Transformer LLM における層単位の FFN 層の重要度検証

矢野一樹¹ 高橋良允¹ 李宰成¹ 柴田圭悟¹ 鈴木潤^{1,2,3} 池田航1 1 東北大学 2 理化学研究所 3 国立情報学研究所 ikeda.wataru.q5@dc.tohoku.ac.jp

概要

Transformer に基づく大規模言語モデル (LLM)の 構成要素の一つであるフィードフォワードネット ワーク (FFN) に着目し、モデル内の配置場所に依存 した重要度を検証する.具体的には、モデル全体の パラメータ数を維持したまま,一部の連続する層で FFN の中間次元を拡大し、残りの層から FFN を除去 したモデル構成を用いて標準的なタスク性能を比較 する. 複数のモデルサイズで評価を行った結果, 全 層数の 70 – 90%の連続した中間から後方の層に FFN を集中配置することで、複数の下流タスクでベース ラインの性能を上回った.この結果から FNN は入 力に近い層より中間から後方の層で特に重要である と示唆される結果が得られた.

はじめに 1

Transformer [1] に基づく大規模言語 モデル(Large Language Model: LLM) では、各モデルによって詳細 設計は異なるが、自己注意機構とフィードフォー ワードネットワーク (Feed-Forward Network: FFN) の二つを主な構成要素として一つの Transformer 層を構成しているモデルが多い [1, 2]. 図 1(a) に Transformer 層を図示する¹⁾. 特に Pre-LN [3] と呼ば れる Transformer のモデル構成の場合, これら2つを 1層として計算された結果を層の数だけ入力層であ るトークンの埋め込みベクトルに加算し、最終的な 隠れ状態ベクトルを算出する計算手順になる.

一般論として、自己注意機構は、主にトークン の埋め込みから得られた情報の混ぜ合わせを担っ ており、FFN は、学習データ内にある知識などを 記憶することが主な役割と説明されることが多 い [4, 5, 6]²⁾. 知識が FFN 層に埋め込まれていると

いう仮説が正しいとした場合, FFN が知識を獲得す る上で最もよい形式なのか,実際に Transformer 内 の複数ある層の中でどのあたりに知識が埋め込まれ るのかなど未解明な事象も多い. そこで、本研究で は FFN に着目し, FFN を削除したモデルや FFN を 大きくした設定など、いくつかの通常とは違う設定 で LLM の事前学習を実施することで、FFN の役割 や機能の一部を検証する.

2 関連研究

Transformer に基づく LLM における FFN の役割や 機能を検証している論文は多く存在している。文 献 [4, 5, 6] では, FFN が知識の記憶装置として機能 することを示し、特定のニューロンが事実的知識の 表現や想起に重要な役割を果たしていることを明ら かにした.一方, 文献 [7] では, 分析を通じて, FFN が層正規化と合わせて入力の文脈化に寄与するとい う FFN の役割に関する新たな解釈を与えた.

このように、幾つもの重要な知見が得られている が、これらの知見は、事前学習済み LLM に対する 分析結果から得られたものであり,標準的なモデ ル構成に限定される.本稿では, Transformer 各層の FFN の一部を削除したり、次元数を増やしたりする など、モデル構成そのものを変化させた際の影響を 検証するという点で、従来研究とは違う検証方法と なっている.また、従来と異なる検証方法を採用す ることで、FFN の機能や役割に関する新たな知見を 得ることを試みている.

モデル設計と層の構成 3

3.1 ベースラインモデル

本研究ではベースラインモデルとして, LLaMA [2] で提案されたモデル設計を用いて実験を行う. LLaMA では, FFN には SwiGLU 活性化関数 [8] が採 用されており、FFN は入力ベクトル $x \in \mathbb{R}^d$ を受け 取り、内部で中間表現の次元数 df まで拡張して処

¹⁾ 本稿では、自己注意機構と FFN をまとめた Transformer 層 を略して「層」と表記する.また,層正規化は本研究では重 要な要素ではないので, 簡略化のため説明から除外する.

²⁾ 現在は、それ以外の様々な機能や効果があることが検証に より示されている [7].



図1 ベースラインモデルと検証モデルの異なる層の構成 通常の LLaMA の層の積み重ねであるベースラインモデルに 対し (上段左),検証モデルでは一部の層の FFN の中間表現の次元を拡張し ((b)FFN 拡張層),残りの層から FFN を除去す る ((c)FFN 不活性層) ことにより (上段右),ベースラインモデルに対して全体のパラメータ数は維持しつつ FFN の計算能 力 (=パラメータ数) を特定の層に集中させたモデルを実現する.

理を行う ($W_{\text{gate}}, W_{\text{up}} \in \mathbb{R}^{d_{\text{f}} \times d}, W_{\text{down}} \in \mathbb{R}^{d \times d_{\text{f}}}$):

$$FFN(x) = W_{down}(Swish(W_{gate}x) \otimes W_{up}x)$$
(1)

$$Swish(x) = x\sigma(x)$$
 (2)

3.2 検証モデル

ベースラインモデルで採用した標準的な LLaMA モデルの層を FFN 拡張層,または,FFN 不活性 層で置き換えたモデルを検証モデルとする.

FFN 拡張層は,標準的な自己注意機構に加え,拡 張された中間次元数を持つ FFN_{expanded} を配置する (図 1 (b) 参照).ここで,FFN_{expanded} は以下のように 定義する:

$$FFN_{expanded}(x) = W'_{down}(Swish(W'_{gate}x) \times W'_{up}x)$$
(3)

ただし、 W'_{gate} , $W'_{up} \in \mathbb{R}^{d'_{f} \times d}$, $W'_{down} \in \mathbb{R}^{d \times d'_{f}}$ であり、 中間次元数 d'_{f} はベースラインの層の中間次元数 d_{f} より大きい値をとる. この FFN 拡張層の中間次元 数 d'_{f} は、ベースラインのすべての層を FFN 拡張層 か後述する FFN 不活性層に置き換えてもモデル全 体のパラメータ数がほぼ一致するように決定する.

FFN 不活性層は,前述の標準的な LLaMA モデル の層から FFN を除去し,自己注意層のみとした層で ある(図1(c)参照).

本研究では、ベースラインモデルの層を、特定の 配置で、これら FFN 拡張層または FFN 不活性層に て置き換えた検証モデルを設定し、FNN の位置に依 存した重要度の検証実験に用いる.

4 実験

本研究では,前節で説明したベースラインモデル と検証モデルを事前学習した後に標準的なタスク性



図2 FFN 拡張層の異なる配置位置 検証モデルの FFN 拡 張層を入力層に近い位置(first),中間層(middle),出力 層に近い位置(final)のいずれかに配置し,各位置におけ る効果を検証した.

能を評価し,どの位置にある FFN を除去しても性 能が劣化しないか,或いは,逆に性能が向上するか などの挙動を調査し,その結果から位置に依存した FFN の重要度を検証する.

4.1 ベースラインモデルの設定

本研究では、285M および 570M パラメータの 2 つの異なるサイズのベースラインモデルを用いる. 285M パラメータモデルは 12 層で構成され、隠れ層 の次元数 *d* は 1280 とする.570M パラメータモデル では、285M と同じ隠れ層の次元数を維持しながら、 層数を 24 層に拡張する.また、FFN における中間 次元数 *d*_f は各モデルサイズのベースラインモデル 共通であり、全ての層で 4480 とする.

4.2 検証モデルの設定

検証モデルにおいても,層数や隠れ層の次元数を 含めて基本的なモデル設定はベースラインと共通 とする.検証モデルでは,以下の2つの要素に従っ てベースラインモデルの層をFFN 拡張層,または, FFN 不活性層に置き換える.



図3 FFN 拡張層の割合に対するタスク性能の推移上段 (a)~(d),下段 (e)~(g) はそれぞれモデルサイズ 285M,570M に での各タスクの評価結果.横軸:FFN 拡張層の割合,縦軸:ベースラインに対する相対的な改善度,異なる色のプロッ ト:FFN 拡張層の異なる配置位置 (first, middle, final),水平赤点線:相対的改善度が0 (ベースラインと同等の性能).

- 全層数に対する FFN 拡張層の割合 r%. ただし、 r ∈ {10, 30, 50, 70, 90, 100} を用いる. FFN 拡張層 の総数は層数 L との積の小数点以下を切り捨て た値([rL/100])となる.
- FFN 拡張層を配置する位置 {fisrt, middle, final}.図2に示す通り, FFN 拡張層の配置位 置として first (第1層から後続する層に配置), middle (中間層 (L/2 層目)を起点に対称になる ように配置), final (最終層から先行する層に 配置)を設定した³⁾.

以上の2つの要素の組み合わせにより,計15種 類(FFN 拡張層の割合5通り×配置位置3通り)の 配置パターンの検証モデルを設定する.

4.3 事前学習と評価

ベースラインモデルと検証モデルの事前学習に は標準的な事前学習法を用いる⁴⁾.また,事前学習 済みモデルの評価については下流タスク性能およ びモデルの知識量の観点から評価を行う⁵⁾.モデ ルの知識量の評価に関しては,2節で示した,FFN が知識を記憶するという先行研究の知見に基づき, Zero-Shot Relation Extraction (zsRE) データセット [9] を用いた知識量測定タスク (zsRE タスク) で評価す る [10, 11]⁶⁾.

各タスクの評価結果について、検証モデルとベー スラインとの比較を簡単にするために以下で定 義される、ベースラインに対する相対的な改善度 (Relative Improvement:RI)を算出する:

$$RI(m,T) = \frac{score(m,T) - score(baseline,T)}{score(baseline,T)} \times 100[\%]$$
(4)

ここで, score(*m*,*T*) はモデル*m*のタスク*T*における 評価スコア, score(baseline,*T*) はベースラインモデ ルの評価スコアを表す. Perplexity など値が小さい ほど性能が良いとされる評価指標については,符号 を反転させている. 相対的改善度が0の場合はベー スラインと同等の性能,正の値はベースラインより 性能が高く,負の値は性能が低いことを示す.

5 実験結果および考察

異なる配置パターン(FFN 拡張層の割合 × 配置 位置⁷⁾)の検証モデルに対する各タスクの評価結果

³⁾ なお本研究では,連続する層に FFN 拡張層を配置する設 定に限定して検証を行う.実際には飛び飛びの層に FFN 拡 張層を配置するなど配置パターンは無数にあるが,これらの 検証は今後の課題とする.

⁴⁾ 事前学習の詳細設定については付録 A に記載

⁵⁾ 具体的な評価タスクは付録 B に記載

⁶⁾ zsRE タスクの詳細は B.2 に記載

⁷⁾ FFN 拡張層の割合が 100%のモデルはベースラインモデル と一致する.図3中のこのモデルのプロットは FFN 拡張層 の配置位置に関わらず全て赤点線上(ベースラインと同等の)

表1 各モデルサイズにおける上位5モデルの平均相対 的改善度(RI) 平均 RI は4つの評価タスク(Wikitext, LAMBADA, HellaSwag, zsRE)の平均値を示す.

| 285M | モデル | 570M モデル | | |
|-----------|-----------|-----------|-----------|--|
| モデル | 平均 RI (%) | モデル | 平均 RI (%) | |
| final_90 | +3.37 | middle_70 | +1.35 | |
| middle_90 | +2.57 | final_90 | +0.85 | |
| middle_70 | +2.14 | middle_50 | +0.03 | |
| final_70 | +1.03 | middle_90 | -0.00 | |
| first_90 | +0.54 | final_70 | -0.52 | |

表 2 モデルサイズ 2B における検証モデルの相対的改善 度の平均 全ての評価タスクの平均を算出.評価の詳細は C を参照.

| | 平均 RI (%) | | | |
|-----------|-----------|--|--|--|
| middle_70 | -0.71 | | | |
| final_90 | +1.06 | | | |

を図 3 に示す⁸⁾.

5.1 FFN 拡張層の割合の影響

FFN 拡張層の割合が性能に与える影響に観察する と,FFN 拡張層の割合が低い(10 – 30%)モデルで は,図3(c)のfinal_30_285M⁹⁾の条件を除くすべて の下流タスクおよび知識量に関するタスクにおい て,ベースラインを大きく下回る性能を示した.つ まり,FFN を一部の層に極端に集中させる設定はモ デルの性能を著しく低下させる.

一方, FFN 拡張層の割合を増やすと性能は徐々に 改善し, FFN 拡張層の割合が 70 – 90% の高い範囲 では, すべてのタスクにおいて, 少なくとも1つの 検証モデルでベースラインを上回る性能が達成さ れた.興味深い点として, ベースラインのように 全ての層に均等に FFN を配置するよりも, 一部の 少数の層から FNN を除去し,残りの層(全層数の 70 – 90%)に緩やかに集中させるとタスク性能が顕 著に向上する可能性があることが示唆された.

5.2 FFN 配置位置の影響

middle 設定は,図3(e)のmiddle_90_570Mを除き, FFN 拡張層の割合が70-90%の範囲ですべての評価タスクにおいてベースラインを上回る性能を示し た.一方,firstやfinalでは,同じFFN 拡張層の 割合でもより多くのタスクでベースラインを下回っ た.また,特に下流タスク性能において(Wikitext, LAMBADA, HellaSwag),firstのモデルが middle や final と比較して大きく性能が下回る傾向が顕著 に見られた(図 3(a),(b),(c),(e),(f),(g)).これ らの観測から,中間部分から後半部分にあるFFN が より効果的に機能していることが示唆され,逆に前 半部分は効果が薄い可能性が示唆された.

5.3 各検証モデルの個別の評価

表1に各モデルサイズにおける平均の相対的改善 度上位5モデル構成を示す.final_90(285Mで1 位(+3.37%),570Mで2位(+0.85%))とmiddle_70 (285Mで3位(+2.14%),570Mで1位(+1.35%))が 全体の中で一貫して良い性能を示した.

5.4 より大きなモデルサイズでの再現性

285M および 570M の両方で安定して高い性能を 示した前述の final_90 と middle_70 のモデル構成 について, 2B パラメータモデルの設定でも追加の検 証を行った.表2 にモデルサイズ 2B における評価 結果を示す¹⁰⁾.表2 より明らかなように final_90 がすべての評価タスクを平均してベースラインモデ ルを上回っており,モデルサイズを拡張しても表1 で示した結果と概ね一貫していることが観測され た.これは,モデルサイズが変わっても傾向が大き くずれない可能性が高いことを示唆しており,良い 性質と言える.

6 おわりに

Transformer に基づく大規模言語モデル (LLM)の 構成要素の一つである FFN に対して,モデル全体で の配置場所に依存した重要度を検証した.複数のモ デルサイズで評価を行い,全層数の 70 – 90%の連続 した中間から後方の層に FFN を集中配置すること で,複数の下流タスクでベースラインの性能を上回 る結果を得た.この結果から FNN は入力に近い層 より中間から後方の層で,より効果を発揮する可能 性が高いことが示唆された.FFN の配置を工夫する ことで,手軽に標準的な Transformer モデルを改善 できるのであれば,モデル構成のチューニングなど 新たな進展や可能性が期待できる.

性能)であることを確認されたい.

⁸⁾ 一部の評価データではベースラインモデルがチャンスレートを下回ったため本稿の議論からは除外した.

⁹⁾ 以降,個別の検証モデルを"[FFN 拡張層の配置位置]_[FFN 拡張層の配置割合]_[モデルサイズ]"の形式で表す。例として、FFN 拡張層の配置位置が middle,配置割合が 50%、モデ ルサイズが 570M の場合は "middle_50_570M" と表す。

^{10) 2}B モデルの評価の詳細については付録 C に記載.

謝辞

本研究は、JST ムーンショット型研究開発事業 JPMJMS2011-35 (fundamental research),および、文部 科学省の補助事業「生成 AI モデルの透明性・信頼 性の確保に向けた研究開発拠点形成」の支援を受け たものです.

本研究は九州大学情報基盤研究開発センター研究 用計算機システムの一般利用を利用しました.

本研究成果(の一部)は、データ活用社会創成 プラットフォーム mdx [12] を利用して得られた物 です.

参考文献

- A Vaswani. Attention is all you need. Advances in Neural Information Processing Systems, 2017. [1]
- Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, [2] and Guillaume Lample. Llama: Open and efficient foundation language models. CoRR, Vol. abs/2302.13971, , 2023.
- models. CoRR, Vol. abs/2302.13971, , 2023. Ruibin Xiong, Yunchang Yang, Di He, Kai Zheng, Shuxin Zheng, Chen Xing, Huishuai Zhang, Yanyan Lan, Liwei Wang, and Tie-Yan Liu. On layer normalization in the transformer architecture. In Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event, Vol. 119 of Proceedings of Machine Learning Research, pp. 10524–10533. PMLR, 2020. Mor Geva, Roei Schuster, Jonathan Berant, and Omer Levy. Transformer feed-forward layers are key-value memories. In Marie-Eranzine Moene, Xu.
- [4] feed-forward layers are key-value memories. In Marie-Francine Moens, Xureceiver variables are key-value memories. In Marie-Francine Moens, Au-anjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Domini-can Republic, 7-11 November, 2021, pp. 5484–5495. Association for Computational Linguistics, 2021.
- [5]
- Computational Linguistics, 2021. Damai Dai, Li Dong, Yaru Hao, Zhifang Sui, Baobao Chang, and Furu Wei. Knowledge neurons in pretrained transformers. In Smaranda Muresan, Preslav Nakov, and Aline Villavicencio, editors, **Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), ACL 2022, Dublin, Ireland, May 22-27, 2022**, pp. 8493–8502. Association for Computational Linguistics, 2022. Kevin Meng, David Bau, Alex Andonian, and Yonatan Belinkov. Locating and editing factual associations in GPT. In Sanmi Koyejo, S. Mohamed, A. Agarwal, Danielle Belgrave, K. Cho, and A. Oh, editors, Advances in Neural Information Processing Systems 35: Annual Conference on Neural Information Processing Systems 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022, 2022. Goro Kobayashi, Tatsuki Kuribayashi, Sho Yokoi, and Kentaro Inui. Analyz-ing feed-forward blocks in transformers through the lens of attention maps. In **The Twelfth International Conference on Learning Representations**, [6]
- [7] The Twelfth International Conference on Learning Representations, ICLR 2024, Vienna, Austria, May 7-11, 2024. OpenReview.net, 2024. Noam Shazeer. GLU variants improve transformer. CoRR, Vol. [8]
- Noam Shazeer. GLU abs/2002.05202, , 2020. [9]
- Omer Levy, Minjoon Seo, Eunsol Choi, and Luke Zettlemoyer. Zero-shot relation extraction via reading comprehension. **CoRR**, Vol. abs/1706.04115, 2017
- Eric Mitchell, Charles Lin, Antoine Bosselut, Chelsea Finn, and Christo-pher D. Manning. Fast model editing at scale. In The Tenth International Conference on Learning Representations, ICLR 2022, Virtual Event, [10] April 25-29, 2022. OpenReview.net, 2022. Nicola De Cao, Wilker Aziz, and Ivan Titov. Editing factual knowledge in
- [11] Hotora De Cao, Wilker AZZ, and Yan Hov. Editing factual knowledge in language models. In Marie-Francine Moens, Xuanjing Huang, Lucia Specia, and Scott Wen-tau Yih, editors, Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021, pp. 6491–6506. Association for Computational Linguistics, 2021. Toyotaro Suzumura, Akiyoshi Sugiki, Hiroyuki Takizawa, Akira Imakura,
- [12] Hiroshi Nakamura, Akriyoshi Sugiki, Hiroyuki Takizawa, Akifa Imawua, Hiroshi Nakamura, Kenjiro Taura, Tomohiro Kudoh, Toshihiro Hanawa, Yuji Sekiya, Hiroki Kobayashi, Yohei Kuga, Ryo Nakamura, Renhe Jiang, Junya Kawase, Masatoshi Hanai, Hiroshi Miyazaki, Tsutomu Ishizaki, Daisuke Shimotoku, Daisuke Miyamoto, Kento Aida, Atsuko Takefusa, Takashi Ku-rimoto, Koji Sasayama, Naoya Kitagawa, Ikki Fujiwara, Yusuke Tanimura, Takayuki Aoki, Toshio Endo, Satoshi Ohshima, Keiichiro Fukazawa, Susumu Data, and Tashibiro Itohibayashi mdy. A cloud platform for supmorting data Date, and Toshihiro Uchibayashi. mdx: A cloud platform for supporting data science and cross-disciplinary research collaborations. In 2022 IEEE Intl Conf on Dependable, Autonomic and Secure Computing, Intl Conf on Pervasive Intelligence and Computing, Intl Conf on Cloud and Big Data Computing, Intl Conf on Cyber Science and Technology Congress (DASC/PiCom/CBDCom/CyberSciTech), pp. 1–7, 2022.
- Anton Lozhkov, Loubna Ben Allal, Leandro von Werra, and Thomas Wolf. Fineweb-edu: the finest collection of educational content, 2024. [13]

- [14] Jordan Hoffmann, Sebastian Borgeaud, Arthur Mensch, Elena Buchatskaya, Trevor Cai, Eliza Rutherford, Diego de Las Casas, Lisa Anne Hendricks, Johannes Welbl, Aidan Clark, Tom Hennigan, Eric Noland, Katherine Millican, George van den Driessche, Bogdan Damoc, Aurelia Guy, Simon Osindero, Karen Simonyan, Erich Elsen, Oriol Vinyals, Jack W. Rae, and Laurent Sifre. An empirical analysis of compute-optimal large language model training. In An eniprical analysis of compare opiniar angle angle angle in the opiniar angle in the second [15]
- 2022, NeurIPS 2022, New Orleans, LA, USA, November 28 December 9, 2022, 2022. Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019. Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. Ope-nAl blog, Vol. 1, No. 8, p. 9, 2019. Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi Charles Foster Laurence Golding. Leffrey Hu, Alain Le Noac'h
- [16]
- [17] Leo Gao, Jonathan Tow, Baber Abbasi, Stella Biderman, Sid Black, Anthony DiPofi, Charles Foster, Laurence Golding, Jeffrey Hsu, Alain Le Noac'h, Haonan Li, Kyle McDonell, Niklas Muennighoff, Chris Ociepa, Jason Phang, Laria Reynolds, Hailey Schoelkopf, Aviya Skowron, Lintang Sutawika, Eric Tang, Anish Thite, Ben Wang, Kevin Wang, and Andy Zou. A framework for few-shot language model evaluation, 07 2024. Denis Paperno, Germán Kruszewski, Angeliki Lazaridou, Quan Ngoc Pham, Raffaella Bernardi, Sandro Pezzelle, Marco Baroni, Gemma Boleda, and Raquel Fernández. The lambada dataset, Aug 2016. Stenhen Merity Caiming Xiong James Bradbury and Richard Socher
- [18]
- Stephen Merity, Caiming Xiong, James Bradbury, and Richard Socher. Pointer sentinel mixture models. In 5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings. OpenReview.net, 2017. Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winnermote, An advergarial unpaged acheme abellance at earle. In The [19]
- [20] Winogrande: An adversarial winograd schema challenge at scale. In The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial In-telligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020, pp. 8732–8740. AAAI Press, 2020 2020
- [21] Yonatan Bisk, Rowan Zellers, Ronan Le Bras, Jianfeng Gao, and Yejin Choi. Piqa: Reasoning about physical commonsense in natural language. In Thirty-
- Figure Addition of the State of [22] Linguistics, 2019.
- Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind Tafjord. Think you have solved question an-swering? try arc, the Al2 reasoning challenge. **CoRR**, Vol. abs/1803.05457, [23] 2018.
- Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity? a new dataset for open book question answering. In **EMNLP**, 2018. [24]

表3 モデルサイズ 2B における検証モデルのベースラインモデルに対する相対的改善度(%)

| モデル | Wikitext | LAMBADA | ARC-e | ARC-c | Winogrande | PIQA | OBQA | HellaSwag | 平均 |
|-----------|----------|---------|-------|-------|------------|-------|-------|-----------|-------|
| middle_70 | +1.60 | -9.01 | 0.00 | +6.32 | -1.61 | +1.11 | -4.35 | +0.24 | -0.71 |
| final_90 | +1.07 | -7.88 | -0.56 | +5.17 | +2.68 | +0.55 | +6.52 | +0.94 | +1.06 |

A 事前学習の詳細設定

事前学習には、FineWeb-Edu [13] データセットを使用する. 学習トークン数は、Chinchilla の法則 [14] に従い、285M、570M、2B パラメータモデルに対してそれぞれ、5.7B、11.4B、40B トークンとする. バッチサイズと最大シーケンス長 (1024 トークン) から算出される 1 ステップあたりのトークン数をもとに、総ステップ数を 20,000 ステップと設定する. 事前学習の設定として、オプティマイザには AdamW [15] を使用し、パラメータは $\beta_1 = 0.9$ 、 $\beta_2 = 0.95$ とする. 重み減衰 係数は 0.1 を採用する. 学習率のスケジューリングにはコサインスケジューラを採用し、最初の 1000 ステップで線形に 増加させて最大学習率 3e-4 に到達させた後、コサイン波形に従って減衰させる. トークナイズには GPT-2 のボキャブラ リを使用する [16].

B 評価手法の詳細

B.1 下流タスク性能の評価

下流タスク性能については lm-evaluation-harness フレームワーク [17] を利用し,複数タスクから多角的な評価を実施する.評価には以下のタスクを使用する:

LAMBADA[18] は文章全体の文脈を理解し最後の単語を予測するタスク,Wikitext[19] はWikipedia 記事を用いた言語モ デリングタスクである。Winogrande[20] は常識的推論に基づく2択形式のタスク,PIQA[21] は日常生活における物理的 な常識の理解を評価する2択形式のタスクである。HellaSwag[22] は文脈から最も自然な文を選択する4択形式のタスク である。ARC[23] は科学的知識と推論能力を評価する4択形式のタスクで,比較的単純なEasy Set (ARC-e) と深い推論が 必要な Challenge Set (ARC-c) で構成される。OpenBookQA(OBQA)[24] は科学的知識を用いた応用的な推論能力を評価 する4択形式のタスクである。

評価指標として、2 択形式の Winogrande と PIQA (チャンスレート 50%)、4 択形式の HellaSwag, ARC, OBQA (チャ ンスレート 25%) では正解率 (Accuracy) を用いる. LAMBADA と Wikitext では正解率に加えて、語彙サイズに依存する Perplexity も測定する.

B.2 モデルの知識量の評価

Zero-Shot Relation Extraction (zsRE) データセット [9] の各事例は、知識に基づく質問文とその答えのペアで構成される. 知識量を測定するタスク (zsRE タスク)[10, 11] では、評価時に、質問文のみ、または質問文と答えの一部を入力として与 え、モデルに次のトークンを生成させる. 具体的には、まず質問文のみのプロンプトから1トークン生成させ、答えの最 初のトークンとの一致を調べる.次に答えの最初のトークンを最初のプロンプトに付加したものを2番目プロンプトと し、1トークン生成させ、答えの2番目のトークンとの一致を調べる、という手順を繰り返す. 答えの全トークンに対す る一致率をその事例に対する正解率とし、全事例(約20,000件)の正解率の平均をモデルの知識量の指標とする.

C 2B モデルの評価詳細

表 3 に 2B モデルにおける評価タスク毎のベースラインに対する相対的な改善度(Relative Improvement:RI)を示す. middle_70と final_90 の評価結果から、いくつかの興味深い知見が得られた. Wikitext や PIQA, HellaSwag では両モデ ルともわずかな性能向上(+0.24%~+1.60%)を示した一方で、文脈全体の理解を必要とする LAMBADA では大幅な性能 低下(-9.01%, -7.88%)が観察された. 特筆すべき点として、科学的知識と複雑な推論を必要とする ARC-c では両モデル とも顕著な性能向上(+6.32%, +5.17%)を示し、OBQA では final_90 で+6.52%の大幅な改善がみられた. 総合的に見る と、final_90 は Winogrande(+2.68%)を含む多くのタスクで性能向上を示し、平均でも+1.06%の改善を達成した一方、 middle_70 は一部のタスクで性能低下が見られ平均で-0.71%となった.