

深層学習モデルを用いた双方向形態屈折の検証

深津聡世* 原田宥都* 大関洋平

東京大学大学院 総合文化研究科言語情報科学専攻

{akiyofukatsu, harada-yuto, oseki}@g.ecc.u-tokyo.ac.jp

概要

人間が形態処理を行うにあたって必要な知識は、規則か類推か、あるいはその両方かという議論は、言語学の形態論において現在も続いている。この議論は「過去時制論争 (Past Tense Debate)」と呼ばれ、近年ではニューラルネットワークを用いた形態処理のモデル化による検証が行われている。本研究では、その形態的な複雑さから過去時制論争において重要とされる日本語動詞の屈折について、類推のモデルである深層学習モデルを用いて、時制について双方向の形態屈折の学習を行った。どちらの時制方向の処理にモデルがより適するかを検証した結果、訓練データのサイズや性質が結果に大きく影響することが示唆された。

1 はじめに

人間が形態処理を行う際に用いられているのは、規則的な処理なのか、類推的な処理なのか、あるいはその両方を用いた処理なのかということについて、言語学の形態論において長く議論されており、この議論は過去時制論争 (Past Tense Debate) と呼ばれる。その中でも重要な現象の一つとして、日本語母語話者が実在語にある屈折パターンと同じように非実在語を屈折させることができない、というものがあ、これは多くの研究で観察されている [1, 2, 3, 4]。もし形態屈折の処理が規則的なものであるならば、非実在語に対しても同じ屈折パターンを適用できるはずであるということから、この現象は、形態処理が類推的であることの証左として言及されてきた。これに対して Oseki ら [5] では、現在形から過去形の方へ規則が派生する、ということがそもそも自明ではなく、日本語においては過去形から現在形の方へ規則が派生しているのではないかとの仮説を、自身の獲得研究での観察をもとに提案した。その仮説に基づき、規則ベースのモ

デルを用いて、非実在語の屈折を双方向（現在→過去、過去→現在）に検証したところ、規則ベースのモデルは過去形から現在形への方により適していることが様々な評価尺度から示された。この結果は、日本語動詞の屈折処理は、過去形から現在形への方においては規則的である、ということを示唆する。

そこで本研究では、類推に基づくモデルである深層学習モデルを用いて、日本語動詞の屈折を双方向（現在→過去、過去→現在）に学習させた。複数の実験設定を用意することで、類推ベースの深層学習モデルがどちらの時制方向の処理により適しているのかを複数の実験設定で検証した。

2 先行研究

過去時制論争は、Rumelhart ら [6] が類推ベースのモデルを提案し、彼らのモデルに対し Pinker ら [7] が問題点を指摘したことで始まった。最近では、Kirov ら [8] が当時最先端のモデルであったアテンション付きリカレントニューラルネットワーク (RNN) を形態屈折課題に応用したことで議論が再燃している [9, 10]。

一方で、日本語研究においては、母語話者は規則を用いないとする立場が一般的である。表 1 は、動詞の屈折パラダイムを示したものである。現在形の語幹末音には 11 個の子音と母音が存在し、それぞれ異なる屈折パターンをもつ。先行研究では、日本語母語話者がこの屈折パターンに則って動詞を屈折できないことが報告されている [1, 2, 3, 4]。

しかしながら、これらの先行研究では、英語と同様に日本語動詞が現在形から過去形へ屈折することが前提とされており、これは再検討が必要な点である。Oseki ら [5] は、言語獲得の観察から過去形が屈折の基底形であると仮定し、規則ベースのモデル [12, 13] を用いて非実在語を用いる形態屈折課題である wug テスト [14] を双方向（現在形→過去形と過去→現在）で行った。実験の結果、人間の正答率は

* 共同第一著者。

表 1 日本語における動詞の屈折パラダイム. 形態素境界は [11] による分析に基づく.

動詞タイプ	語幹末音	非過去	過去
母音語幹	-i	<i>mi-ru</i>	<i>mi-ta</i>
	-e	<i>tabe-ru</i>	<i>tabe-ta</i>
子音語幹	-k	<i>kak-u</i>	<i>kai-ta</i>
	-g	<i>oyog-u</i>	<i>oyoi-da</i>
	-m	<i>yom-u</i>	<i>yon-da</i>
	-b	<i>yob-u</i>	<i>yon-da</i>
	-n	<i>shin-u</i>	<i>sin-da</i>
	-r	<i>kaer-u</i>	<i>kaet-ta</i>
	-t	<i>kat-u</i>	<i>kat-ta</i>
	-w	<i>aw-u</i>	<i>at-ta</i>
	-s	<i>kas-u</i>	<i>kas-ita</i>

現在→過去方向よりも過去→現在方向の方が高く、また、モデルと人間の相関も過去→現在方向の方が高かった。さらに、現在→過去よりも過去→現在方向の方がモデルの複雑性の評価がより単純であることもわかった。このことから、Oseki ら [5] は過去→現在方向では規則による処理、現在→過去方向では類推による処理が行われている可能性があると考えしている。その場合、類推ベースのモデルである深層学習モデルの精度は過去→現在方向より現在→過去方向の方がより高くなると予測される。

3 方法

3.1 深層学習モデル

過去時制論争において、ニューラルネットワークは類推的な処理を行うモデルとして位置付けられるが、近年の形態屈折を扱う深層学習モデルには複数のアーキテクチャが存在する。複数のモデルで実験を行うことで、モデルの種類がどのように結果に影響を及ぼすかを確認するために、本研究では、形態屈折を扱うモデルとしてよく用いられる以下の2種類で実験を行った。

アテンション付き RNN 形態屈折は、機械翻訳と同様に系列変換課題として扱うことができるため、機械翻訳の分野で登場したアテンション付き RNN [15] はそのまま Kann ら [16] によって形態屈折に応用されている。その後 Kirov ら [8] によって英語における動詞の形態屈折に応用され、この文献が過去時制論争の再燃のきっかけとなった。本研究で

は Kirov ら [8] のアテンション付き RNN を再実装し、日本語動詞の屈折のために用いている。

Transformer 形態論の研究において、これまでアテンション付き RNN を用いるのが主流となっていたが、近年になり Transformer [17] を用いた研究が行われている [18, 19]。これらの研究では、モデルの縮小化やハイパーパラメータの設定など、文処理よりもデータ数の少ない形態処理を Transformer で扱うための提案が行われており、これらの工夫によって Transformer でもアテンション付き RNN の性能を超えることができるようになった。本研究では、Wu ら [18] の実装を用いて実験を行っている。これは4層のエンコーダ-デコーダ層と4つのセルフアテンションヘッドからなる、小さいモデルサイズの Transformer である。

3.2 データセットの作成

本研究では、実在語のデータセットと非実在語のデータセット (wug) の2種類を用意した。実在語のデータセットについては2つのコーパスから動詞を抽出し、それらの動詞を2種類の文字表記に変換したものを実験に用いた。この節では、動詞の抽出方法とそれらの文字表記の変換方法について述べる。

3.2.1 動詞の抽出

本研究では、2種類のコーパス (京都大学テキストコーパス, IPA 辞書) から抽出した動詞をもとに現在形と過去形のペアを作成した。これらのコーパスから抽出された動詞を合わせ重複を除いた結果、5,502 個の現在形と過去形のペアが得られた。また、実在語で訓練したモデルのテストデータとして、非実在語 (wug) のデータセットも用意した。これにより、規則ベースのモデルで実験を行った Oseki ら [5] との結果の比較が可能になる。

以下に各データセットの概要を示す。

京都大学テキストコーパス 京都大学テキストコーパス [20] は、1995 年に出版された毎日新聞と社説、それぞれ2万文に対して、形態素解析システム JUMAN、構文解析システム KNP で自動解析を行い、人手で修正が加えられたものである。このコーパスからは約 1,300 個の動詞が得られた。

IPA 辞書 IPA 辞書は、情報処理振興事業協会 (IPA) で設定された IPA 品詞体系 (THiMCO97) に基づいて作成された日本語の形態素解析用辞書であ

る [21]. 動詞を抽出した結果, 約 5,300 個の動詞が得られた.

wug 動詞 wug テストには, Oseki ら [5] で作成された 64 個の wug 動詞を用いた. これらの動詞の「語根」は, Suski[22] の動詞コーパスから抽出した動詞 1,269 個から CV 音節を抽出し, 無作為に 2 つの CV 音節を繋げることで, 基底形となる現在形と過去形がそれぞれ 32 個作成されている. ただし, n で終わる語根を持つ実在動詞は「死ぬ (*shin-u*)」のみであることから, 語根が n で終わる現在形 2 個はデータセットから除いた.

3.2.2 文字表記

本研究では, 動詞を 2 種類の文字表記に変換し実験を行った. 1 つは, アルファベットに変換したものであり (以下, latin), この変換には, Pykakasi¹⁾を用いた. もう 1 つは, アルファベット表記をさらに国際音声記号に変換したものであり (以下, IPA), この変換には phonemizer (espeak)²⁾を用いた.

3.3 実験設定

実在語のみを用いた条件では, 訓練データとテストデータを 8:2 の割合で分割した. さらに, 実在語で訓練したモデルを用いて, wug テストを行った. 正解の判定は, 実在語によるテストの場合は実在語と一致する屈折形を正解とし, 非実在語を用いた wug テストの場合は, 実在語から予測される可能な屈折形をすべて正解とした (例えば, wug 動詞 *kuhan-da* に対する正解は *kuham-u* または *kuhab-u*). 正答率は, 5 回の試行の平均を算出した. 各モデルのハイパーパラメータの設定を以下に示す.

アテンション付き RNN アテンション付き RNN のハイパーパラメータについては, Kirov ら [8] の実験設定を参照し, 設定を合わせている. 学習時の最適化アルゴリズムには AdaDelta を用い, 単語の埋め込み次元数は 300, LSTM のサイズは 100 に設定した. バッチサイズは 20 で実験を行った.

Transformer Transformer のハイパーパラメータは Wu ら [18] の実験を踏襲し, 埋め込み次元数は 256 に設定し, 学習時の最適化アルゴリズムには Adam を用いた. ただし, バッチサイズは本研究と同等のデータ数での検証を行った Ma ら [19] の実験を参照し, 32 に設定した.

1) <https://codeberg.org/miurahr/pykakasi>

2) https://github.com/bootphon/phonemizer?ref=moriogh.com&utm_source=moriogh.com

4 結果

本研究では, 2 種類のモデル, 2 種類の表記, 2 種類のテストデータセットを用いて, 双方向 (現在→過去, 過去→現在) 形態屈折の検証を行った.

はじめに, 実在語のみを用いた検証の結果を以下に示す.

表 2 実在語を用いた双方向形態屈折課題の結果

モデル	表記	正解率 (%)	
		現在→過去	過去→現在
アテンション付き RNN	latin	96.79	97.74
アテンション付き RNN	IPA	95.51	96.79
Transformer	latin	92.24	92.72
Transformer	IPA	93.37	92.98

実在語のみで検証を行ったところ, 4 条件中 3 条件で過去→現在方向の方が正答率が高い結果となった. 現在→過去方向でより正答率が高かったのは IPA 条件の Transformer のみだった. モデル間で正答率に差があり, 実在語条件では, すべての条件においてアテンション付き RNN の方が正答率が高かった.

次に, 実在語で訓練したモデルを用いて wug テストを行った. 表 3 はその結果である. 末行には Oseki ら [5] の研究で報告された人間の正答率を提示している.

表 3 双方向 wug テストの結果

モデル	表記	正解率 (%)	
		現在→過去	過去→現在
アテンション付き RNN	latin	78.98	86.80
アテンション付き RNN	IPA	81.16	94.44
Transformer	latin	93.36	92.98
Transformer	IPA	92.44	92.27
日本語母語話者 [5]	ひらがな	48	72

wug テストを行った結果, 4 条件中 2 条件で現在→過去の方が正答率が高かった. 過去→現在方向より現在→過去方向の方が正答率が高かったのは Transformer だった. モデル間でも正答率の差があり, IPA 表記で過去→現在方向に訓練・テストした条件を除いてすべての条件で Transformer のほうがアテンション付き RNN より正答率が高かった.

実験の結果, 8 条件中 5 条件において現在→過去方向の方が過去→現在方向よりも正答率が高くなることがわかり, 仮説は支持されなかった. アテンション付き RNN を用いて, 様々なデータセットの

条件を設定し検証を行った深津ら [23] の研究では、データセットのサイズに比例して正答率が高くなる結果が得られたことから、モデルの結果が訓練データの影響を受けていることが示唆された。そこで、追加の検証として、すべての条件において現在→過去方向よりも過去→現在方向においてより正答率が高かったアテンション付き RNN を対象に、kyodai データのみ、ipadic のみのデータセットを作成しモデリングを行った (表 4)。

表 4 より小さなデータセットで訓練されたモデルによる双方向 wug テストの結果

訓練データ	表記	正解率 (%)	正解率 (%)
		現在→過去	過去→現在
kyodai のみ	latin	78.00	61.87
kyodai のみ	IPA	86.00	80.00
ipadic のみ	latin	96.00	95.93
ipadic のみ	IPA	92.66	94.38

追加検証を行った結果、4 条件中 3 条件において、過去→現在方向よりも現在→過去方向でモデルの正答率がより高くなることがわかった。

5 考察

5.1 データサイズによる影響

本研究では、8 条件中 5 条件で現在→過去方向より過去→現在方向の方が正答率が高い結果が得られた。このことから、現在→過去方向の形態屈折がより類推ベースのモデルに適しているという仮説は支持されなかった。

しかしながら、より小さなデータサイズでモデリングを行った追加の検証では、4 条件中 3 条件で現在→過去方向の方が過去→現在方向よりも高い正答率となった。本研究では、5,502 個の動詞を用いて検証を行ったものの、子どもが形態屈折の獲得時に聞くとと思われる動詞の数はもっと少ない。例えば、英語の動詞屈折を規則ベースのモデルを用いて検証した Yang[24] の研究では、子どもの自然発話コーパスである CHILDES[25] から動詞の抽出した 1,042 個の動詞を用いている。そのため、限られたデータで学習する場合には、現在から過去の方向に学習するほうが類推的処理を行うのにより適している可能性がある。

5.2 データの質による影響

より小さなデータサイズでモデリングを行なった追加の検証において、正答率は、kyodai (n=1,300) > ipadic-kyodai (n=5,502) > ipadic (n=5,300) の順に高くなっており、kyodai を含むデータセットの方が含まないものより、正答率が高い。このことから、データサイズだけではなくデータの質も学習に影響を与えていることが示唆された。

例えば、モデルの正答率に影響を与えた要因として、訓練データに含まれる語彙の違いが挙げられる。kyodai は毎日新聞のデータに基づく京大テキストコーパスから抽出したものであり、より日常的な語彙で構成される。一方、ipadic は形態素解析辞書から抽出したものであり、中には複合語 (例:「明かし暮らす」) や古語 (例:「掻き暗す」)、あまり使われない語彙 (例:「こんがらがる」の類語である「こんぐらがる」) などが含まれる。このように、ipadic に比べ kyodai にはより日常的に用いられる基本的な動詞が含まれていることもモデルの正答率に影響を与えたと考えられる。

これらを踏まえて、今後の研究では、CHILDES コーパスなどをもとに、より基本的な動詞を含むサイズの小さいデータセットを作成し検証を行うことが必要である。

6 おわりに

本研究では、過去→現在方向よりも現在→過去方向での処理が類推ベースのモデルにより適していると予想し実験を行ったものの、その予想を完全に支持する結果は得られなかった。一方で、サイズの異なるデータを用いて行った追加実験では仮説が支持される結果となった。このことから、訓練データのサイズや、含まれる語彙の種類は仮説検証に大きく影響する要素であると考えられる。今後は、言語獲得研究における知見を活かしながら、実験設定やデータセットの内容をより細かく統制することで、実際に人間が形態処理を行う際の条件に近いモデリングを目指す。

謝辞

本研究は、JST さきがけ JPMJPR21C2 の支援を受けたものです。また、Transformer を用いた実験を行うにあたり、東京大学大学院の吉田遼氏にご協力いただきました。感謝申し上げます。

参考文献

- [1] Timothy J Vance. **An introduction to Japanese phonology**. State University of New York Press, Albany, NY, 1987.
- [2] Timothy J Vance. A new experimental study of Japanese verb morphology. **Journal of Japanese Linguistics**, Vol. 13, pp. 145–166, 1991.
- [3] Terry Klafehn. **Properties of Japanese Verbal Inflection**. PhD thesis, University of Hawai’i, August 2003.
- [4] Terry Klafehn. Myth of the wug test: Japanese speakers can’t pass it and English speaking children can’t pass it either. **Proceedings of the 37th Annual Meeting of the Berkeley Linguistics Society**, pp. 170–184, 2013.
- [5] Yohei Oseki, Yasutada Sudo, Hiromu Sakai, and Alec Marantz. Inverting and modeling morphological inflection. **Proceedings of the 16th Workshop on Computational Research in Phonetics, Phonology, and Morphology**, pp. 170–177, August 2019.
- [6] David E. Rumelhart and James L. McClelland. On learning the past tenses of English verbs. In David E. Rumelhart, James L. McClelland, and CORPORATE PDP Research Group, editors, **Parallel Distributed Processing: Explorations in the Microstructure of Cognition, Vol. 2: Psychological and Biological Models**, Vol. 2, chapter On Learning the Past Tenses of English Verbs, pp. 216–271. MIT Press, 1986.
- [7] S Pinker and A Prince. On language and connectionism: analysis of a parallel distributed processing model of language acquisition. **Cognition**, Vol. 28, No. 1-2, pp. 73–193, Mar 1988.
- [8] Christo Kirov and Ryan Cotterell. Recurrent neural networks in linguistic theory: Revisiting Pinker and Prince (1988) and the past tense debate. **Transactions of the Association for Computational Linguistics**, pp. 651–665, 2018.
- [9] Maria Corkery, Yevgen Matushevych, and Sharon Goldwater. Are we there yet? Encoder-decoder neural networks as cognitive models of English past tense inflection. **Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics**, pp. 3868–3877, July 2019.
- [10] Kate McCurdy, Sharon Goldwater, and Adam Lopez. Inflecting when there’s no majority: Limitations of Encoder-Decoder neural networks as cognitive models for German plurals. **Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics**, pp. 1745–1756, July 2020.
- [11] Bernard Bloch. Studies in colloquial Japanese I inflection. **Journal of the American Oriental Society**, Vol. 66, No. 2, pp. 97–109, April-June 1946.
- [12] Adam Albright and Bruce Hayes. Modeling English past tense intuitions with Minimal Generalization. **Morphological and Phonological Learning: Proceedings of the 6th Workshop of the ACL Special Interest Group in Computational Phonology (SIGMORPHON)**, pp. 58–69, July 2002.
- [13] Adam Albright and Bruce Hayes. Rules vs. analogy in English past tenses: a computational/experimental study. **Cognition**, Vol. 90, No. 2, pp. 119–61, Dec 2003.
- [14] Jean Berko. The child’s learning of English morphology. **Word**, Vol. 14, No. 2-3, pp. 150–177, August 1958.
- [15] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. **arXiv preprint arXiv:1409.0473**, 2014.
- [16] Katharina Kann and Hinrich Shültze. Single-model Encoder-Decoder with explicit morphological representation for reinflection. **Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics**, pp. 555–560, August 2016.
- [17] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, L ukasz Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, **Advances in Neural Information Processing Systems**, Vol. 30. Curran Associates, Inc., 2017.
- [18] Shijie Wu, Ryan Cotterell, and Mans Hulden. Applying the transformer to character-level transduction, 2020.
- [19] Xiaomeng Ma and Lingyu Gao. How do we get there? evaluating transformer neural networks as cognitive models for english past tense inflection. In **Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)**, pp. 1101–1114, Online only, 2022. Association for Computational Linguistics.
- [20] Sadao Kurohashi and Makoto Nagao. **Treebanks: Building and Using Parsd Corpora**, Vol. 14, chapter Building a Japanese parsed corpus while improving the parsing system, pp. 249–260. Kluwer Academic Publishers, 2003.
- [21] 浅原正幸, 松本裕治. ipadic version 2.7.0 ユーザーズマニュアル, 2003. <http://manual.freeshell.org/chasen/ipadic-ja.pdf>.
- [22] Peter M. Suski. **Conjugation of Japanese Verbs in the Modern Spoken Language**. P. D. and Ione Perkins, 1942.
- [23] 深津聡世, 原田宥都, 関澤瞭, 田村鴻希, 大関洋平. 深層学習モデルによる日本語動詞の双方向形態屈折の検証. 日本言語学会第 165 回大会予稿集, pp. 193–199, 2022.
- [24] Charles Yang. **The Price of Productivity**. MIT Press, 2016.
- [25] Brian MacWhinney. **The CHILDES project: Tools for analyzing talk**. Erlbaum, New Jersey, 1991.