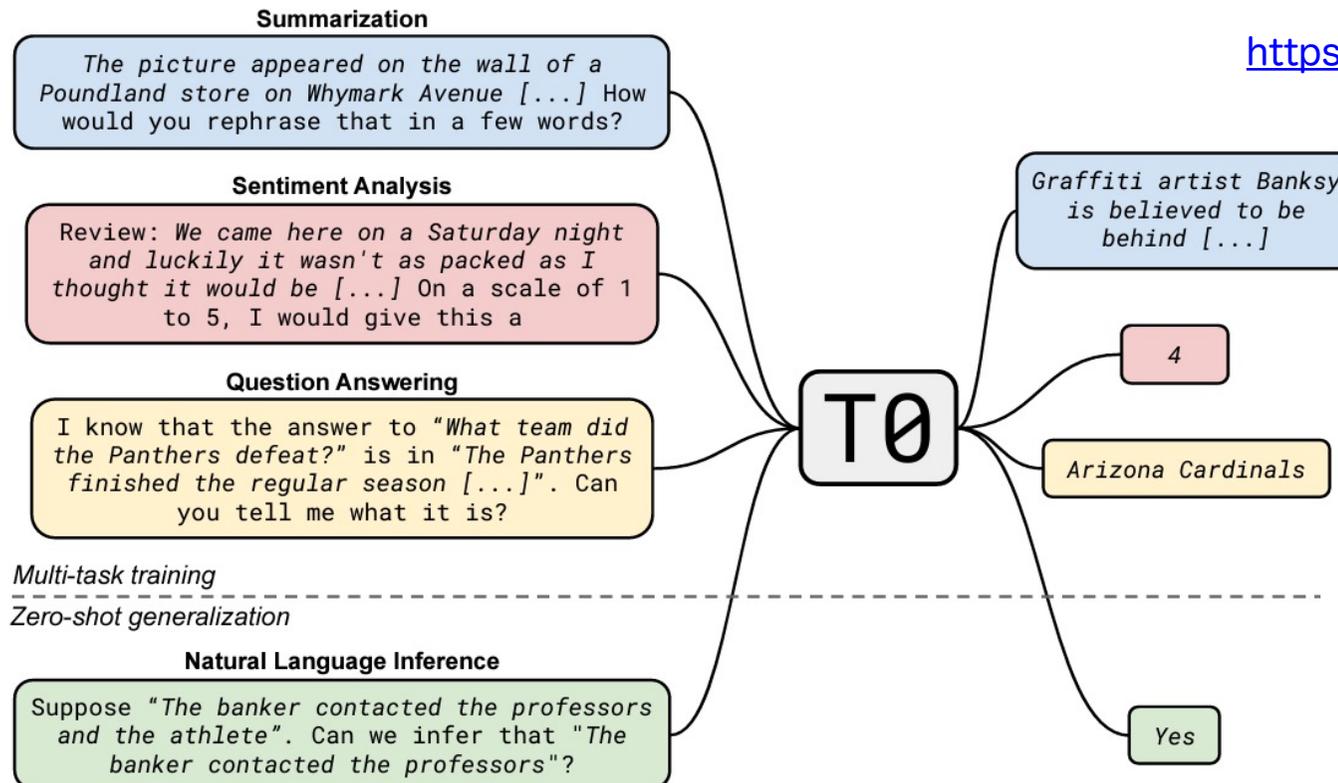
<https://arxiv.org/abs/2005.14165>

指示文設計 (プロンプトエンジニアリング)

- T0: 指示文をそのまま利用

Multitask Prompted Training Enables Zero-Shot Task Generalization

<https://arxiv.org/abs/2110.08207>



おおむね現在使われるプロンプトに近い状態

指示文設計 (プロンプトエンジニアリング)

- Chain-of-thought: 解き方を教える

Chain-of-Thought Prompting
Elicits Reasoning in Large
Language Models

<https://arxiv.org/abs/2201.11903>

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Example
(One-shot)

Question

Model Output

A: The answer is 27. ❌

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. ✅

まとめ

- プロンプト発展の経緯
=> 一つのモデルで多数のタスクに対応したい
- 時代とともに微妙に「プロンプト」が指すものが広がっている（気がする）

- (LLM構築の取り組み)



The screenshot shows the LLM-jp website with a navigation bar (News, 趣旨説明, 資料, メンバー, 参加申請, 連絡先) and a main content area. The main content area features a blue abstract graphic with vertical lines and a title 'LLM 勉強会'. Below the title is a paragraph of text and a bulleted list of activities.

LLM 勉強会

本勉強会では、自然言語処理および計算機システムの研究者が集まり大規模言語モデルの研究開発について定期的に情報共有を行っています。

具体的には、以下の目的で活動しています。

- オープンソースかつ日本語に強い大規模モデルの構築とそれに関連する研究開発の推進
- 上記に関心のある自然言語処理および関連分野の研究者によるモデル構築の知見や最近の研究の発展についての定期的な情報交換
- データ・計算資源等の共有を前提とした組織横断的な研究者間の連携の促進
- モデル・ツール・技術資料等の成果物の公開

<https://llm-jp.nii.ac.jp/news>

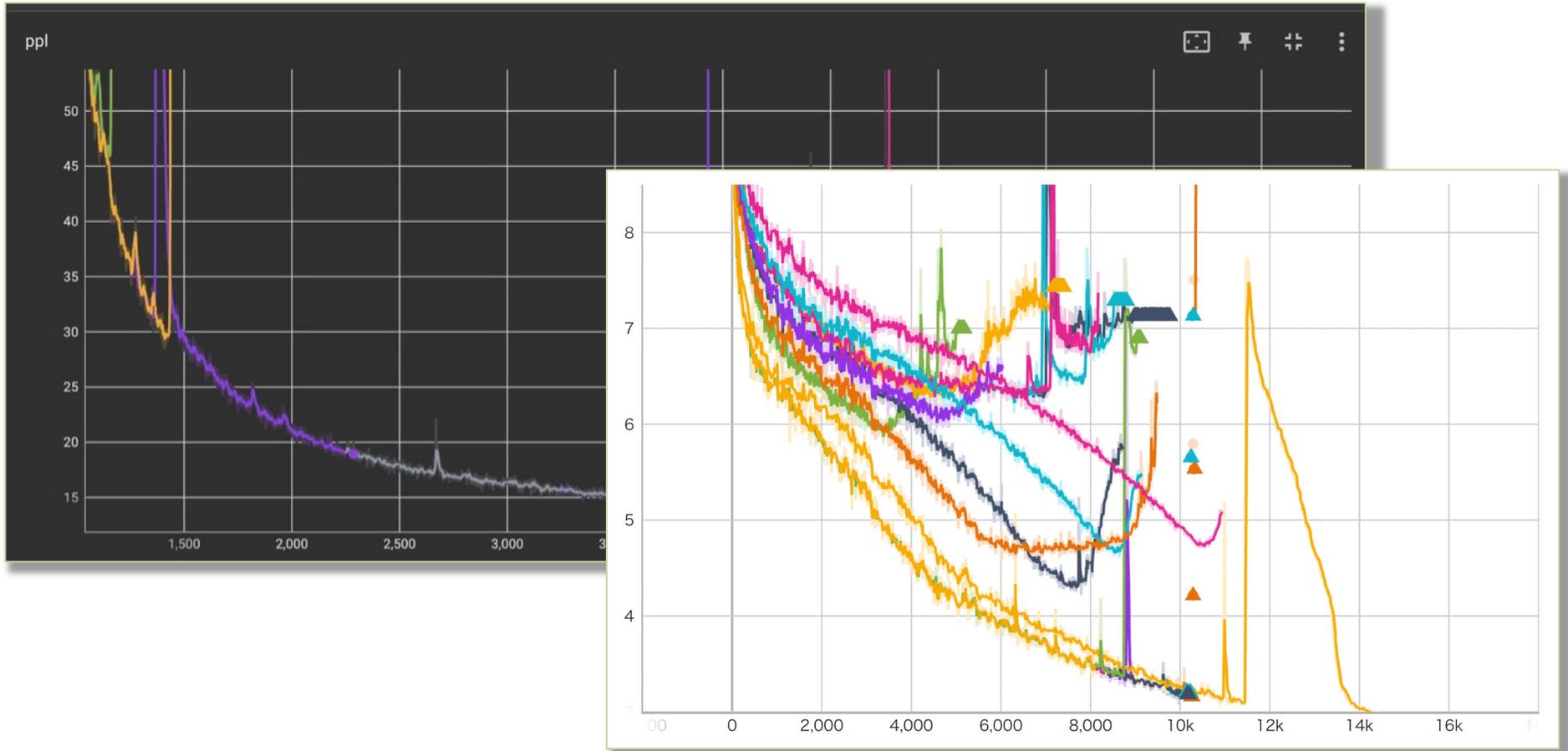
モデル構築検討WG幹事として活動

- LLM 勉強会
 - 本勉強会では、自然言語処理および計算機システムの研究者が集まり大規模言語モデルの研究開発について定期的に情報共有を行っています。
 - 具体的には、以下の目的で活動しています。
 - オープンソースかつ日本語に強い大規模モデルの構築とそれに関連する研究開発の推進
 - 上記に関心のある自然言語処理および関連分野の研究者によるモデル構築の知見や最近の研究の発展についての定期的な情報交換
 - データ・計算資源等の共有を前提とした組織横断的な研究者間の連携の促進
 - モデル・ツール・技術資料等の成果物の公開

LLMの構築

- LLMの学習はそれほど単純ではない

https://github.com/facebookresearch/metaseq/blob/main/projects/OPT/chronicles/10_percent_update.md



<https://huggingface.co/blog/bloom-megatron-deepspeed#bf16optimizer>

- まとめ

まとめ (1/2)

- ChatGPT (対話型文章生成AI) とは？

=> 対話形式の指示を受け付け

その指示に適した文章を生成する文章生成器

- 技術

- 基盤：言語モデル
- ① 深層ニューラルネットワーク (DNN) の利用
- ② 大規模化 (パラメタ数/データ量)
- ③ 指示文設計 (プロンプトエンジニアリング)
- ④ 指示文 (+対話文) チューニング
 - ⑤ 人手点数付け結果の活用

まとめ (2/2)

- 本日の話題
 - 文章生成AI (言語モデル) がうまくいく第一要因は？
 - 言語モデルに Transformer はたまたま？
 - プロンプト (汎用/非特化型) のきっかけは？
 - (LLM構築の取り組み)