

# Compare Results

Old File:

**2024.emnlp-main.1055.pdf**

16 pages (14.23 MB)

versus

New File:

**2024\_emnlp-main\_1055.pdf**

15 pages (252 KB)

2/8/2026 4:47:04 AM

**Total Changes**

**38**

**Content**

|    |              |
|----|--------------|
| 9  | Replacements |
| 14 | Insertions   |
| 15 | Deletions    |

**Styling and Annotations**

|   |             |
|---|-------------|
| 0 | Styling     |
| 0 | Annotations |

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



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

Abhishek Purushothama<sup>1</sup> Adam Wiemerslage<sup>2</sup>

Katharina von der Wense<sup>2,3</sup>

<sup>1</sup> Georgetown University

<sup>2</sup> University of Colorado Boulder

<sup>3</sup> Johannes Gutenberg University Mainz

[abhishek@cs.georgetown.edu](mailto:abhishek@cs.georgetown.edu)

## Abstract

Pretrained transformers such as BERT (Devlin et al., 2019) have been shown to be effective in many natural language tasks. However, they are underexplored 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 setups 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—not just from additional data. We also find that this does not hold true for specific combinations of secondary task and training setup, which has interesting implications for unsupervised training and denoising objectives in character- level tasks.

## 1 Introduction

Transformers have been shown to be an effective architecture for various natural language processing tasks (Vaswani et al., 2017), facilitating the ubiquitous method of pretraining on some unsupervised task with an abundance of data and then finetuning to a specific supervised task. Transformers have

also been shown to be an effective architecture for character- level tasks such as grapheme- to- phoneme conversion (G2P) and morphological inflection (Wu et al., 2021).

However, very little work has explored the application of pretrained models to character- level tasks, which likely require different inductive biases than the more semantically- oriented tasks where pretraining is typical. For instance, Xue et al. (2022, ByT5), a multilingual pretrained transformer using byte inputs, showed impressive performance on several semantically- oriented benchmarks, as well as on some character- level tasks including morphological inflection. However, it still under- performs the best two shared task submissions for the inflection benchmark (Vylomova et al., 2020).

The computational morphology community is frequently interested in low- resource languages - languages that do not have sufficient data available to apply standard NLP techniques. This is harder for morphologically complex languages, where the large set of inflectional patterns lead to an explosion in possible words, which become difficult to model with a small dataset. For these reasons, there is interest in building tools to aid in expanding morphological resources for language education tools, research, and documentation. Using NLP methods to build systems for analyzing and applying morphology in generalizable way to unseen words is thus a useful



goal. Several shared tasks have been held to this end (Cotterell et al., 2016, 2018; Vylomova et al., 2020; Pimentel et al., 2021; Kodner et al., 2022), where a machine learning model that performs well can be seen as competently representing the underlying system of morphology for a given language.

In this work, we explore utilizing secondary unsupervised tasks - tasks similar to language modeling which can serve as auxiliary tasks in a multitasking setup or pretraining tasks in a pretraining setup - when training encoder- decoder transformers for the task of morphological inflection. We investigate the benefits of pretraining (PT) beyond expanding the vocabulary distribution during training and also compare it to multi- task learning (MTL). Following Kann and Schütze (2017), we use autoencoding (AE) as an unsupervised secondary task and additionally compare it to the denoising task of character- level masked language modeling (CMLM) (Wiemerslage et al., 2023; Devlin et al., 2019). We explore these methods in data- scarce settings to investigate their potential impact in the low- resource setting. Our data samples and code are available publicly.<sup>1</sup>

We specifically investigate the following research questions:

- RQ1: Is training on secondary unsupervised tasks an effective method for low-resource inflection, even without introducing any new words to the dataset? This allows us to measure the impact that unsupervised tasks have on a model outside of the obvious benefit of increasing data diversity.
- RQ2: Are denoising tasks a better alternative to autoencoding for morphological inflection?
- RQ3: When training a model for the given target task, does multi-task learning outperform pre-training?

Our results show that both unsupervised PT and MTL are effective for morphological inflection, even with samples prepared exclusively from the supervised data itself. We find that simply autoencoding

the training words is more effective than CMLM in these data- scarce settings. Though the best method on average seems to be MTL with AE in our experiments, this is not consistent across every language. We also find that, in the MTL setup, CMLM actually performs worse than the baseline—though this is quickly reversed if we use out- of- distribution data for the secondary task.

## 2 Background Work

### 2.1 Character-level Sequence-to-Sequence Tasks

Character- level sequence- to- sequence tasks, sometimes referred to as character transduction tasks, are a special case of neural sequence- to- sequence learning problems that deal with approximately word-sized sequences. They are characterized by small vocabularies  $\Sigma^*$  and short source and target strings. Given source strings  $S \in \Sigma^*$ , target strings  $Y \in \Sigma^*$ , and optionally some features  $\tau$  to condition on, the goal of this task is to learn a mapping

$$f(S, \tau) \rightarrow Y \quad (1)$$

where  $f(\cdot)$  is typically parameterized by a neural network. In this work, we focus on morphological inflection: a character- level task where a particular  $s \in S$  is typically a lemma,  $t \in \tau$  is a bundle of tags specifying inflectional features, and  $y \in Y$  is a surface word of the lemma that expresses the specified morphological features, e.g.,:

$$f(\text{cry}, \text{PST}) \rightarrow \text{cried}$$

Morphological inflection is an active area of research in NLP. Many shared tasks in the computational morphology community (Cotterell et al., 2017; Goldman et al., 2023) have spurred progress on this task, which can be considered a good proxy for measuring the extent to which machine learning models can acquire the system of morphology in a language. Wu et al., 2021 trained a transformer (Vaswani et al., 2017) for several character- level transduction tasks resulting in state- of- the- art results. We follow

---

<sup>1</sup>1



their training methodology for inflection models as our baseline in this work.

## 2.2 Transfer Learning

Additional data for tasks different from the target task can be used to learn representations that benefit some target task via transfer learning. This often entails training on an unsupervised secondary task like language modeling, due to the large availability of unannotated text and the high cost of attaining annotations for specific target tasks. There has also been a great deal of research in transfer learning with supervised tasks (Bingel and Søgaard, 2017; Phang et al., 2018; Pruksachatkun et al., 2020).

We explore two different setups for this, both of which are unsupervised. Multi- task learning (Caruana, 1997, MTL) refers to training some task(s) together with the target task by including samples from both in a single training run and combining the loss from each (Luong et al., 2016). Intuitively, a well- chosen secondary task will benefit the target task by encouraging a model to learn a representation that minimizes the loss for both tasks simultaneously (Fifty et al., 2021). Pretraining (PT) refers to an alternative training setup in which models are first trained solely on secondary task(s) to encourage learning representations independent of the target task and then finetuned to some target task (Peters et al., 2018). Though both setups are similar, MTL relies on the joint optimization of multiple objectives, requiring a model to resolve all tasks at the same time. On the other hand, PT attempts to learn a representation that can be finetuned to a task later, by way of leveraging general encodings, or drawing upon an inductive bias learned in the pre- training phase.

## 2.3 Secondary Tasks

We also explore two secondary tasks: Autocoding (AE) is a simple and surprisingly effective method for representation learning. Here, an input is encoded with a model, and then decoded back to its original form. For word level tasks such as inflection, this

means sampling a word, and then simply predicting that same word, e.g.:

$$\text{tried} \rightarrow \text{tried}. \quad (2)$$

Denoising methods involve adding some noise to an input and then decoding the original form as it was before the noising step (Vincent et al., 2010), e.g.: given tried, we might have

$$\text{tr}@e@ \rightarrow \text{tried}, \quad (3)$$

where @ is a noise token that is applied in a data preprocessing step, and which the model must learn to replace with the original token. Many denoising strategies have been proposed for pretraining language models (Devlin et al., 2019; Raffel et al., 2019; Lewis et al., 2020), which may have advantages for particular downstream tasks.

## 2.4 Transfer Learning for Character-level Tasks

Kann and Schütze (2017) investigated the effectiveness of AE in an MTL setup by autoencoding with additional out- of- distribution words along with the target inflection task. Recently, Wiemerslage et al. (2023) pretrained various neural models on a character- level masked language modeling (CMLM) task, which follows the objective from Liu et al. (2019, RoBERTa), finding it can increase robustness to noise in the training data without the addition of new words. We follow them and use CMLM as the denoising task in our experiments. Similarly, Dong et al. (2022) pretrained a transformer encoder with a grapheme- based masking objective before finetuning to a downstream grapheme- to- phoneme (G2P) task and showed improvements for some datasets (Ashby et al., 2021).

## 2.5 Data Diversity and Multi-task Learning

The (word- level) token distribution for data in an MTL setup has been shown to have a strong impact on model performance (Martinez Alonso and Plank, 2017). In an exploration of supervised secondary

tasks, Bingel and Sogaard (2017) found that, when training with MTL for many NLP tasks, the out-of-vocabulary rate in the auxiliary task is positively associated with performance. This can also translate to unsupervised training for character-level tasks, where external data can positively impact model training regardless of the task for training on that data. Bjerva et al. (2019) perform MTL on many supervised tasks annotated for the same input examples. They train on the predictions for auxiliary tasks on the test set in a transductive learning setup, which increases performance. Krishna et al. (2023) found reusing downstream task data for unsupervised pretraining - which they refer to as self pretraining - to be an effective alternative to pretraining on external data. In experiments, they show that this often outperformed finetuning on off-the-shelf model that was pretrained on external data.

Similarly, in this work, we explore how data diversity impacts performance. That is, we compare secondary task words drawn from the target task to external data. This isolates secondary task impact from the effect of increased data diversity.

## 3 Architecture and Training

In this section we discuss our training methodology including architecture, training setups, and tasks.

### 3.1 Architecture

All of our experiments utilize the character encoder-decoder transformer from Wu et al. (2021). We use 4 encoder and 4 decoder layers, 4 attention heads, embedding size 256, and a feed-forward layer with hidden size 1024. We also follow their methodology for selection of the best checkpoint, where the highest accuracy on a validation set is selected out of 50 checkpoints. For all hyperparameters, refer to Wu et al. (2021).

## 3.2 Training Tasks

### 3.2.1 Morphological Inflection

In this work, morphological inflection is the only supervised task considered, and it is the target task for all experiments. We formulate the inflection task identically to prior work (Kann and Schütze, 2016; Wu et al., 2021).

### 3.2.2 CMLM

We follow Wiemerslage et al. (2023) in implementing CMLM for the denoising secondary task, where masking hyperparameters follow RoBERTa, though we increase the mask sampling rate. Specifically, we sample  $m = 20\%$  of all input characters for masking. Then, for each character, with probability  $p_m = 0.8$  we replace it with a special mask token, with probability  $p_r = 0.1$  we replace it with another character randomly sampled from the vocabulary, and with probability  $p_i = 0.1$  we leave the character unchanged.

### 3.2.3 AE

We additionally compare to autoencoding as a secondary task, in which we do no denoising at all: the source and target word are identical.

## 3.3 Training Setups

We compare three different training setups: supervised-only, pretrain-finetune (PT) and multitask learning (MTL).

### 3.3.1 Supervised- only

This is identical to the training setup from (Wu et al., 2021), where a model is trained only for the morphological inflection task. We follow them in training the model on the target-task data for 800 epochs and the best of 50 checkpoints by validation accuracy is chosen.

### 3.3.2 Pretrain- Finetune (PT)

We first pretrain an encoder- decoder model on an unsupervised secondary task and then train it on supervised data in a finetuning stage. We train the encoder- decoder fully in both the pretraining and finetuning stages. The finetuning stage is nearly identical to the supervised training setup, except we train from a pretrained checkpoint instead of training from scratch. We train both stages for 800 epochs. Since this is a two- stage setup, we apply model selection criteria twice. In the pretraining stage, the best checkpoint is chosen by minimizing evaluation loss on the secondary unsupervised task. This means that in the pretraining stage the model is motivated to learn representations over the character sequences from the vocabulary. The finetuning stage model selection remains identical to the supervised setup.

### 3.3.3 Multi- task Learning (MTL)

Similar to the setup in Kann and Schütze (2017), models are trained simultaneously for the target task and an unsupervised secondary task. We assign a fixed task weight factor  $\alpha$  for the unsupervised secondary task and  $\beta$  for the target inflection task. For all experiments, we set  $\alpha = 1$  and  $\beta = 1$  , and compute loss as the weighted sum of the two:

$$L(\theta) = \alpha \sum l_1(g(i), o) + \beta \sum l_2(f(s, t), y) \quad (4)$$

where  $f$  is the inflection task as in Section 2.1,  $g(I)$  is the unsupervised secondary task function,  $i \in I$  and  $o \in O$  are the unsupervised source and target, and  $l_1$  and  $l_2$  are loss functions for the for the two tasks, respectively. In initial experiments, we tried varying the tasks weights and found little impact on performance.

Although the training objective is to minimize  $L(\theta)$  , the best model is selected as in the previous setups with the best evaluation accuracy on the target task after training for 800 epochs. We added specific task identifiers (i.e., [TASK1], [TASK2]) to the input during training and inference. These identifiers are part of the input, however separated from the source (and

features) with a start token. This way the model can identify the relevant task for the sample.

## 4 Data

### 4.1 Target-task Data

Morphological inflection training data is sampled from the 2023 shared task on morphological inflection (Goldman et al., 2023). This supervised dataset consists of triples comprising (lemma, feature set, inflected form). It consists of 10k train samples and 1k each of development and test samples for 26 languages and an additional unvocalized variant (heb\\_unvoc) of Hebrew (heb). We differentiate Hebrew variants in our experiments and results, although we refer to it collectively as a language. In order to simulate a data-scarce setting, we randomly subsample the train split to 1k samples, as in the medium setting of the SIGMORPHON 2017 shared task (Cotterell et al., 2017). We also flatten the hierarchical features following most submissions to the 2023 shared task. This is performed by parsing the features during pre-processing and combining the multi-level features with special characters to make combined features. Consequently, our task data consists of the development and test splits and a subsampled 1k train split, all with flattened features. We inherit the fact that the shared task partitions lemmas between the 3 splits, which means all experiments require generalizing to unseen lemmas.

### 4.2 Extracted Data

We experiment with secondary-task data taken exclusively from the training data. That is, given a labeled triple from the supervised morphological inflection dataset like (debut,V;PRS;NOM(3,SG), debuts), we make two unsupervised training samples: debut → debut and debuts → debuts.

### 4.3 External Data

We perform an additional analysis with data sampled from a source external to the supervised data, which we refer to as external data. Here, we sample words



from the universal dependencies (UD) treebanks (Zeman et al., 2023). Since the availability of languages in UD does not directly correspond to the 2023 shared task data, we select 19 languages for which treebanks are available. The specific treebank used for dataset creation for each language is mentioned in Table 1. From each language’s treebank, we sample 2k words to use for secondary tasks. For details on how words are sampled, see appendix (Section A.3).

## 5 Experiments

### 5.1 Experimental Setup

We compare five model variants: baseline refers to the supervised model following Wu et al. (2021). We refer to PT-CMLM for models pretrained on the extracted data with the CMLM objective and then finetuned to the supervised data, whereas MTL-CMLM models train both tasks in MTL setup. PT-AE and MTL-AE reflect the same respective training setups, but use autoencoding as the secondary task. With these variants, we can compare all models to the baseline to answer RQ1, and we can compare across training setups and secondary tasks to answer RQ2 and RQ3, respectively.

### 5.2 Results

In Table 2 we present the main results: the accuracy of all five model variants averaged over all 27 languages on each of the development and test set. For a per-language results breakdown, see Table [MISSING]. For all comparisons, we focus on average accuracy on the test set.

The baseline is outperformed by almost all model variants that have been trained on secondary tasks. This means that secondary unsupervised tasks are beneficial even when no new data is introduced (RQ1). PT-CMLM outperforms the baseline by 1.84 absolute accuracy, only performing worse than the baseline on 6 languages: deu, ita, jpn, rus, sme, sqi. PT-AE performs even better, outperforming the baseline by 3.16 absolute accuracy, but performs worse than the baseline in 5 languages: bel, dan, jpn,

mkd, rus. We perform a paired permutation test and find all comparisons to the baseline to be statistically significant ( $p < 0.03$ ).

A comparison across unsupervised objectives shows that AE outperforms CMLM (RQ2). Although on average the difference is small (1.32) in the PT setup, AE outperforms CMLM substantially by 10.9 absolute accuracy in the MTL setup on the test set. Overall, MTL-AE is the best performing model, which indicates that MTL is a better setup for this task than PT (RQ3). However, this is not true when using the denoising objective. Only on 6 languages (dan, fra, heb, heb\_unvoc, klr, san) does MTL-CMLM outperform the baseline, and on average it performs worse than the baseline.

#### 5.2.1 Unsupervised Training on the Target-task Data

Most of the models outperform the baseline using strictly extracted finetuning data for unsupervised training with no additional words. This indicates that unsupervised tasks are effective for transfer learning in low-resource scenarios separately from the effect of exposing the model to new data. For PT, we hypothesize that the unsupervised training allows the model to better learn character-level representations before specializing to the inflection task.

## 6 Additional Analysis with External Data

Here all data for unsupervised learning is sampled from a source external to the finetuning data. We use Universal dependencies (Zeman et al., 2023, UD) as the source of external data, which we discuss in more detail in Subsection A.3.

Universal Dependencies Data All inflection task data (Subsection 4.1) is derived from the SIGMORPHON 2023 shared task, which samples its splits from UniMorph (Batsuren et al., 2022)—a type- level multilingual morphological resource for NLP, with labeled morphological paradigms comprising 182 languages, 122M inflections, and 769K derivations extracted semi- automatically. Universal Dependencies



is another multilingual NLP resource consisting of treebanks in 148 languages (as of the 2.13 release), though annotated data comprises token- level corpora. We choose UD as the source of external data in order to simulate a more naturally occurring type distribution than UniMorph. Whereas UniMorph types are likely to (i) be of the same part of speech as the test set, and (ii) represent interesting inflections that may be rare in a realistic low- resource scenario, UD contains types more representative of any arbitrary text. At the same time, unlike raw text scraped from the internet, UD data is relatively clean and has been vetted by experts, which ensures we do not experiment with e.g., data that has been misidentified as the target language or is otherwise contaminated.

Since not all 27 languages have treebanks in UD, we manually select a single treebank in only 19 of the 27 languages for these experiments. All models that use external data for secondary tasks are referred to with the suffix "UD".

## 6.1 Results

In Table 3, we present results for all 19 languages where MTL-CMLM-UD and MTL-AE-UD use external data sampled from UD for the respective secondary task. Using external data results in a 13.24 increase in absolute accuracy over MTL-CMLM, and outperforms the baseline substantially. On the other hand, the external data also leads to improved performance for MTL-AE-UD, but at a much smaller scale of 3.38 absolute accuracy over MTL-AE. On average, MTL-AE and MTL-CMLM-UD perform similarly. In a paired permutation test, all results have a statistically significant increase in performance over the baseline, except for MTL-CMLM which underperforms the baseline ( $p < 0.006$ ). We now focus on the substantial increase for MTL-CMLM-UD. This result supports the hypothesis that jointly optimizing a sufficiently different task from the target task, but on the same data causes issues. Consider the MTL-CMLM-UD model. The denoising task is learning representations over character sequences that are different from those in the target task, allowing the two tasks to update model parameters for separate distributions, and reducing conflicts in the joint- op-

timization. Indeed substituting the extracted data with external data when using the same denoising task leads to a remarkable improvement in performance.

## 6.2 Training Dynamics in MTL

We analyze the training dynamics between both the target and secondary task to further explain the MTL behavior. Bingel and Sogaard (2017) find that features of the learning curves are strong predictors of which secondary tasks lead to the best performance in an MTL setup. They hypothesize that MTL helps most in cases where a target task converges quickly, while the secondary task is still learning, which may help target tasks avoid getting stuck in local minima. We explore this hypothesis by, like them, looking at the gradients of each task’s training loss with respect to epochs, where the losses are recorded at the end of each epoch.

We then check the target task gradients that are  $\geq 0$  within the first 10% – 30% of training epochs, which we can consider to indicate that the task is plateauing early in training. In Figure 1 we provide violin plots of the secondary task gradients for those early target task plateaus in Sami—the language with the highest MTL improvement when UD data is added, and Danish—the language with the lowest improvement. For both languages, AE distributions have small variance around 0, whereas the CMLM plots show wider distributions. This reflects the fact that the CMLM loss is less stable, oscillating much more than the AE loss. More directly addressing the hypothesis about helping the target task recover from local minima, we see distributions that are either top- heavy, or normal for Danish, where no secondary task leads to a very large increase in performance over the baseline. On the other hand, the CMLM-UD distribution is more bottom- heavy for Sami, indicating that there are more negative gradients, and thus more epochs where the model is still learning this task when the target task seems to plateau. The AE distribution, while still low variance around 0, also have lower negative gradients compared to Danish.

This small analysis suggests two things. First, we



IMAGE NOT PROVIDED

Figure 1: The distribution of secondary task gradients between 20% and 30% training as in Bingel and Sogaard (2017) for cases in which the target task gradients are  $\geq 0$ . A negative number indicates the model is still improving upon the secondary task.

have weak support for the hypothesis that MTL helps when the secondary task continues to converge when the target task plateaus early. We see more negative values in the Sami distribution where MTL is more helpful, especially in the CMLM-UD secondary task when compared to the CMLM without UD data. Second, AE, typically the best secondary task in our experiments, appears to have a lower variance in gradients, indicating that the training loss is more stable. Indeed, the variance for CMLM gradients is larger in Sami, where CMLM hurts performance, and the variance is smaller in Sami when we add the UD data, which has a large positive impact.

## 7 Conclusion

In this work, we explored multiple methods for transfer learning for morphological inflection, many of which showed remarkable performance for a large set of languages. We investigated two different training methods: pretraining- finetuning and multi- task learning, and two different secondary tasks: denoising and autoencoding. In a low- resource setting, we found that secondary unsupervised tasks are effective even without the addition of any new vocabulary items beyond the finetuning dataset. While pretraining is an effective setup for improving morphological inflection without any external data, multi- task learning with an autoencoding objective is the best setup in all experiments. On the other hand, multi- task learning with the CMLM denoising objective is the worst performing setup, performing below the baseline on average. In further analysis, we found that performing CMLM on external data that is separate from the finetuning data solves this issue, resulting in significantly better performance.

The success of denoising objectives such as MLM cannot be denied for large- scale training and semantically oriented tasks. Our experiments and results show that similar tasks are effective in data- scarce settings for character- level tasks like morphological inflection. In practice, it seems that low- resource character- level tasks should always consider training in a multi- task setup with an autoencoding secondary task even if the supervised training data is the only available data - and exploring denoising objectives if unsupervised data from an external source is available.

## 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.

## 9 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.

## 10 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



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

| ISO-639-2 | Language            | UD Treebank used      |
|-----------|---------------------|-----------------------|
| afb       | Arabic, Gulf        | Arabic-PADT           |
| amh       | Amharic             | Amharic-ATT           |
| arz       | Arabic, Egyptian    | [MISSING]             |
| bel       | Belarusian          | Belarusian-HSE        |
| dan       | Danish              | Danish-DDT            |
| deu       | German              | German-GSD            |
| eng       | English             | English-Atis          |
| fin       | Finnish             | Finnish-FTB           |
| fra       | French              | French-GSD            |
| grc       | Ancient Greek       | Ancient_Greek-Perseus |
| heb       | Hebrew              | Hebrew-HTB            |
| heb.unvoc | Hebrew, Unvocalized | [MISSING]             |
| hun       | Hungarian           | Hungarian-Szeged      |
| hye       | Eastern Armenian    | Armenian-ArmTDP       |
| ita       | Italian             | Italian-ISDT          |
| jpn       | Japanese            | Japanese-GSD          |
| kat       | Georgian            | [MISSING]             |
| klr       | Khaling             | [MISSING]             |
| mkd       | Macedonian          | [MISSING]             |
| nav       | Navajo              | [MISSING]             |
| rus       | Russian             | Russian-GSD           |
| san       | Sanskrit            | Sanskrit-UFAL         |
| sme       | Sami North          | North_Sami-Giella     |
| spa       | Spanish             | Spanish-AnCora        |
| sqi       | Albanian            | [MISSING]             |
| swa       | Swahili             | [MISSING]             |
| tur       | Turkish             | Turkish-Atis          |

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] Jacob Devlin, Ming- Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre- training

of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4171–4186, 2019.

- [2] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [3] Shijie Wu, Ryan Cotterell, and Mans Hulden. Applying the transformer to character-level transduction. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 1901–1907, 2021.
- [4] Linting Xue, Aditya Barua, Noah Constant, Rami Al-Rfou, Sharan Narang, Mihir Kale, Adam Roberts, and Colin Raffel. ByT5: Towards a token-free future with pre-trained byte-to-byte models. *Transactions of the Association for Computational Linguistics*, 10:291–306, 2022.
- [5] Ekaterina Vylomova, Jennifer White, Elizabeth Salesky, Sabrina J. Mielke, Shijie Wu, Edoardo Maria Ponti, Rowan Hall Maudslay, Ran Zmigrod, Josef Valvoda, Svetlana Toldova, Francis Tyers, Elena Klyachko, Ilya Yegorov, Natalia Krizhanovsky, Paula Czarnowska, Irene Nikkarinen, Andrew Krizhanovsky, Tiago Pimentel, Lucas Torroba Hennigen, Christo Kirov, Garrett Nicolai, Adina Williams, Antonios Anastasopoulos, Hilaria Cruz, Eleanor Chodroff, Ryan Cotterell, Miikka Silfverberg, and Mans Hulden. 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, 2020.

- [6] Ryan Cotterell, Christo Kirov, John Sylak-Glassman, David Yarowsky, Jason Eisner, and Mans Hulden. The SIGMORPHON 2016 shared Task- Morphological reinfection. In *Proceedings of the 14th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 10–22, 2016.
- [7] Ryan Cotterell, Christo Kirov, John Sylak-Glassman, Geraldine Walther, Ekaterina Vylomova, Arya D. McCarthy, Katharina Kann, Sabrina J. Mielke, Garrett Nicolai, Miikka Silfverberg, David Yarowsky, Jason Eisner, and Mans Hulden. The CoNLL- SIGMORPHON 2018 shared task: Universal morphological reinfection. In *Proceedings of the CoNLL- SIGMORPHON 2018 Shared Task: Universal Morphological Reinfection*, pages 1–27, 2018.
- [8] Tiago Pimentel, Maria Ryskina, Sabrina J. Mielke, Shijie Wu, Eleanor Chodroff, Brian Leonard, Garrett Nicolai, Yustinus Ghanngo Ate, Salam Khalifa, Nizar Habash, Charbel El- Khaissi, Omer Goldman, Michael Gasser, William Lane, Matt Coler, Arturo Oncevay, Jaime Rafael Montoya Samame, Gema Celeste Silva Villegas, Adam Ek, Jean- Philippe Bernardy, Andrey Shcherbakov, Aziyana Bayyrpool, Karina Sheifer, Sofya Ganieva, Matvey Plugarov, Elena Klyachko, Ali Salehi, Andrew Krizhanovsky, Natalia Krizhanovsky, Clara Vania, Sardana Ivanova, Aelita Salchak, Christopher Straughn, Zoey Liu, Jonathan North Washington, Duygu Ataman, Witold Kieraś, Marcin Woliński, Totok Suhardijanto, Niklas Stoehr, Zahroh Nuriah, Shyam Ratan, Francis M. Tyers, Edoardo M. Ponti, Grant Aiton, Richard J. Hatcher, Emily Prud'hommeaux, Ritesh Kumar, Mans Hulden, Botond Barta, Dorina Lakatos, Gábor Szolnok, Judit Ács, Mohit Raj, David Yarowsky, Ryan Cotterell, Ben Ambridge, and Ekaterina Vylomova. SIGMORPHON 2021 shared task on morphological reinfection: Generalization across languages. In *Proceedings of the 18th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 229–259, 2021.
- [9] Jordan Kodner, Salam Khalifa, Khuyagbaatar Batsuren, Hossep Dolatian, Ryan Cotterell, Faruk Akkus, Antonios Anastasopoulos, Taras Andrushko, Aryaman Arora, Nona Atanalov, Gabor Bella, Elena Budianskaya, Yustinus Ghanngo Ate, Omer Goldman, David Guriel, Simon Guriel, Silvia Guriel- Agiashvili, Witold Kieraś, 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. 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, 2022.
- [10] Katharina Kann and Hinrich Schütze. 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, 2017.
- [11] Adam Wiemerslage, Changbing Yang, Garrett Nicolai, Miikka Silfverberg, and Katharina Kann. An investigation of noise in morphological inflection. In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 3351–3365, 2023.
- [12] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. arXiv preprint arXiv:1907.11692, 2019.
- [13] Joachim Bingel and Anders Sogaard. Identifying beneficial task relations for multi- task learning in deep neural networks. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 164–169, 2017.



- [14] Jason Phang, Thibault Fevry, and Samuel R Bowman. Sentence encoders on stilts: Supplementary training on intermediate labeled- data tasks. arXiv preprint arXiv:1811.01088, 2018.
- [15] Yada Pruksachatkun, Jason Phang, Haokun Liu, Phu Mon Htut, Xiaoyi Zhang, Richard Yuanzhe Pang, Clara Vania, Katharina Kann, and Samuel Bowman. Intermediate- task transfer learning with pretrained language models: When and why does it work? In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5231–5247, 2020.
- [16] Rich Caruana. Multitask learning. *Machine Learning*, 28(1):41–75, 1997.
- [17] Thang Luong, Quoc V. Le, Ilya Sutskever, Oriol Vinyals, and Lukasz Kaiser. Multi- task sequence to sequence learning. In *International Conference on Learning Representations*, 2016.
- [18] Christopher Fifty, Ehsan Amid, Zhe Zhao, Tianhe Yu, Rohan Anil, and Chelsea Finn. Efficiently identifying task groupings for multi- task learning. In *Neural Information Processing Systems*, 2021.
- [19] Matthew E. Peters, Mark Neumann, Mohit Iyyer, Matt Gardner, Christopher Clark, Kenton Lee, and Luke Zettlemoyer. Deep contextualized word representations. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2227–2237, 2018.
- [20] Pascal Vincent, H. Larochelle, Isabelle Lajoie, Yoshua Bengio, and Pierre- Antoine Manzagol. Stacked denoising autoencoders: Learning useful representations in a deep network with a local denoising criterion. *Journal of Machine Learning Research*, 11:3371–3408, 2010.
- [21] Colin Raffel, Noam M. Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text- to- text transformer. *Journal of Machine Learning Research*, 21(140):1–67, 2019.
- [22] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 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, 2020.
- [23] Lu Dong, Zhi- Qiang Guo, Chao- Hong Tan, Ya- Jun Hu, Yuan Jiang, and Zhen- Hua Ling. Neural grapheme- to- phoneme conversion with pre- trained grapheme models. In *ICASSP 2022 - 2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 6202–6206, 2022.
- [24] Lucas F.E. Ashby, Travis M. Bartley, Simon Clematide, Luca Del Signore, Cameron Gibson, Kyle Gorman, Yeonju Lee- Sikka, Peter Makarov, Aidan Malanoski, Sean Miller, Omar Ortiz, Reuben Raff, Arundhati Sengupta, Bora Seo, Yulia Spektor, and Winnie Yan. Results of the second SIGMORPHON shared task on multilingual grapheme- to- phoneme conversion. In *Proceedings of the 18th SIGMORPHON Workshop on Computational Research in Phonetics, Phonology, and Morphology*, pages 115–125, 2021.
- [25] Hector Martinez Alonso and Barbara Plank. 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, 2017.
- [26] Johannes Bjerva, Katharina Kann, and Isabelle Augenstein. Transductive auxiliary task self- training for neural multi- task models. In *Proceedings of the 2nd Workshop on Deep Learning Approaches for Low- Resource NLP (DeepLo 2019)*, pages 253–258, 2019.



- [27] Kundan Krishna, Saurabh Garg, Jeffrey Bigham, and Zachary Lipton. 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 12207–12222, 2023.
- [28] Omer Goldman, Khuyaqbaatar Batsuren, Salam Khalifa, Aryaman Arora, Garrett Nicolai, Reut Tsarfaty, and Ekaterina Vylomova. 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, 2023.
- [29] Khuyaqbaatar Batsuren, Omer Goldman, Salam Khalifa, Nizar Habash, Witold Kieraś, Gábor Bella, Brian Leonard, Garrett Nicolai, Kyle Gorman, Yustinus Ghanggo Ate, Maria Ryskina, Sabrina Mielke, Elena Budianskaya, Charbel El-Khaissi, Tiago Pimentel, Michael Gasser, William Abbott Lane, Mohit Raj, Matt Coler, Jaime Rafael Montoya Samame, Delio Siticon-tazi Camaiteri, Esaú Zumaeta Rojas, Didier López Francis, Arturo Oncevay, Juan López Bautista, Gema Celeste Silva Villegas, Lucas Torroba Hennigen, Adam Ek, David Gurie, Peter Dirix, Jean-Philippe Bernardy, Andrey Scherbakov, Aziyana Bayyr-ool, Antonios Anatasopoulos, Roberto Zariquey, Karina Sheifer, Sofya Ganieva, Hilaria Cruz, Ritván Karahóga, Stella Markantonatou, George Pavlidis, Matvey Plugaryov, 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 Czarnowska, Irene Nikkarinen, Aelita Salchak, Brijesh Bhatt, Christopher Straughn, Zoey Liu, Jonathan North Washington, Yuval Pinter, Duygu Ataman, Marcin Wolinski, Totok Suhardijanto, Anna Yablonskaya, Niklas Stoehr, Hossep Dolatian, Zahroh Nuriah, Shyam Ratan, Francis M. Tyers, Edoardo M. Ponti, Grant Aiton, Aryaman Arora, Richard J. Hatcher, Ritesh Kumar, Jeremiah Young, Daria Rodionova, Anastasia Yemelina, Taras Andrushko, Igor Marchenko, Polina Mashkovtseva, Alexandra Serova, Emily Prud'hommeaux, Maria Nepomniashchaya, Fausto Giunchiglia, Eleanor Chodroff, Mans Hulden, Miikka Silfverberg, Arya D. McCarthy, David Yarowsky, Ryan Cotterell, Reut Tsarfaty, and Ekaterina Vylomova. UniMorph 4.0: Universal Morphology. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 840–855, 2022.
- [30] Daniel Zeman, Joakim Nivre, Mitchell Abrams, et al. Universal dependencies 2.12. LINDAT/CLARIAH-CZ digital library at the Institute of Formal and Applied Linguistics (ÚFAL), Faculty of Mathematics and Physics, Charles University, 2023.
- [31] Katharina Kann and Hinrich Schütze. Single-model encoder- decoder with explicit morphological representation for reinfection. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 555–560, 2016.
- [32] Ryan Cotterell, Christo Kirov, John Sylak-Glassman, Geraldine Walther, Ekaterina Vylomova, Patrick Xia, Manaal Faruqui, Sandra Kubler, David Yarowsky, Jason Eisner, and Mans Hulden. CoNLL- SIGMORPHON 2017 shared task: Universal morphological reinfection in 52 languages. In *Proceedings of the CoNLL SIGMORPHON 2017 Shared Task: Universal Morphological Reinfection*, pages 1–30, 2017.
- [33] Jordan Kodner, Sarah Payne, Salam Khalifa, and Zoey Liu. 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, 2023.



- [34] Saliha Muradoglu and Mans Hulden. Eeny, meeny, miny, moe. how to choose data for morphological inflection. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7294–7303, 2022.
- [35] Joakim Nivre, Marie- Catherine de Marneffe, Filip Ginter, Jan Hajic, Christopher D. Manning, Sampo Pyysalo, Sebastian Schuster, Francis Tyers, and Daniel Zeman. Universal Dependencies v2: An evergrowing multilingual treebank collection. In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 4034–4043, 2020.
- [36] William Falcon and The PyTorch Lightning team. PyTorch Lightning. 2019.
- [37] Pauli Virtanen, Ralf Gommers, Travis E. Oliphant, Matt Haberland, Tyler Reddy, David Cournapeau, Evgeni Burovski, Pearu Peterson, Warren Weckesser, Jonathan Bright, Stefan 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, Ilhan 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, Antonio H. Ribeiro, Fabian Pedregosa, Paul van Mulbregt, and SciPy 1.0 Contributors. SciPy 1.0: Fundamental Algorithms for Scientific Computing in Python. *Nature Methods*, 17:261–272, 2020.
- [38] Sina Ahmadi and Aso Mahmudi. 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, 2023.
- [39] Christo Kirov, John Sylak- Glassman, Roger Que, and David Yarowsky. 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, 2016.

## A Data details

### A.1 Limitations of UniMorph and SIGMORPHON

The unimorph project is the primary source for the dataset. It draws heavily from Wiktionary<sup>2</sup> in a semi-automated way based on Kirov et al. (2016). Wiktionary is a collaboratively built resource which, despite processes to promote accuracy, 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 Mahmudi (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 Selection and Sampling

Many features of morphological inflection data, such as overlap and frequency, have been shown to be important factors for model performance (Kodner et al., 2023). (Muradoglu 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 Goldman 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

---

<sup>22</sup>



is maintained in our datasets, and is thus not a relevant point of analysis as in other work (e.g. Kodner 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 19 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<sup>3</sup>, which is built over pytorch- lightning (Falcon and The PyTorch Lightning team, 2019). We utilize yoyodyne’s existing implementation of the Wu 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 to-

---

<sup>3</sup>gether 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 morphological inflection, namely due to its extremely large vocabulary size. In our work this persists 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).

