

検索アイテムの付加価値 推定と LLM での応用

橋本 和真 / Kazuma HASHIMOTO

Google

NLP コロキウム @ 2024 年 10 月 30 日 (水)

公開用・改訂版 (2024 年 11 月 1 日)

自己紹介

- 2012～2018: 東京大学
 - 卒論～博士 (指導教員: 鶴岡 慶雅 先生)
- 2018～2021: Salesforce
 - AI Research
- 2021～現在: Google
 - Research (2021～2024 前半)
 - DeepMind (2024 後半～)

+ 技術系・研究の世界に興味のある学生と研究を行う非営利組織のボランティア活動

- 個人ページ:
<https://hassygo.github.io>
- メール:
kazumah@google.com
- LinkedIn:
Kazuma Hashimoto

アメリカのベイエリアが拠点なので、お越しの際はお気軽に声をかけてください。

今回の発表の流れ

1. 検索ベースの NLP モデルと応用に関する話
2. 企業の LLM 研究の中心地で働くことに関する話

検索ベースの NLP モデル

- 検索 (search, retrieval) を利用した NLP タスク・モデル
 - QA: query -> **document(s)** -> answer
 - MT: source text -> **translation memory** -> translation
 - ...
- RAG (retrieval-augmented generation) という形で更に普及
 - 昨今の LLM (large language model) と相性が良く、様々な問題において効果的に使用可能

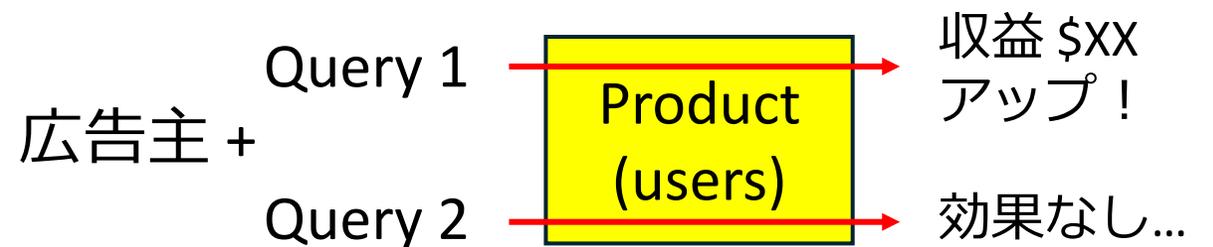
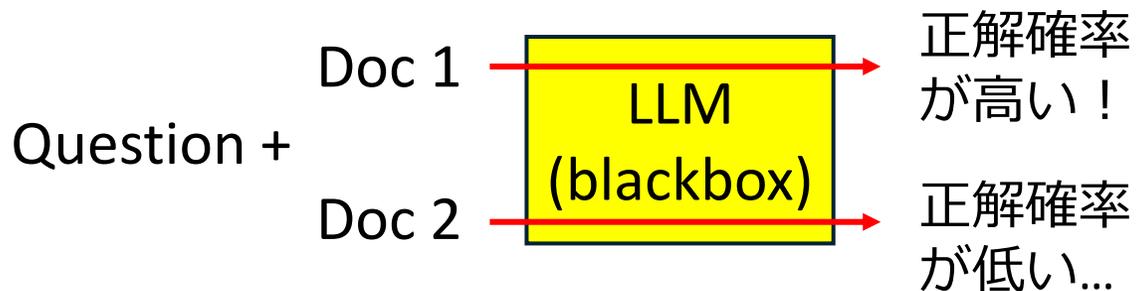
検索クエリ自体や結果の**妥当性 (価値)** はどのように測られるのか？

検索クエリ・結果の価値の推定

➤(半) 人手のアノテーション

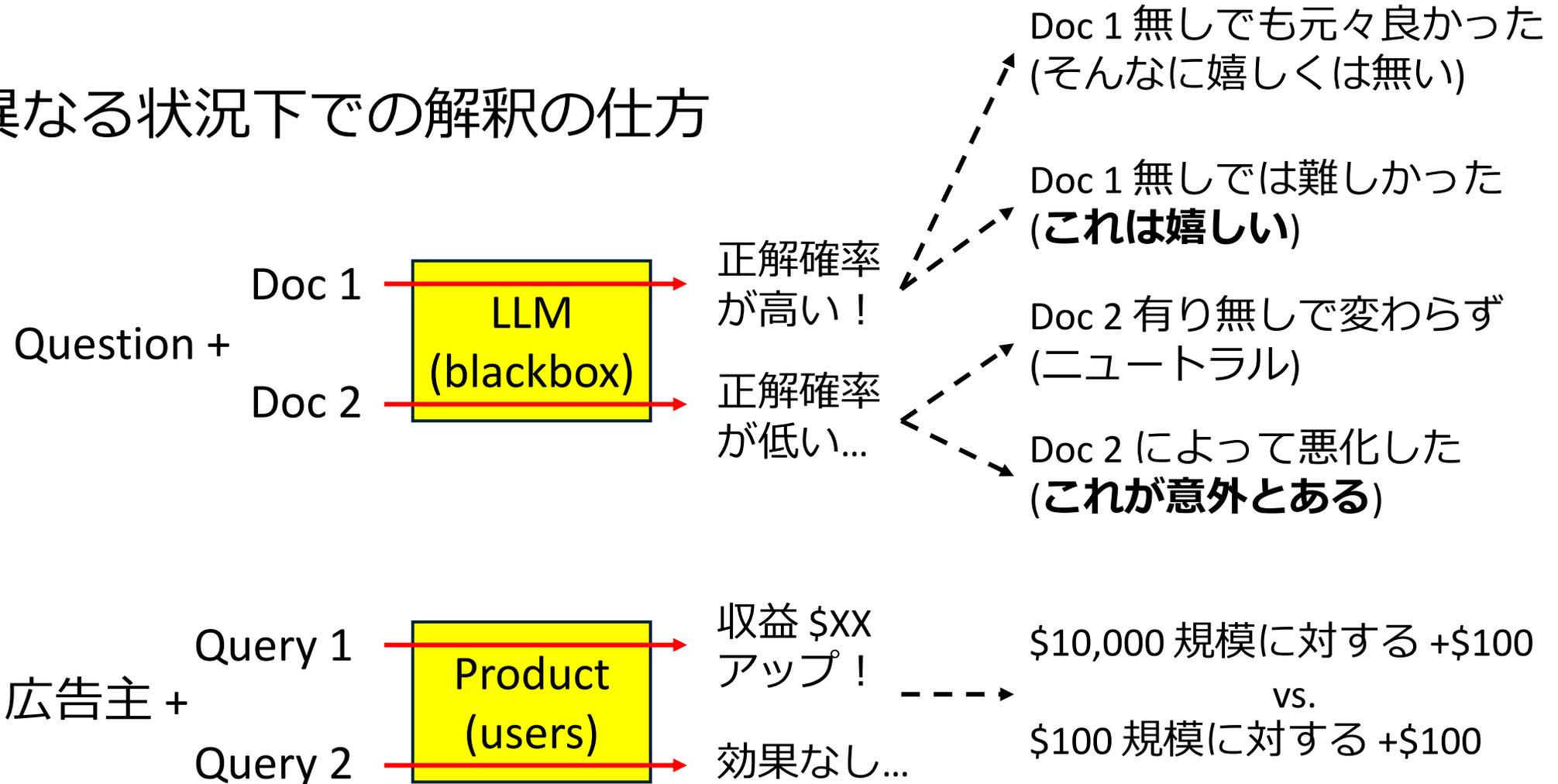
- 研究用のデータセットでは各 query に対して有用なデータのアノテーションが付与されることがある
- 例: question + **gold docs** -> retrieval/reranking モデルの学習

➤最終ゴールの達成度・報酬によるアノテーション (**今回の話題**)



ゴール達成度・報酬の解釈

▶異なる状況下での解釈の仕方

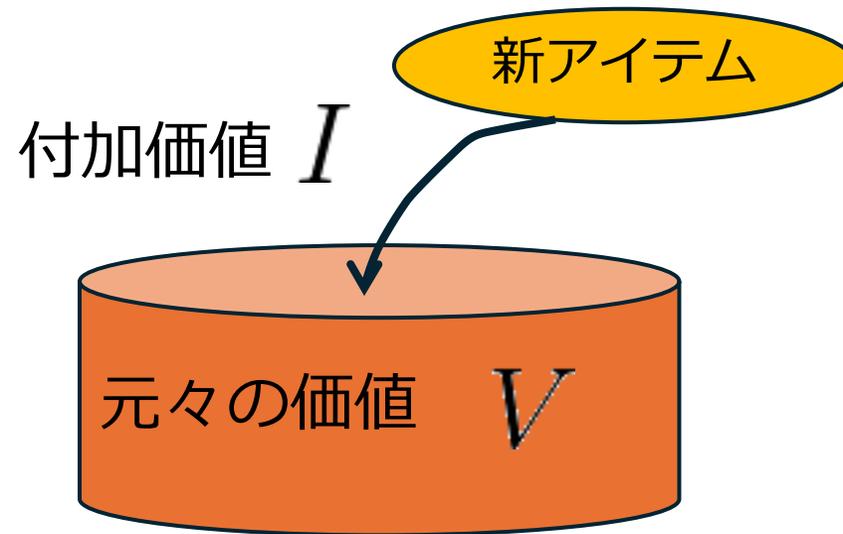


付加価値 (Incremental Utility)

- ゴールの達成度、スコア、報酬などの絶対値ではなく、新しいアイテムによってもたらされた**付加価値の割合**を計算

$$\frac{(V + I) - V}{V + I}$$

問題によっては
• スムージング
• リスケーリング
等を行う。



In-Context Learning での話

➤ ICL (in-context learning)

- LLM に**入出力の事例**を与えることで新しい入力に対処
- 手続き的には RAG に近い (有用な事例を見つけたい)
- 各事例の付加価値を計算したい

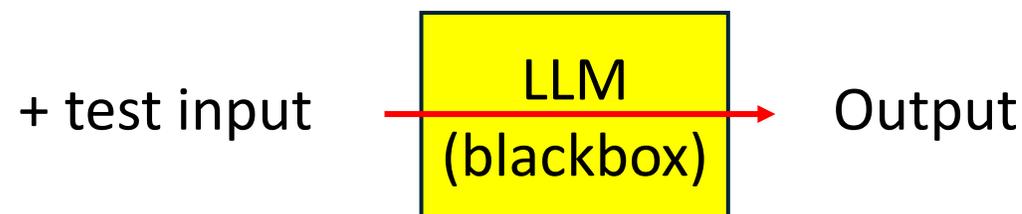
入出力の事例 (demonstrations)

Example input 1 -> output 1

Example input 2 -> output 2

...

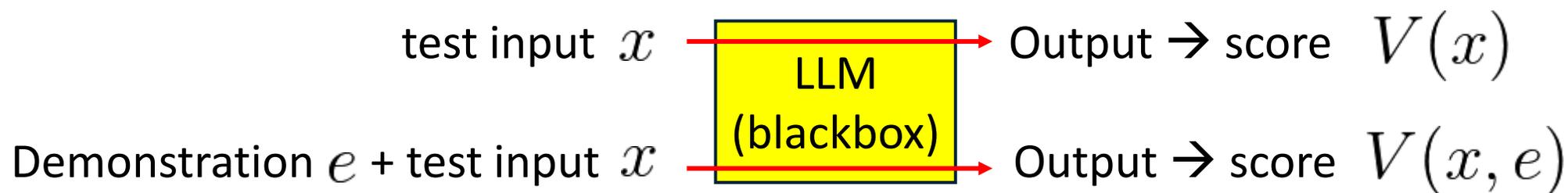
Example input K -> output K



どの事例の価値が高いのか？

入出力事例の付加価値

- 各入出力事例を足した場合の LLM の出力への影響を計算
 - 入出力事例間の**関係**が考慮されていないのでベストではない



$$\frac{V(x, e) - V(x)}{\max(V(x, e), V(x))^\ell}$$

値域は [-1.0, 1.0]

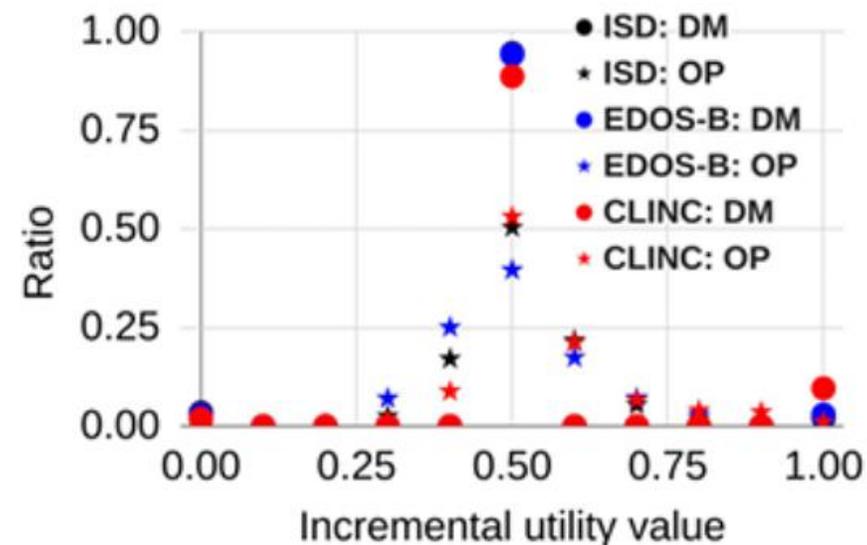
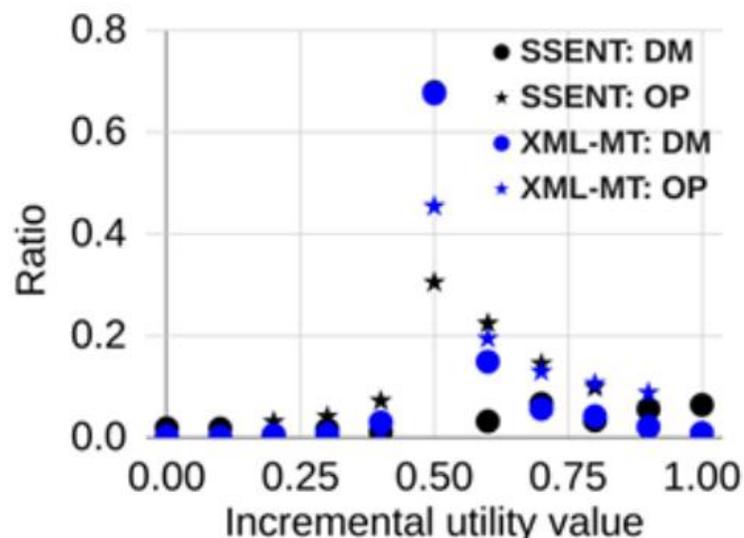
→
値域が [0.0, 1.0] になるように線形に変換

(分母の変形は、**値の対称性**とリスケーリングのため)

付加価値推定の実例

▶ 様々なデータセットで入出力の価値推定をプロット

- DM: V = タスク特有の評価指標 (F1, BLEU, etc.)
- OP: V = 正解例を出力する場合の確率



学習データでこの分析を行い、

その結果を元に**未知のデータ**についても価値の高い入出力事例を選択できるように学習 (NAACL 2024)

マルチタスク・多言語への対応

- 新しい問題に対して毎回この分析 + モデル学習は高コスト
- Demonstration 選択モデルの汎化性能アップの検討
 - Reranking model → dense retrieval model (スケーラビリティ)
 - 1 モデルを **80 以上のデータセット** (タスク) で学習
 - 新規タスクに対応する際のチャレンジ
 - マルチ・クロスリンガル性
 - 最新の Google Translate を利用して **230 の言語へ対応** すべく ICL のデータ拡張

様々なタスクへの対応の難しさ

➤ 80 の多様なデータセットを収集

- 要約、翻訳、NLI、・・・
- マルチリンガル

➤ Retrieval モデルを使う際の前提

- タスクの instruction が簡素に書ける
- 入出力の形式がはっきりしている

➤ 最近の LLM が扱う問題はもっと複雑

- きれいな形式のデータセットで学習したモデルが通用しない (現在の状況)

No.	Name	Type	Languages	Source	Scoring	[Q]	[C]
01	WMT14 en-es (Bojar et al., 2014)	Machine translation	en, es	Link	GLEU	100,000	30,099,732
02	WMT14 fr-es (Bojar et al., 2014)	Machine translation	en, fr	Link	GLEU	100,000	30,099,732
03	WMT16 en-es (Bojar et al., 2016)	Machine translation	de, en	Link	GLEU	60,000	4,143,251
04	WMT16 de-es (Bojar et al., 2016)	Machine translation	de, en	Link	GLEU	60,000	4,143,251
05	WMT16 en-es (Bojar et al., 2016)	Machine translation	en, es	Link	GLEU	30,000	2,296,592
06	WMT16 es-es (Bojar et al., 2016)	Machine translation	en, es	Link	GLEU	30,000	2,296,592
07	ANLI v1 (Nie et al., 2020)	Natural language inference	en [+MT]	Link	Probability	8,473	8,473
08	ANLI v2 (Nie et al., 2020)	Natural language inference	en	Link	Probability	22,730	22,730
09	ANLI v3 (Nie et al., 2020)	Natural language inference	en	Link	Probability	30,000	70,459
10	QNLI (Raghuvar et al., 2018)	Natural language inference	en	Link	Probability	30,000	74,543
11	MNLI (Williams et al., 2018)	Natural language inference	en	Link	Probability	30,000	100,000
12	WNLI (Levesque et al., 2012a)	Natural language inference	en	Link	Probability	317	318
13	MRPC (Dolan and Brockett, 2005)	Paraphrase identification	en	Link	Probability	200	3,268
14	PAWS (Zhang et al., 2019)	Paraphrase identification	en	Link	Probability	30,000	19,401
15	Tatoeba (Arntze and Schwab, 2019)	Translation identification	sq, fr, kur, tur	Link	Probability	30,000	177,554
16	MRPC (Dolan et al., 2011)	Sentiment classification	en	Link	Probability	12,467	12,467
17	SST2 (Socher et al., 2013)	Sentiment classification	en	Link	Probability	30,000	37,149
18	Yelp (Fast AI)	Sentiment classification	en	Link	Probability	30,000	100,000
19	Tweet Sentiment Extraction (Kaggle)	Sentiment classification	en [+MT]	Link	Probability	10,000	17,281
20	AfriSenti (Muhammad et al., 2021a)	Sentiment classification	amh, ha, or, ...	Link	Probability	30,000	33,685
21	TweetEval-emoji (Barbieri et al., 2018)	Emoji classification	en	Link	Probability	20,000	25,000
22	TweetEval-emotion (Muhammad et al., 2018)	Emotion classification	en	Link	Probability	1,600	1,657
23	Dialogintention (Kumar et al., 2021)	Multi-speaker emotion classification	en, hi	Link	F1	700	799
24	Massive-intent (Bergweiler et al., 2021)	Dialog intent classification	de, en, es, fr, ...	Link	Probability	30,000	100,000
25	MTOP-domain (Li et al., 2021)	Dialog domain classification	de, en, es, fr, ...	Link	Probability	30,000	43,928
26	MTOP-intent (Li et al., 2021)	Dialog intent classification	de, en, es, fr, ...	Link	Probability	30,000	43,928
27	ATIS-intent (Pitca, 1990)	Multi-label dialog intent classification	en	Link	F1	2,000	2,189
28	ELIZABETH-reversed (Dulek et al., 2019)	Semantic parsing (text to dict)	en	Link	F1	16,662	16,663
29	WikiSQL (Zhong et al., 2017)	Semantic parsing (textable to SQL)	en	Link	GLEU	20,000	36,353
30	BioSCTER (Li et al., 2016)	Named entity recognition (biomedical)	en	Link	F1	2,000	2,500
31	BioNLP1PC (Ohno et al., 2013)	Named entity recognition (biomedical)	en	Link	F1	1,000	1,499
32	JNLPBA (Huang et al., 2010)	Named entity recognition (biomedical)	en	Link	F1	9,000	9,346
33	MultiCoNER2 (Fetahu et al., 2023)	Named entity recognition	de, fr, fi, ...	Link	F1	30,000	140,824
34	CoNLL2003 (Tjong Kim Sang and De Meulder, 2003)	Named entity recognition	en	Link	F1	7,000	7,041
35	MTOP-slot (Li et al., 2021)	Dialog slot labeling	en, hi, hi	Link	F1	19,000	19,811
36	SNPES-slot (Crawford et al., 2018)	Dialog slot labeling	en	Link	F1	6,000	7,084
37	ATIS-slot (Pitca, 1990)	Dialog slot labeling	en	Link	F1	2,000	2,478
38	SemRel (Hendricks et al., 2010)	Relation classification (common)	en [+MT]	Link	Probability	5,500	4,000
39	DDI3 (Herrero-Zazo et al., 2013)	Relation classification (drugs)	en	Link	Probability	8,000	10,779
40	ChemProt (Islamaj Dogan et al., 2019)	Relation classification (chemical and protein)	en	Link	Probability	9,000	10,460
41	WordSeg (Hofman et al., 2020)	Word segmentation	en	Link	GLEU	30,000	100,000
42	PixPunct (Hofman et al., 2020)	Punctuation fix	en	Link	GLEU	30,000	100,000
43	CoLA (Warstadt et al., 2019)	Linguistic acceptability judgment	en	Link	Probability	4,173	4,176
44	CoNLL2000 (Tjong Kim Sang and Buchholz, 2000)	Named entity recognition	en	Link	F1	4,000	4,936
45	Prosson (Rahman and Ng, 2012)	Conference resolution	en	Link	Probability	561	561
46	WSC (Levesque et al., 2012b)	Conference resolution	en	Link	Probability	252	252
47	Winogrande (Sakaguchi et al., 2019)	Sentence completion	en	Link	Probability	20,999	20,999
48	WIC (Pitohar and Camacho-Collados, 2019)	Word sense disambiguation	en	Link	Probability	2,614	2,614
49	Pythos (Lu et al., 2021)	Code summarization	en	Link	GLEU	30,000	100,000
50	Java (Lu et al., 2021)	Code summarization	en	Link	GLEU	30,000	100,000
51	Go (Lu et al., 2021)	Code summarization	en	Link	GLEU	30,000	100,000
52	PHP (Lu et al., 2021)	Code summarization	en	Link	GLEU	30,000	100,000
53	Gigaword (Napoles et al., 2012)	Text summarization	en	Link	GLEU	30,000	100,000
54	SAMSum (Ghose et al., 2019)	Dialog summarization	en	Link	GLEU	7,366	7,366
55	Debate (Wang and Ling, 2016)	Debate summarization	en [+MT]	Link	GLEU	859	800
56	HateCheck (Klinger et al., 2022)	Hate speech detection/classification	en, hi, hi, ...	Link	Probability	20,053	20,053
57	Toxic (Mueenighoff et al., 2023)	Toxic text detection	en	Link	Probability	24,800	24,800
58	Contrast (O'Neill et al., 2021)	Contrastual review detection	de, en, ja	Link	Probability	7,500	7,718
59	Isoty (Van Hee et al., 2018)	Isoty detection	en	Link	Probability	1,400	1,462
60	Offensive (Zampieri et al., 2019)	Offensive text detection	en	Link	Probability	5,000	6,916
61	Savazon (Abu Farha et al., 2022)	Savazon detection	ar, es	Link	Probability	2,500	3,414
62	SQuAD2 (Raghuvar et al., 2018)	Reading comprehension	en	Link	GLEU	30,000	100,119
63	BioQ (Clark et al., 2019)	Reading comprehension	en [+MT]	Link	Probability	4,611	4,614
64	ERICP (Dua et al., 2019)	Reading comprehension (numerical)	en	Link	Probability	29,635	46,621
65	CopiedQA (Mishra et al., 2018)	Reading comprehension (common sense)	en	Link	Probability	1,475	2,678
66	CoSmos (Huang et al., 2019)	Reading comprehension (common sense)	en	Link	Probability	12,531	12,531
67	SciDocs (Cohan et al., 2020)	Relevance, no-ranking	en	Link	Probability	30,000	99,159
68	HotpotQA (Yang et al., 2018)	Relevance, no-ranking	en	Link	F1	30,000	60,447
69	A12 ABC-easy (Clark et al., 2018)	Close-book question answering	en	Link	Probability	1,025	1,026
70	A12 ABC-challenge (Clark et al., 2018)	Close-book question answering	en	Link	Probability	459	460
71	TriviaQA (Joshi et al., 2017)	Open-book question answering	en	Link	Probability	108,184	100,000
72	Math (Sutton et al., 2019)	Math question answering	en	Link	Probability	30,000	30,000
73	CommonGen (Lin et al., 2020)	Constrained text generation (common sense)	en	Link	GLEU	30,000	37,161
74	SNLI-en (Bowman et al., 2015)	Constrained text generation (restaurant)	en	Link	GLEU	10,112	33,106
75	PRQA-agen (Bisk et al., 2019)	Question/query generation	en [+MT]	Link	GLEU	7,956	7,957
76	arXiv (Mueenighoff et al., 2023)	Multi-label topic/category classification	en	Link	F1	30,000	69,113
77	medRxiv (Mueenighoff et al., 2023)	Topic/category classification	en	Link	Probability	5,000	16,229
78	EBpedia (Lehmann et al., 2014)	Topic/category classification	en [+MT]	Link	Probability	5,000	5,000
79	Yahoo (Zhang et al., 2015)	Topic/category classification	en	Link	Probability	14,575	14,575
80	AG news (Zhang et al., 2015)	Topic/category classification	en	Link	Probability	30,000	89,800
81	TREC (Li and Roth, 2002)	Topic/category classification	en [+MT]	Link	Probability	2,626	2,626

以上を基にした最近の興味

- 従来の事前学習モデル (mT5 等) の retriever に限界を感じる
 - 「要約」「翻訳」などとしたざっくりしたタスク instruction ではなく、「〇〇に特に気を付けた要約/翻訳」といった **ニュアンスを加味したうえで価値のある事例**を選択可能か？
 - より最近の product レベルの LLM を用いることで可能かどうか
- これまでの研究の LLM 学習自体への還元
 - LLM を blackbox として feedback を得るために使うのではなく、それ自体に価値推定などの能力を明示的に学習させる

後半の話題

- LLM 研究の中心により近づくことで研究の仕事が変化
 - Google Research 時代:
 - 研究 <--> Product の繋がりを意識しつつも Publication は推奨されていた

- 現在は研究と競争意識の高い Product (LLM) との距離が急接近
 - 何を考えながら過ごしているか？

Product としての LLM の普及

- LLM に対する様々な使い方と日々の評価
 - 研究・開発に限らない人々が**あらゆる方法**で触る (評価する)
 - テレビ等でも当たり前のように取り上げられている

究極のマルチタスク学習

- 様々な**ユーザーの用途に対応**する必要性
 - 究極のマルチタスク学習
- 従来のマルチタスク学習研究との違い
 - 想定するタスクの数: $O(1) \sim O(10)$ → 無数のユーザー入力
 - 「タスク」という概念自体の曖昧性 (何がどんな形式で聞かれるのか、どう出力すべきか、の想定が難しい)

研究としての難しさ (おもしろさ)

- 一般的な研究サイクル
 - 問題提起 → 解決法の提案 → 改善 (→ 論文)
- 追加で考えること (コアの **LLM 学習** 自体への貢献として)
 - その他・既存の学習データとの共存が大事である
 - 特定の問題を解決した結果、他に悪影響がでるかもしれない

研究者としての身の振り方

- Publication の優先度の (一時的な) 低下
 - 内部インパクト (LLM 貢献度) が大きい = 秘密事項が多い
 - おもしろいコアの部分に近づけば近づくほど、对外発表をひとまず度外視する傾向になる
- どうやってうまくバランスをとるか、が (個人的な) 課題
 - こういった業界に興味を持たれている方にはここが一番重要なポイントかもしれません