

Compare Results

2/9/2026 12:45:21 PM

Summary of Comments on A91r838gh_gil2fb_ugk.tmp

This page contains no comments

Old File:

2024_emnlp-main.1055.pdf
versus
16 pages (14.23 MB)

New File:

2024_emnlp-main_1055.pdf
13 pages (314 KB)
2/8/2026 4:32:39 AM

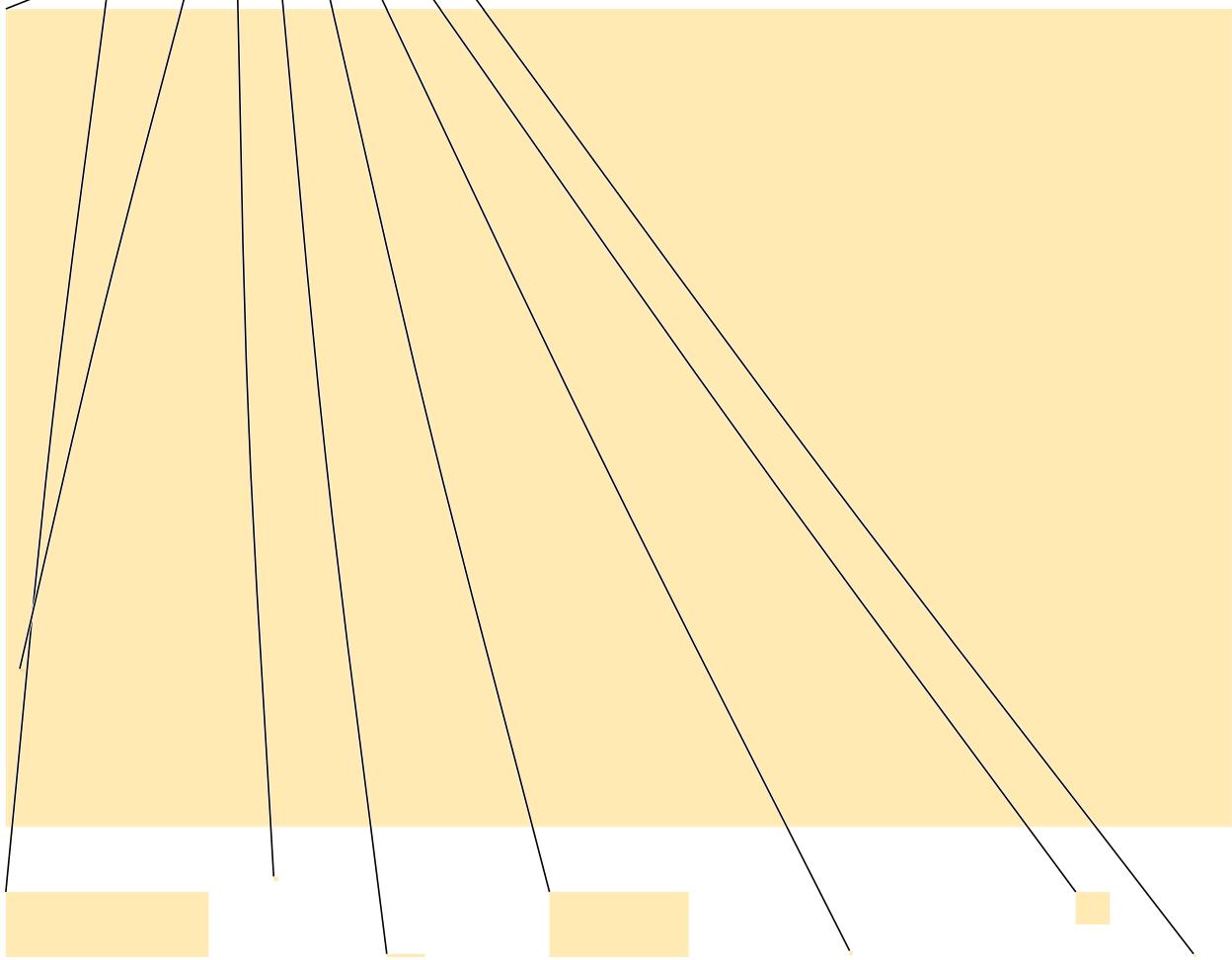
36

Total Changes	Content
9	Replacements
12	Insertions
15	Deletions

[Go to First Change \(page 1\)](#)

Page: 1

Image Replaced



Getting The Most Out of Your Training Data: Exploring Unsupervised Tasks for Morphological Inflection

Katharina von der Wense¹ Georgetown University² University of Colorado Boulder³ Johannes Gutenberg University Mainz⁴

Abhishek Punshotthaman¹ Adam Wiernerslage¹

Abstract

Pretrained transformers such as BERT (Devlin et al., 2019) have been shown to be effective in many natural language tasks. However, they are under-explored for character-level sequence-to-sequence tasks. In this work, we investigate pretraining transformers for the character-level task of morphological inflection in several languages. We compare various training seeds and secondary tasks where unsupervised data taken directly from the target task is used. We show that training on secondary unsupervised tasks increases inflection performance even without any external data, suggesting that models learn from additional unsupervised tasks themselves. In addition, we find that standard denoising tasks can hurt in multi-task setups, but using external data for denoising solves this issue.

1 Introduction

Morphological inflection is the task of generating a word form given a lemma and a set of morphological features. For example, given the lemma *walk* and features *V;PST*, the inflected form is *walked*. Morphological inflection is a core task in computational morphology and has been explored in a variety of settings, including as part of the SIGMORPHON-UniMorph shared tasks (Výbornová et al., 2020; Pinente et al., 2021; Kothiyal et al., 2022; Goldman et al., 2023). In this work, we focus on low-resource morphological inflection, where the amount of training data is limited.

Recent advances in NLP have been driven by pretraining large-scale models on large amounts of text using self-supervised objectives such as masked language modeling (Devlin et al., 2019) and denoising sequence-to-sequence objectives (Lewis et al., 2020; Raffel et al., 2019). While such objectives have been successful for token-level tasks and semantic tasks, they are less explored for character-level sequence-to-sequence tasks such as morphological inflection. At the same time, the computational morphology community is frequently interested in low-resource languages and settings where large external corpora may be unavailable.

We explore whether unsupervised tasks constructed from the available morphological inflection data can improve model performance, and when denoising objectives may hurt in multi-task learning. We investigate:

- Pretraining setups using only the target task data, including two-stage training and multi-task learning
- Secondary tasks: character-level masked language modeling and character-level autocodecoding.
- The effect of using additional unlabeled data from Universal Dependencies (UD) treebanks for denoising tasks.

150-639-2 Language UD treebank used I 150-639-2 Language UD Treebank used arb Arabic – Gulf Arabic-PADT ita Italian Italian-SDT amh Amharic Amharic-ATT jpn Japanese Japasse-GSD arz Arabic_Egyptian kat Georgian bel Belarusian-BSE kir Khaling dan Danish Danish-DDT mkd Macedonian deu German German-GSD nay Navajo eng English English-Atis rus Russian Russian-GSD fin Finnish Finnish-FTB san Sanskrit Sanskrit-UFAL fra French French-GSD sue Sami_North_Sami_Glottal grc Ancient Greek AncientGreek – Persesuspro Spanish-Spanish – AnCorpa heb Hebrew-Hebrew – HTBsAlbanian heb_(nac)Hebrew_Unicode-edsavaSechili hunHungarianHungarian-SzegedtunTurkishTirrkish – Atis hye-EasternArmenian-Armenian – ArmTDP

Table 1: The 27 typologically diverse languages (Subsection 4.1) from the 2023 shared task [M1](#) of which are investigated in this work. We use some UD treebanks for analytical experiments in Subsection 6; the specific treebanks are listed in the final column.

2 Related Work

Multi-task learning (MTL) has long been studied as a way to improve generalization by training on multiple related tasks (Cartanha, 1997; Luong et al., 2016). In NLP, intermediate-task and supplementary training can improve downstream performance (Phang et al., 2018; Prukaschatkun et al., 2020). Identifying beneficial task relations has also been explored (Fringel and Sogard, 2017; Martinez Alonso and Plank, 2017; Fifty et al., 2021).

Denoising objectives such as MLM (Devlin et al., 2019) and sequence-to-sequence denoising models, ByT5 (Xue et al., 2022) explores byte-to-byte pretraining. However, the role of such objectives in low-resource character-level tasks is less clear.

Morphological inflection has a rich history, including shared tasks and neural approaches using encoder-decoder models and transformers (Kann and Schütze, 2016; Wit et al., 2021). Unlabeled data has been used for morphological generation (Kann and Schütze, 2017), and dataset quality and sampling issues have been studied (Koehn et al., 2023; Muradoglu and Hulden, 2022). Noise in morphological inflection has also been investigated (Wienershage et al., 2023).

3 Morphological Inflection

We follow the standard formulation: given a lemma and morphological features, generate the inflected form. Inputs and outputs are treated as character sequences. We focus on typologically diverse languages from the SIGMORPHON-UniMorph 2023 shared tasks (Goldman et al., 2023), and create low-resource training subsets.

4 Data

We use the 27 languages from the SIGMORPHON-UniMorph 2023 shared task (Goldman et al., 2023). For each language, we subsample training data to simulate low-resource scenarios.

5 Baseline Model

Our baseline is a transformer encoder-decoder model operating at the character level (Wu et al., 2021), implemented with *yogodyne*.

6 Training Methods

We evaluate several training setups.

6.1 Baseline

We train the model on the supervised morphological inflection task only.

6.2 Two-stage pretraining (PT)

We perform two-stage training: first train on an unsupervised objective using unlabeled data derived from the tasks data, then fine-tune on supervised morphological inflection.

6.3 Multi-task learning (MTL)

We train on the supervised task jointly with an unsupervised auxiliary task. For MTL, the loss is a combination:

$$\mathcal{L} = \mathcal{L} * \text{SUP} + \lambda \mathcal{L} * \text{UNSUP}. \quad (1)$$

We set λ to [ILLEGIBLE].

6.4 Auxiliary tasks

We use two auxiliary tasks:

- Character-level masked language modeling denoising (CMLM): apply noise to the input sequence and train to reconstruct the original.

- Autoencoding (AE): reconstruct the original sequence from itself.

6.5 Noise

We use a span-based corruption process. Let x be the input sequence. We sample spans and apply replacements [ILLEGIBLE]. The mask sampling rate is a hyperparameter [ILLEGIBLE]. We also explore using external unlabeled data from UD treebanks.

7 Results and Analysis

We report development and test accuracies for each language and model variant.

7.1 When Does Denoising Hurt MTL?

There is a remarkable gap in performance between MTL-AE and MTL-CMLM. The CMLM denoising objective is the worst performing setup, performing below the baseline on average. In further analysis, performing CMLM on external data that is separate from the finetuning data solves this issue, resulting in significantly better performance.

	Baseline	PT-CMLM	PT-AE	MTL-CMLM	MTL-AE	Languages	rso	639-2	Dev	Test	Dev	Test	Dev	Test	Arabic	Gulf	afb	68.8	69.4	72.2	70.5	72.7	t	7	t		
66.8	67.8	72.7	77.72.7	72.7	72.7	Ammaric	amn	44.6	42.9	48.0	50.8	56.5	60.9	34.9	36.7												
56.5	61.4	Arabic	Egyptian	atz	82.8	82.5	83.1	83.9	82.3	84.3	80.7	81.4															
83.6	83.8	Belarusian	bel	61.2	59.0	62.9	68.8	61.5	58.7	59.8	56.5	64.4															
61.7	Danish	dan	81.7	80.	78.1	81.2	79.9	80.0	80.7	83.2	82.5																
German	68.2	77.2	70.3	68.7	74.4	74.4	73.4	65.8	65.1	77.7	74.3	73.2	English														
eng	91.6	88.2	91.5	88.6	91	89.0	3	89.5	87.2	92.3	90.9	Finnish	fin	74.6													
56.7	77.7	61.6	78.2	61.8	58.9	44.0	81.4	68.6	French	fra	15.2	66.7	76.9														
68.0	80.6	68.9	69.9	67.7	70	81.1	73.6	Ancient	Greek	c54.0	63.160.441.352.834.5	+	-1.329.656.640.7	Hebrewhet	74.272.176.676.637.77.676.6372.272.6180.377.95	Hebrew	Uncocalizydhew	voc81.5									

Table 2: The development and test accuracies of the 5 model variants, for all the 27 languages. For each language, the highest development accuracy is underlined and highest test accuracy is bolded.

IMAGE NOT PROVIDED

Figure 1: Figure 1: The distribution of performance (test set accuracy) for each model variant on the various data sizes. Distributions are plotted as violin plots, with box plots visualizing the mean, first and third quartile, and min and max values.

7.2 External Data for Denoising

We prepare additional unlabeled data from UD treebanks and use it for denoising in MTL. Results are shown below.

8 Future Work

The denoising tasks requires hyperparameters for the instrumentation of the noise. Due to this, further work is required in exploring these tasks under different hyperparameter settings with multiple methods to shed light on their sensitivity and ability to improve models for character-level tasks such as morphological inflection and G2P. Future work should also consider exploring more secondary tasks, especially based on particular morphological phenomenon in diverse languages.

Limitations

- Our work is limited to the character-level task of morphological inflection. Thus, findings may not hold for other similar tasks such as G2P and interlinear glossing.
- Considering the sensitivity of training methods to vocabulary and data sizes, it is unclear whether these results can be extrapolated to different scenarios.
- Our work does not explore the disparity of performance of the methods across languages and requires expert analysis over various of linguistic features.

Baseline MTL-QMLM MTL-AE MTL-CMLM, UD MTL-AE-UD Language ISO 639-2 Dev Test	
Dev Test	Dev Test Dev Test Arabic Gulf abf 68.8 69.4 68.8 67.8
Dev Test	Dev Test Dev Test Amharic ethi 64.6 42.9 34.9 36.7 56.5 65.3
5	5.7.7 61.0 66.6 Belarusian bel 67.2 59.0 59.8 56.5 64.4 67.7 64.2 61.5 65.3
62.2	Danish dan 81.7 80.1 80.0 80.7 83.2 82.5 82.3 80.8 83.7 82.9 German deu
68.2	77.2 65.8 65.7 74.3 73.2 75.4 74.4 75.4 76.3 English eng 91.6 88.2 89.8
87.2	92.3 90.9 91.3 88.5 97.4 96.7 98.9 Finnish fin 74.6 56.7 58.9 44.0 81.4 66.6
81	5.7.8 82.7 73.6 French fra 75.2 68.2 69.9 67.0 81.1 73.6 82.8 75.2
85.8	74.1. Ancient Greek grc 54.5 33.1 43.3 32.8 56.6 40.7 64.2 46.5 63.5
47.t	Hungarian hun 75.7 65.7 65.4 61.3 80.4 77.1 81.1 75.2 83.6 78.1 Hebrew heb 74.2 77.2 72.2 78.3 77.5 79.3 75.7 Eastern Armenian
lyre	'79.2 79.4 76.8 76.0 86.9 89.5 90.5 89.0 98.4 93.0 I J t a a p l i a a n n
esse	j i a ta p 9 1 0 5 . . 5 8 8 2 5 0 . . 7 1 8 4 3 . 1 . 3 7 5 t . . 8
4.9	1 4 5 . . 0 4 2 9 0 t . . 9 4 9 3 4 4 . . 8 3 8 3 8 2 . . 7 2 9 4 4
4.	3 1 9 4 3 2 . . 8 3 Russian rus 78.7 72.7 72.7 72.7 72.7 72.7 72.7 72.7 72.7 81
.7	80.1 81.8 82.9 Sanskrit san 55.0 49.0 47.6 50.5 63.4 56.4 65.4 57.9 65.7
58.3	Sami sme 57.3 43.9 44.2 33.8 70.0 60.4 70.2 66.7 74.8 66.3 Spanish spa
88.2	88.0 19.3 78.9 91.6 90.9 91.5 90.3 91.8 91.8 Turkish tur 85.3 86.
76.4	73.4 89.7 88.5 87.5 85.9 89.6 89.9 Avg 69.51 64.39 62.45 68.98 74.88
'7t, 28	76.3t 72.22 78.09 74.66

Table 3: Results for our models by language from the experiments with external data, reporting development and test accuracy. For each language, the highest development accuracy is underlined and highest test accuracy is bolded. Note: results for non-UD models are identical to Table 2.

Acknowledgments

We thank the anonymous reviewers for their useful suggestions and feedback and the NALA Lab at the University of Colorado Boulder. This work utilized the Blanca condo computing resource at the University of Colorado Boulder. Blanca is jointly funded by computing users and the University of Colorado Boulder.

References

- [1] Sina Ahmadi and Aso Malapodi. 2023. Revisiting and amending Central Kurdish data on UniMorph 4.0. In Proceedings of the 20th SIGMORPHON workshop on Computational Research in Phonetics, Phonology, and Morphology, pages 38–48, Toronto, Canada. Association for Computational Linguistics.
- [2] Lucas F. E. Ashby, Travis M. Bartley, Simon Clematide, Luca Del Signore, Cameron Gibson, Kyle Gorman, Leonju Lee-Silkka, Peter Makarov, Aidan Malanowski, Sean Miller, Ohnar Ortiz, Reuben Raiff, Arunudhati Sengupta, Bora Soe, Yulia Spector, and Winnie Yan. 2021. Results of the second SIGMORPHON shared task on multilingual grapheme-to-phone conversion. In Proceedings of the 18th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology, pages 115–125, Online. Association for Computational Linguistics.
- [3] Khuyaebtaar Batturen, Omer Goldman, Salam Khalifa, Nizar Habash, Witold Kiereg, Gébor Bella, Brian Leonard, Garrett Nicolai, Kyle Gorman, Yustimus Ghangoo Ate, Maria Ryskina, Sabrina Mielke, Elena Budianskaya, Charbel El-Khaissi, Tiago Pimentel, Michael Gasser,

- William Abbott Lane, Mohit Raj, Matt Coler, Jaime Rafael Montoya Samane, Delio Soticonauzi Camaiteri, Esaf Ziu-maeuta Rojas, Didier Lopez-Francés, Arturo Oncelay, Juan Lopez-Bautista, Gema Celeste Silva Villegas, Lucas Tonoba Henning, Adán Ek, David Gurie, Peter Dirix, Jean-Philippe Bernardi, Andrey Scherbakov, Aziyana Bayyr-zool, Antonios Anastasopoulos, Roberto Zariquey, Karina Scheifer, Sofya Canieva, Hilaria Cruz, Ruiyiin Karaifla, Stella Markantonatou, George Pavlidis, Matvey Plugarov, Elena Klyachko, Ali Salehi, Candy Angulo, Jatayu Baxi, Andrew Krizhanovsky, Natalia Krizhanovskaya, Elizabeth Salesky, Clara Vania, Sardana Ivanova, Jennifer White, Rowan Hall, Maudslay, Josef Valvoda, Ran Zmigrod, Paula Czamowska, Irene Nikkarinen, Adita Saichak, Brijesh Bhatt, Christopher Strangehn, Zoey Liu, Jonathan North Washington, Yival Pinter, Diygu Atamian, Marein Wolinski, Totok Suthardijanto, Anna Yablonitskaya, Niklas Stoehr, Hossep Dolatian, Zahoh Nurjiah, Shyam Ratan, Frances M. Tyers, Edoardo M. Ponti, Grant Aiton, Aryaman Arora, Richard J. Hatchell, Ritesh Kumar, Jereniak Young, Daria Rodionova, Anastasia Yemelina, Taras Andrusko, Igor Marchenko, Polina Maskovtseva, Alexandra [ILLEGIBLE].
- [4] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of deep bidirectional transformers for language understanding. [ILLEGIBLE].
- [5] Falcon and The PyTorch Lightning team. 2019. PyTorch Lightning.
- [6] Christopher Fifty, Ehsan Amid, Zhe Zhao, Tianhe Yu, Rohan Anil, and Chelsea Finn. 2021. Efficiently identifying task groupings for multi-task learning. In *NLP Information Processing Systems*.
- [7] Omer Goldman, Khuyaagbaatar Batsuren, Salam Khalifa, Aryaman Arora, Garrett Nicolai, Reut Tsarfaty, and Ekaterina Vylomova. 2023. SIGMORPHON-UniMorph 2023 shared task 0: Typologically diverse morphological inflection. In *Proceedings of the 20th SIGMORPHON workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 117–125, Toronto, Canada. Association for Computational Linguistics.
- [8] Katharina Kann and Hinrich Schütze. 2016. Single-model encoder-decoder with explicit morphological representation for reinflection. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 555–560.
- [9] Katharina Kann and Hinrich Schütze. 2017. Unlabeled data for morphological generation with character-based sequence-to-sequence models. In *Proceedings of the First Workshop on Subword and Character Level Models in NLP*, pages 76–81, Copenhagen, Denmark. Association for Computational Linguistics.
- [10] Christo Kirov, John Stake-Glassman, Roger Que, and David Yarowsky. 2016. Very-large scale parsing and normalization of Wiktionary morphological paradigms. In *Proceedings of the Tenth International Conference on Language Resources and Evaluation (LREC’16)*, pages 3121–3126, Portorož, Slovenia. European Language Resources Association (ELRA).
- [11] Jordan Koehler, Salam Khalifa, Khuyaagbaatar Batsuren, Hossep Dolatian, Ryan Cotterell, Faruk Akkus, Antonios Anastasopoulos, Taras Andrusko, Aryaman Arora, Nona Atamalov, Géfió Bela, Elena Budianskaya, Yustinus Ghanegoo Ae, Omer Goldman, David Gurie, Simon Guriel, Silvia Gurie-[Gashvili], Witold Kiera5, Andrew Krizhanovsky, Natalia Krizhanovsky, Igor Marchenko, Magdalena Markowska, Polina Mashkovtseva, Maria Nepomniashchaya, Daria Rodionova, Karina Scheifer, Alexandra Sorova, Anastasia Yemelina, Jeremiah young, and Ekaterina Vylomova. 2022. SIGMORPHON-UniMorph 2022 shared task 0: Generalization

- and typologically diverse morphological inflection. In Proceedings of the 19th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology, pages 176–203, Seattle, Washington. Association for Computational Linguistics.
- [12] Jordan Kothe, Sarah Payne, Salam Khalifa, and Zoey Litr. 2023. Morphological inflection: A reality check. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 6082–6101, Toronto, Canada. Association for Computational Linguistics.
- [13] Kundan Krishna, Saarabh Garg, Jeffrey Bigham, and Zachary Lipson. 2023. Downstream datasets make surprisingly good pretraining corpora. In Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 12201–12222, Toronto, Canada. Association for Computational Linguistics.
- [14] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 7871–7880, Online. Association for Computational Linguistics.
- [15] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach.
- [16] Thang Luong, Quoc V. Le, Ira Sutskever, Oriol Vinyals, and Lukasz Kaiser. 2016. Multi-task sequence-to-sequence learning. In International Conference on Learning Representations.
- [17] Héctor Martínez-Alonso and Barbara Plank. 2017. When is multitask learning effective? semantic sequence prediction under varying data conditions. In Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers, pages 44–53, Valencia, Spain. Association for Computational Linguistics.
- [18] Saliha Muradoglu and Mans Hulden. 2022. Eny, meny, many, more: how to choose data for morphological inflection. In Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing, pages 294–303, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- [19] Joakim Nivre, Marie-Catherine de Marneffe, Filip Ginter, Jan Hajd, Christopher D. Manning, Saúlao Pysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. 2020. Universal Dependencies v2: An evergrowing multilingual treebank collection. In Proceedings of the Twelfth Language Resources and Evaluation Conference, pages 4034–1043, Marseille, France. European Language Resources Association.
- [20] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. 2018. Deep contextualized word representations. In Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers), pages 2227–2237, New Orleans, Louisiana. Association for Computational Linguistics.
- [21] Jason Phang, Thibault F6vry, and Samuel R. Bowman. 2018. Sentence encoders on stiffs: Supplementary training on intermediate labeled-data tasks. arXiv preprint arXiv:1811.01088.

Page Inserted
Page Deleted

- [22] Tiago Pimentel, Maria Ryskina, Sabrina J. Mielke, Shijie Wu, Eleanor Chodroff, Brian Leonard, Garrett Niclai, Yustinius Ghanggo Aye, Salam Khalifa, Nizar Habash, Charbel El-Khaïssi, Omer Goldman, Michael Gasser, William Lane, Matt Coler, Arturo Oncevay, Jaine Rafael Montoya Samanie, Gema Celeste Silva Villegas, Adam Ek, Jean-Philippe Bernardy, Andrey Sheherbekov, Azizava Bayyrzhol, Karina Sheifer, Sofya Ganjeva, Matvey Plugarov, Elena Klyachko, Ali Salehi, Andrew Krizhanovsky, Natalia Krizhanovsky, Clara Vania, Sardana Ivanova, [ILLEGIBLE], and Ekaterina Vlyonova. 2021. SIGMORPHON 2021 shared task on morphological reinflection: Generalization across languages. In Proceedings of the 18th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology, pages 229–239, Online. Association for Computational Linguistics.
- [23] Yada Pruksachatkun, Jason Phang, Haokun Liu, Phu Mon Hrut, Xiaoyi Zhang, Richard Yuanzhe Pang, Clara Vania, Katharina Kann, and Samuel Bowman. 2020. Intermediate-task transfer learning with pre-trained language models: When and why does it work? In Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, pages 5231–5241.
- [24] Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *J. Mach. Learn. Res.*, 21:140:1–140:67.
- [25] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems*, 30.
- [26] Pascal Vincent, H. Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre-Antoine Manzagol. 2010. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *J. Mach. Learn. Res.*, 11:3371–3408.
- [27] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stéfan J. van der Walt, Matthew Brett, Joshua Wilson, K. Jarrod Millman, Nikolay Mayorov, Andrew R. J. Nelson, Eric Jones, Robert Kern, Eric Larson, C. J. Carey, illian Polat, Yu Feng, Eric W. Moore, Jake VanderPlas, Denis Laxalde, Josef Perktold, Robert Cimrman, Ian Henriksen, E. A. Quintero, Charles R. Harris, Anne M. Archibald, Antoniú H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. 2020. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272.
- [28] Ekaterina Vlyonova, Jennifer White, Elizabeth Salesky, Sabrina J. Mielke, Shijie Wu, Edoardo Maria Ponti, Rowan Hall, Maudslay, Ran Zingrid, Josée Valova, Svetlana Toldova, Francis Tyers, Elena Klyachko, Illya Yegorov, Natalia Krizhanovsky, Paula Carnowska, Irene Nikkarinen, Andrey Krizhanovsky, Tiago Pimentel, Lucas Torralba Hennigen, Christo Kirov, Garrett Nicolai, Adina Williams, Antonios Anastasopoulos, Hilaria Cruz, Eleanor Chodroff, Ryan Cotterell, Miljka Silfverberg, and Mans Hulden. 2020. SIGMORPHON 2020 shared task 0: Typologically diverse morphological inflection. In Proceedings of the 17th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology, pages 1–39, Online. Association for Computational Linguistics.
- [29] Adam Wiemerslage, Changling Yang, Garrett Nicolai, Miljka Silfverberg, and Katharina Kann. 2023. An investigation of noise in morphological inflection. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3351–3365, Toronto, Canada. Association for Computational Linguistics.

[30] Shijie Wu, Ryan Cotterell, and Mans Hulden. 2021. Applying the transformer to character-level translation. In Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume, pages 1901–1907. Online. Association for Computational Linguistics.

- [31] Linting Xie, Aditya Barua, Noah Constant, Rani Al-Rfou, Sharar Narang, Mihir Kale, Adam Roberts, and Colin Raffel. 2022. ByT5: Towards a token-free future with pre-trained byte-to-byte models. *Transactions of the Association for Computational Linguistics*, 10:291–306.
- [32] Daniel Zeman, Joakim Nivre, Mitchell Abrams, Elia Ackermann, Noemi Aeppli, Hamid Aghaei, Zeljko Agic, Amir Ahmadi, Lars Alrenberg, Chika Kennedy Ajede, Salih Furkan Akkurt, Gabriele Aleksandraviciute, Ika Alfina, Aymeric Algoan, Chiara Alzetta, Erik Andersen, Lene Antonsen, Tatsuya Ayama, Kaito Aplonova, Angelina Aquino, Carolka Aragon, Glyd Aragues, Maria Jesus Aranza-be, Bilge Nas Aucan, Ronnun Amardottir, Greshaw Arntie, Jessica Naraiswari Arwidarasti, Masayuki Asahara, Kayla Asgeisafotir, Deniz Baran Aslan, Cengiz Asnazodlu, Luma Ateyah, Funkan Attia, Mohammed Attia, Aitzazur Aturxa, Liesbeth Augustinus, Mariana Avelds, Elena Badmaeva, Keerthan Balasubramani, Miguel Ballisteros, Esila Bauerjee, Sebastian Bank, Virginica Barbu Mititelu, Stephan Barkarson, Rodolfo Basile, Victoria Basnay, Colin Batchelor, John Bauer, Sayritalka Bedir, Shahnam Belizad, Kepa Bengoechea, Ibrahim Benli, Yilai Ben Moshe, Gedze Berk, Riyaz Ahmad Bhat, Erica Biagetti, Eckhard Bick, Agne Bielinskasienė, Kristin Biapadafitir, Roger Blokland, Victoria Bobicev, Loïc Bozon, Emanuel Borges Vdker, Carl Björstell, Cristina Bosco, Gosse Bonnia, Sam Bowman, Adriane Boyd, Anouck Bragaar, Auféno Branco, Kristina Brodaite, Aljoscha Burchardt, Marisa Camp-Pos, Marie Candito, Bernard Caron, Gauthier Caron, Catarina Carvalheiro, Rita Carvalho, Lauren Cassidy, Maria Clara Castro, Sfoglio Castro, Tatiana Caval-canti, Gilgen Cebiroglu Eryilm, Flavio Massimiliano Cecolini, Giuseppe G. A. Celano, Slavomir Čipčík, Nelišhan Cesur, Savas Çelik, Önder Oqtanlu, Fabri-icio Chaibh, Liyanage Chamila, Shweta Chatiani, Ethan Choi, Taishii Chikka, Yongseok Cho, Jinho Choi, Jayee Chin, Jayeon Chung, Alessandra T. Cizmarová, Silvie Cirková, Andree Collomb, Qalir Qollikin, Miriam Connor, Daniela Corbetta, Frausto Costa, Marine Courtin, Mihuela Cristescu, Ingerid Ipyning Dale, Philipeon Daniel, Elizabeth Davidson, Leonel Figueredo de Alencar, Mathieu Dehonick, Martine de Laurentiis, Marie-Catherine de Marinelle, Valeria de Paiva, Mehmet Ozan De-erim, Flávia de Souza, Aranaa Diaz de Ilarratza, Carly Dickerson, Arawinda Dinakaranani, Elisa Di Nuovo, Bamba Dione, Peter Dirix, Kaja Dobrovolc, Adrian Doyle, Timothy Dozat, Kira Proganova, Punet Dwivedi, Christian Ebert, Henne Eckhoff, Federica Gamba, Marcos Garcia, Moa Gilardens fors, Fabrizio Ferraz Gerardi, Kim Gendes, Luke Gessler, Filip Zinner, Sandra Eich, Marhaba Eli, Ali Elkahly, Binyam Ephrem, Olga Ejina, Tomasz Erjavec, Farah Essadi, Mine Etienne, Wogaine Eveyni, Sidney Facundes, RichSrd Farkas, Federica Favero, Jamail Ferdaousi, Marilia Fernanda, Hector Fernandez Al-calde, Anal Fethi, Jennifer Foster, CIS/udia Freitas, Kazunori Fujita, Katarína Gažišová, Daniel Galbraith, Federica Gamba, Marcos Garcia, Moa Gilardens fors, Fabrizio Ferraz Gerardi, Kim Gendes, Luke Gessler, Filip Zinner, Gustavo Godoy, Jakes Goenaga, Koldo Gojenola, Memduh Gökmen, Yoav Goldberg, Xavier Gómez Guinovart, Berta González Saavedra, Bernadete Gracifte, Matias Grioni, Ilic Grobel, Noémunds Gruzitis, Bruno Guillot-Barbance, Tonga Giinged, Nizar Halash, Hinrik Halsteinsson, Jan Hajid, Ian Hajid jr., Mika Hiimilli-inen, Linh HANH, Na-Rae Han, Mohammad Yedid-Yaakov, Takaaki Harada, Sam Hardwick, Kim Harris, Dag Haug, Johannes Heinecke, Oliver Hellwig, Felix Hennig, Barbara Hoadl6, Jaroslava Hlaváčková, Florinel Hociung, Petter Hohle, Marvel Huerta-Mendez, Jenia Hwang, Takumi Ikeda, An-ton Karl Ingason, Radu Ion, Elena Irinina, Olijide Ishola, Artan Islamaj,

Kaoru Ito, Stratun Janata, Tom65 Jelfuek, Apoora Jha, Katharine Jiang, An-ders Johansen, Hildur Jónsdóttir, Fredrik Jor-gensen, Markus Juntura, Hinne Kadraba, Nadezhda Karabéeva, Sylvain Kahane, Hiroshi Kanayama, Jenna Kanerva, Nesilhan Kara, Rirviin Karash6ia, Andrej Kasen, Tolga Kayadelen, Sarvesvaran Keigath-haraiyer, Yelvaya Kettnerovo, Jesse Kirchner, Elena Klementieva, Elena Klyachko, Arne Kdhn, Abdül-latif Kolossal, Kamil Kopacewicz, Tino Korkakangas, Mehmet Kdse, Alexey Koshevoy, Natalia Kosyba, Jolanta Kowalevskaite, Simon Krek, Parameswari Krishnamurtti, Sandra Kubler, Adrian Kuqi, Oluzhan Kuyruklu, Ash Kuzgun, Sookyung Kwak, Kris Kyle, Veronika Laippala, Lorenzo Lamberti, Ta-tiana Landu, Septima Dian Larasati, Alexei Lavrentiev, John Lee, Phng Lo Hdng, Alessandro Lenci, Sarah Lerpradit, Hernan Leung, Maria Levine, Lauren Levine, Cheuk Ying Li, Josie Li, Keying Li, Yixuan Li, Yuan Li, KyungTae Lin, Bruna Lima Padovani, Yi-Ju Jessica Lin, Kristfer Lind6in, Yang Janet Liu, Nikola Ljub6idi, Olga Loganova, Ste-fano Lusito, Andry Luthfi, Mikko Luukko, Oleg Lyashkovskaya, Teresa Lynn, Vivien Mackenz, Menel Mahmut, Jean Maillard, Ilya Makarshuk, Alibek Makazhanov, Michael Mandl, Christopher Manning, Ruli Mamrung, Bi-igra Margan, Cetlina Mch5n-dic, David Marecek, Katrin Macheinecke, Stella Markantonatou, Héctor Martínez Alonso, Lorena Martín Rodríguez, Aydito Martins, Cl6udia Mar-tins, Jan Mašek, Hiroshi Matsuda, Yuiji Matsunoto, Alessandro Mazzei, Ryan McDonald, Sarah McGuinness, Gustavo Mendonga, Tatiana Merzhevich, Yko Miecka, Aaron Miller, Karina Mischenkova, Anna Misselli, C5tilin Mittleh, Maria Mito-ofan, Yusuke Miyao, AmirHossein Mojiri Foroushani, Judit Mol-niir, Amirsaeid Moloodi, Zinonetta Montenagni, Amir More, Laura Monzo Romero, Giovanni Moretti, Shinsuke Mori, Tomohiko Morioka, Shigeaki Moro, Bjartur Mortensen, Boldan Moskalevskyi, Kadri Muusinek, Robert Munro, Hugo Murawski, Kaili Müürisepp, Pihlea Naimannai, Marija Naklić, Juan Ignacio Navarro Horflacke, Anna Nedoluzhko, Gunta Nejbore-Berzkalne, Janelia Nevaci, Ling Nguyen Thi, Huynh Nguyen Thi Minh, Yoshi-hiro Nikaido, Vitaly Khokarev, Rattima Nitis-aroj, Alireza Nourian, Hanna Nurmi, Sanna Ojala, Atul Kr. Ojha, Juulda Glad6ittir, Ad6layo Ohidkun, Mai Onnra, Enoka Onwugbuzia, Noam Ordan, Pepe Osenoro, Robert Osthang, Lrlja Øvreid, azyleB-estlO-zateg, MerveO-gekik, ArzucanO-ghir, Balkr-OzilirRagaran, TeresaPuccosi, AlessioPalme Passos, Ghi-lia Pedanose, Angelika Pejajo, Raji ska, Siyao Peng, Sjogolagan Peng, Rita Pereira, SFlaviaF AugustoPerez, NataliaPerkova, GunPerter, SlavPetrov, DariaPetron, AndreaFeverelli, JasonPhe-lan, JussiPitulainen, YavalPinter, CeliaPinto, TommiAPirinen, EmelyPitter, MagdalenaPlanada, Barker maker, MizanurRahman, TarakaRama, LoganathanRamasamy, JoanaRamos, FanRashel, Moham-madSadighRosoli, VinilRavishankar, LucyReal, PetruRebeja, SivaReddy, MathildeRegnault, GeorgRehr son, ManuelaSangalli, EgiSauntry, DagsStry, MartaSartor, BarbaraSaurie, YannisSci-er, EinerFreym, Sigurdsson, JodoSilva, AlineSil-vera, NataliaSilveira, SaraSilveira, MariaSimi, Radu-nava, TelStther, MariaSkachekabova, AaronSmith, IsabelaSoares — Bastos, Per ErikSöder, BarbaraSommerhauser, ShafiqSourov, RachelleSprungoli, Vi-vianStanou, StephanSteiggrimsso, AntonStella, AbishekStephen, MilanStraka, EmmettStrickland, J barun, MaryJAnC.Tan, TokaokiTanaka, DipaliTanyaa, MirkaTavoni, SamsonTeila, IsabelleTelier, Mar ers, SvenbjörnThiordarson, VilhjalmurIvarsson, SunireUmatua, RomanUnilon, ZdeňkaTrčová, La-jananNord, ViktorVarga, UlianaVedenna, GiuliaVenturi, VeronikaVincze, NataliaVlasova, AgaWaka gailWofsy, JonathanNorthWashington, MaximilanWendt, PaulWidmer, ShiraliYilderson, Sv-HartariW terp, ZsuzsannaWoldenariam, Tak-snannWong, AlinaWrbilewska, MaryYako, KajoyYannista, NokiYana

A Data details

A.1 Limitations of UniMorph and SIGMORPHON

The unimorph project is the primary source for the dataset. It draws heavily from Wiktionary¹[(<https://www.wiktionary.org/>)]. Wiktionary is a collaboratively built resource in a semi-automated way based on Kirov et al. (2016). Wiktionary is not a linguistic resource that is considered as gold-standard data. The semi-automated methodology, sources, and broad mandate limits the utility and effectiveness of the dataset. A notable example is Ahmadi and Mahnfi (2023), which discusses this in the context of Sorani (ckb) also known as Central Kurdish (not one of the 27 languages in this work). The limitations of the dataset used in this work, being only very recently released, are not well-studied, and consequently also apply to our work.

A.2 [ILLEGIBLE]

Selection and Sampling

Many features of morphological inflection data, such as overlap and frequency, have been shown to be important factors for model performance (Kodher et al., 2023). (Muradetlu and Hulden, 2022) demonstrated how data could be sampled using active learning methods to improve model performance. Since we investigate training methods rather than data methods, we perform analysis on data which has been selected specifically for benchmarking purposes. We recommend the readers check Section 4 “Data preparation” of the shared task paper Goldzman et al. (2023) for more information on the data methods used for target-task data selection and splits. We discuss details relevant to our selection and sampling below.

Lemma Overlap The 2023 shared task dataset was specifically designed to prevent lemma overlap between any of dev, train, and test. Since we only sub-sample from train, the lack of lemma overlap is maintained in our datasets, and is thus not a relevant point of analysis as in other work (e.g. Kodher et al. (2023)).

A.3 Preparing Additional Data from UD Treebanks

With a fixed seed, we randomly sample words from the selected UD Treebank to prepare an unlabeled training set of size 2k for each language. We perform sampling only after filtering out NUM and PUNCT tagged and tokenized words (Nivre et al., 2020). We do not otherwise use the token-level annotations from UD, simulating a more realistic data setting than the one UniMorph words represent. Table 1 shows the 10 languages from the shared task for which UD was used for additional training data in our investigation of the denoising task in the MTL setup. We list the specific treebanks used in order to encourage reproducibility. We preserve both the data and corpus information for the selected words. Specifically, we have also collected the token frequency, UPOS frequency, and character frequency for each of the additional data sampled, to be made available with the code for future analysis.

B Models and Experimental Details

B.1 Implementation

All models are implemented with a fork of yoyodyne, which is built over pytorch-lightning (Falcon and The PyTorch Lightning team, 2019). We utilize yoyodyne’s existing implementation of the Wu

¹<https://www.wiktionary.org/>

et al., 2021 models. We additionally implemented the CMLM objective, two stage training for PT setup, and the MTL setup including data and loss combination using the framework.

B.2 Compute and Infrastructure

For reproducibility, we utilize only Nvidia V100 GPUs for our experiments. The reported models together required \sim 180 hours of GPU time.

B.3 Reproducibility

In addition to using a consistent GPU architecture, we use a fixed random seed of 1 for all our model experiments. We also maintain copies of the specific data.

B.4 Morphological Inflection in Japanese

Organizers of the 2023 shared task note “the challenges that Japanese presents in morphologized inflection, namely due to its extremely large vocabulary size. In our work this ~~presented~~ as most models perform poorly on Japanese and do not meaningfully improve upon the baseline.”

C Significance Testing

In order to analyze the significance of our results, we perform a paired permutation test between test accuracies of all the models compared to the baseline. For all these tests we use the null-hypothesis that the mean difference between the test accuracies for these pairs is 0 and run the tests with 100k sampled permutations of the differences using SciPy (Virtanen et al., 2020).