# How Low is Too Low?
# A Computational Perspective on Extremely Low-Resource Languages

**Rachit Bansal**[1]     **Himanshu Choudhary**[1]     **Ravneet Punia**[1]
**Niko Schenk**[2†]     **Jacob L Dahl**[3]     **Émilie Pagé-Perron**[3]
[1] Delhi Technological University     [2] Amazon Berlin, Germany     [3] University of Oxford
{rachitbansal2500, himanshu.dce12, ravneet.dtu}@gmail.com
nikosch@amazon.com
{jacob.dahl, emilie.page-perron}@wolfson.ox.ac.uk

## Abstract

Despite the recent advancements of attention-based deep learning architectures across a majority of Natural Language Processing tasks, their application remains limited in a low-resource setting because of a lack of pre-trained models for such languages. In this study, we make the first attempt to investigate the challenges of adapting these techniques for an extremely low-resource language – Sumerian cuneiform – one of the world's oldest written languages attested from at least the beginning of the 3rd millennium BC. Specifically, we introduce the first cross-lingual information extraction pipeline for Sumerian, which includes part-of-speech tagging, named entity recognition, and machine translation. We further curate *InterpretLR*, an interpretability toolkit for low-resource NLP, and use it alongside human attributions to make sense of the models. We emphasize on human evaluations to gauge all our techniques. Notably, most components of our pipeline can be generalised to any other language to obtain an interpretable execution of the techniques, especially in a low-resource setting. We publicly release all software, model checkpoints, and a novel dataset with domain-specific preprocessing to promote further research.

## 1 Introduction

Sumerian is one of the oldest written languages, attested in the cuneiform texts from around 2900 BC and possibly the language of even older proto-cuneiform texts from the second half of 4th millennium BC (Englund, 2009). Specialists in Assyriology have recently worked to digitize Sumerian scripts, annotate, and translate a part of them to modern-day languages like English and German.

---

Datasets and training subroutines are available at
linktr.ee/rachitbansal

†Work was done prior to joining Amazon at Goethe University Frankfurt



obverse.
1. `1(disz) kusz udu niga`
1 hide, grain-fed sheep;

2. `1(disz) kusz masz2 niga`
1 hide, grain-fed goat;

3. `kusz udu sa2-du11`
sheep hides, regular offerings,

4. `ki {d}iszkur-illat-ta`
from Adda-illat,

reverse.
1. `a-na-ah-i3-li2`
Anah-ili;

2. `szu ba-an-ti`
did receive.

3. `iti ezem-an-na`
Month: An-festival,

4. `mu na-ru2-a-mah`
`mu-ne-du3`
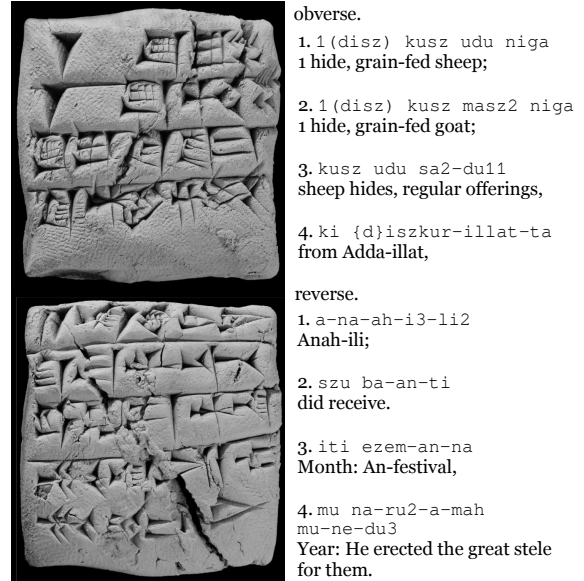Year: He erected the great stele for them.

Figure 1: Tablets inscribed with Sumerian cuneiform script, their corresponding digitized transliterations, and human-translated English text for each line.

In this work, we attempt to create the first information extraction and translation pipeline for Sumerian. Specifically, we focus on machine translation from Sumerian to English, and sequence labeling tasks of Named Entity Recognition (NER) and Part of Speech Tagging (POS).

Figure 1 shows a sample of our raw data where the Sumerian text has been derived from the tablet-inscribed cuneiform script along with its human-interpreted English translations. Creating an annotated corpus for such a language is a tedious task. Thankfully, we obtain our data from openly available sources and corpora, painstakingly annotated and translated by human experts. Yet, for languages like Sumerian, which are not fully-understood by humans themselves, transferring knowledge and patterns to learning algorithms from this limited data becomes extremely difficult. The consequent challenge posed for NER and POS is evident. Lack of annotated data and fuzzy character-level text

makes it hard for a model to generalise, irrespective of its size.

In case of machine translation, the labeled data is composed of incomplete and short phrase-like sentences, specially on the target-side. This makes the context largely ambiguous. Moreover, we find that for a majority of medieval languages the target-side translated text is highly incoherent with modern-day English language text, making it impossible to use the latter in semi-supervised and unsupervised settings.

Throughout this study, we elaborate on such challenges induced when working with low-resource languages, and talk about what makes some of these languages like Sumerian 'extremely' low-resource. Through extensive experimentation, evaluation, and analysis we further introduce specific algorithms and modifications to work around them.

In all, our contribution is three-fold:

1. Building and analyzing a variety of algorithms on the unexplored human-annotated Sumerian dataset for sequence labeling tasks of POS Tagging and NER. (§3)
2. Introducing the problem of *Target-side Incoherence* for low-resource settings and its effect on semi-supervised and unsupervised machine translation (§4.2). Further investigating specific modifications and methodologies to cope-up with these constraints. (§4)
3. Introducing *InterpretLR*, a generalisable toolkit to interpret low-resource NLP. We apply to additionally study, compare, and evaluate all of the proposed techniques for machine translation and sequence labeling. (§7)

Throughout this work, we have conducted human studies and evaluation for our models, in addition to automated metrics. For studying our models with *InterpretLR*, we have made use of human annotations.

## 2 Background

### 2.1 Data

Sumerian is an ancient language from Iraq that was written using the cuneiform script. While Basque and Turkish display some similarities (split-ergativity, agglutinativity), it is a language isolate (Englund, 2009). We have found artifacts dating to around 2900 BC with Sumerian texts inscribed until the first century AD. Most of the Sumerian texts found to this day are administrative in nature

as, during the third dynasty of the Ur III Period, the state administration swell to an unprecedented level of activity which was not seen again later in the history of Mesopotamian culture. All through this study, our evaluation sets are composed of Ur III Admin text only and it acts as our in-domain data.

Part of the datasets we used were assembled from the Cuneiform Digital Library Initiative (CDLI)[1], Machine Translation and Automated Analysis of Cuneiform languages (MTAAC) project (Pagé-Perron et al., 2017)[2] and The Electronic Text Corpus of Sumerian Literature (ETCSL) dataset [3]. CDLI and MTAAC datasets contain the Ur III Administrative (Admin) texts[4] which are preserved by the CDLI[5]. The MTAAC and ETCSL corpora were both manually annotated for morphology by cuneiform linguistics.

We divided the data between training and testing sets, and then to reduce the data sparsity, we performed text augmentation using a set of labeled named entities for these sets separately. This increased our combined number of phrases from $25,000$ to $48,000$, representing our final dataset for sequence labeling. Figures 2 and 3 provide the distribution of word tokens in our final pre-annotated dataset. The corpus consists of phrases with lengths ranging from 1 to 19 words. These phrases are small since they are translated line by line from the scripts. Around $2,500$ phrases were used for testing, while the $45,500$ were employed for training purposes.

For machine translation, the final dataset summarizes as (i) $10,520$ parallel phrases from the Ur III administrative corpus; (ii) $88,460$ parallel phrases, all genres combined; and (iii) all monolingual Sumerian data ($1.43$ million phrases). In all cases, phrases are short, generally ranging from 1 to 5-word tokens.

### 2.2 Related Work

Past work aimed at machine translation of Sumerian-English (Pagé-Perron et al., 2017; Punia et al., 2020) have used the minimal bitext upon a variety of general statistical and neural supervised

---

[1] https://cdli.ucla.edu
[2] https://cdli-gh.github.io/mtaac/
[3] http://http://etcsl.orinst.ox.ac.uk/
[4] The Third Dynasty of Ur is a cultural and temporal period ranging in $\sim 2112 - 2004$ BC, in Mesopotamia
[5] https://github.com/cdli-gh/data, https://github.com/cdli-gh/mtaac_gold_corpus/tree/workflow/morph/to_dict

techniques. However, they do not handle the text-level peculiarities any differently than one would do for a high-resource language, thus, often failing to capture context, resulting in poor and inconsistent translations. Techniques, learning algorithms, and architectures that optimally use the vast monolingual data and parallel sentences while keeping in mind the several linguistic limitations are motivated in such a scenario. Thus, we experiment on semi-supervised and unsupervised techniques across the three categories of Data Augmentation (Sennrich et al., 2016; He et al., 2016), Knowledge Transfer (Zoph et al., 2016), and Pre-training (Conneau and Lample, 2019; Song et al., 2019).

In the past, Pagé-Perron et al. (2017) applied statistical models for morphological analysis and information extraction for Sumerian. Although, due to the unavailability of annotated data, these models could not generalise well. Liu et al. (2015) and Luo et al. (2015) used an unsupervised approach for NER with the help of domain experts and used Contextual and Spelling rules to build the model. They also post-processed their outputs automatically, which enhanced their results. In this work, we thoroughly investigate a wide range of algorithms for these sequence labeling tasks and consequently take a first step towards effective information extraction for Sumerian.
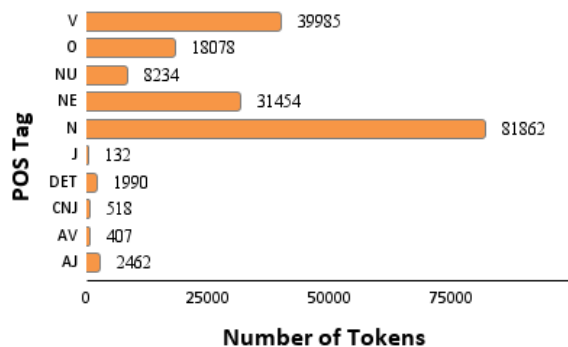


Figure 2: Composition of the POS tagging dataset. Here, "NE" stands for named entities, "O" stands for unstructured words. Other tags are in accordance with ORACC.

## 3 Part of Speech Tagging and Named Entity Recognition

In this section, we talk about the various algorithms that we investigated to carry out the sequence labeling tasks of POS and NER for Sumerian. The subsequent experimental results are described and discussed in Section 6.
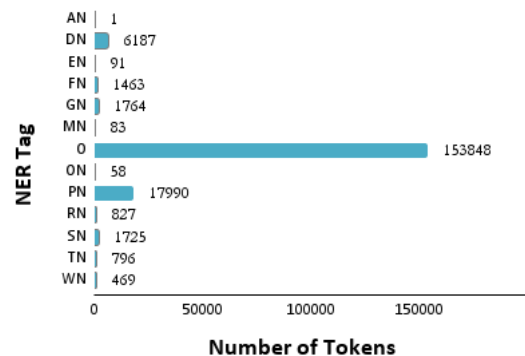


Figure 3: Composition of our NER dataset. Tags are in accordance with ORACC.

**Conditional Random Fields** CRF (Lafferty et al., 2001) is a discriminative probabilistic classifier, which optimises the weights or parameters in order to maximize the conditional probability distribution $P(y \mid x)$. They take set of input features (language or domain specific) into account, using the learned weights associated with these features and previous labels to predict the current label. Since CRFs use feature sets (rules) which are language-specific, it makes the model more robust specially for very low-resource languages. In our case we developed domain specific rules with the help of previous studies (Liu et al., 2015; Luo et al., 2015) and language experts. A set of these rules are mentioned in the Appendix.

**Bi-directional LSTM** We also experiment across Recurrent Neural Networks (RNNs) to deal with the sequential text input. We employ Bi-LSTM (Hochreiter and Schmidhuber, 1997; Schuster and Paliwal, 1997) in particular. As in Huang et al. (2015), an additional CRF layer is used for efficient usage of sentence level tag information and past input features by LSTM cells.

**FLAIR** Akbik et al. (2018) introduced a Contextual String Embedding for Sequence Labeling, FLAIR, which has shown great success for various languages for NER (Akbik et al., 2019b). We make use of the two distinct properties of its embeddings: (i) training without any explicit notion of words and fundamentally modeling the words as a sequence of characters, and (ii) deriving and using the context from surrounding tokens. We train the bi-directional character language model using the Sumerian monolingual phrases and retrieve the contextual embedding for each word which we then pass into the vanilla Bi-LSTM CRF model.

**RoBERTa** We also investigate the transformer-based language model, RoBERTa (Liu et al., 2019). The encoder is first pre-trained on our Sumerian monolingual data, and then fine-tuned on our downstream sequence labeling tasks using the labeled data.

## 4 Machine Translation

In this section, we present our experiments for machine translation, primarily focusing on specific data and algorithmic modeling techniques which may be generalised for any extremely low-resource language that may or may not suffer from *Target-side Incoherence*, a phenomenon which we also introduce herein. All results are summarised in Table 1.

### 4.1 Supervised NMT

In order to create a benchmark for the semi-supervised and unsupervised approaches, we perform supervised machine translation using the limited bitext available ($\sim$10,000 phrases). We perform experiments on a variety of data configurations which are given by:

1. `UrIIISeg`: Follows the format as present in the original texts provided by Assyriologists and used in the past attempts for Sumerian-English machine translation (Pagé-Perron et al., 2017). It contains only the in-domain Ur III Admin Data with line-by-line translated segments of 1-5 words each, amounting to 10528 segments.
2. `UrIIIComp`: Also contains the in-domain data only, but multiple segments are concatenated together to form complete sentences. The 'completeness' of a sentence is ensured using punctuation marks. It comprises of only 4792 sentences.
3. `AllSeg`: Contains out-of-domain Sumerian text segments in addition to in-domain Ur III Admin alone. The additional text varies across a wide range of genres such as literary, lexical, ritual and legal, resulting into a corpus size of 88466 segments.
4. `AllComp`: Combines the additional features of 2. and 3., thus comprising of a total of 32694 complete text sentences from all genres.

We make use of the vanilla transformer encoder and decoder architecture (Vaswani et al., 2017) for all our supervised machine translation experiments

over these three different bitext configurations. Noting the results as in 1, the `AllComp` text configuration is used for all other experiments. The computational configurations are mentioned in Section 5.

### 4.2 Semi-Supervised and Unsupervised NMT

We observed that one of the primary reasons for the lack of success of semi-supervised and unsupervised algorithms for low-resource settings, specially for medieval languages, is *the lack of coherence between monolingual texts in the modern-day corpora to the target-side text in the available parallel corpora*. We refer to this as the **Target-side Incoherence** (*TSIC*) problem for such languages.

As can be seen from Figure 1, the transliterated English text in our parallel corpora is vastly different from general modern-day English texts. In Sumerian, this is because the text has been human-translated to English on the level of words and small segments due to insufficient knowledge of the language. This results into a contextually distorted English language text, as compared what we see in general corpora. This leads to multiple pitfalls. Most significantly, the colossal monolingual data available for a data-rich target-side language (i.e., English in this case) can no longer be used. This *Target-side Incoherence* holds true for most medieval language texts like Sumerian, which makes them 'extremely' low-resource.

In this section, we elaborate on the problems caused due to *TSIC* and further present our hypothesis on adapting various semi-supervised and unsupervised NMT techniques to deal with them.

**Forward Translation** Back-translation (BT) (Sennrich et al., 2016) has been widely used and analysed for NMT across a large set of language pairs. BT uses a reverse model, Sumerian $\leftarrow$ English trained on the existing parallel corpora, when the task is to translate from Sumerian $\rightarrow$ English, and applies it on the target-side monolingual corpus. The synthetic samples thus generated are added to the source-side corpus and a new reverse model is trained on the augmented dataset. It has been shown to outperform its forward counterpart, Forward Translation (FT) (Zhang and Zong, 2016; Burlot and Yvon, 2018), which instead uses a forward (Sumerian $\rightarrow$ English) model to augment

the target-side of the bitext.

However, due to *TSIC*, the target-side monolingual data falls into a completely different distribution than what a Sumerian ← English model is trained on. Using back-translation in such a scenario results into a poor source-side augmentation, doing more harm than good. Keeping this in mind, we rely on forward-translation (FT), thus using the Sumerian monolingual text.

We divide the Sumerian monolingual data into 8 shards, each containing ∼100,000 monolingual `AllComp` sentences each. The FT process takes place for each shard and the Transformer model is trained after each shard is forward-translated.

Large scale studies (Edunov et al., 2018; Wu et al., 2019) have shown the heavy dependency of BT and FT on aspects like sampling methods and the amount of parallel data. The performance with non-MAP (where, MAP stands for *maximum a posteriori*) estimation methods like Nuclear Sampling (Holtzman et al., 2018) and Beam Search with Noise improves almost-linearly with the amount of bitext, and thus, for low-resource settings (∼80,000 sentence pairs), MAP methods have been shown to give better results. This was also observed in our experiments and the reported results are obtained using Beam Search (§5).

**Cross-Lingual Language Model Pre-training** We further make use of XLM (Conneau and Lample, 2019) to carry out a wide range of experiments for both unsupervised and semi-supervised fine-tuning techniques. Considering the lack of original target monolingual text due to *TSIC*, the following target data configurations were used for pre-training the XLM:

1. `WMT`: Ignoring the lack of coherence between general English texts and the evaluation + training texts, to compose the entire target monolingual data with the WMT '18 English Texts. Amounts to a total of 20M sentences.
2. `Orig`: Composed of all the English side texts in `UrIIISeg`, `UrIIIComp`, `AllSeg` and `AllComp` combined. Contains only ∼60,000 sentences.
3. `Mixed`: This consists of all of 2. and as many sentences as 1. through which the net size of the corpus equalizes the Sumerian monolingual, i.e., 1.5M.

In the pre-training phase, we perform various experiments over different combinations of

MLM and TLM. It is further fine-tuned on a denoising auto-encoding objective for Unsupervised while cross-reference machine translation objective over the parallel data for semi-supervised training. BT steps are also performed in both cases.

**Data Augmentation** In order to further reduce the effect of *TSIC* on the model performance and to allow the model to attend to a larger and more diverse volume of target text during pre-training, we make use of the following data augmentation techniques:

1. `BERT`: Replacing words by the spatially closest words measured by Cosine Similarity in BERT (Devlin et al., 2019) Embeddings, with a threshold of 0.8.
2. `WordNet`: Replacing words with WordNet (Miller et al., 1990) synonyms.
3. `CharSwap`: Introduces certain character-level perturbations in the text by substituting, deleting, inserting, and swapping adjacent character tokens.

Different combinations of these techniques have been used to augment the `Orig` type target monolingual data. The resultant target-side corpora sizes are summarised in Figure 4.
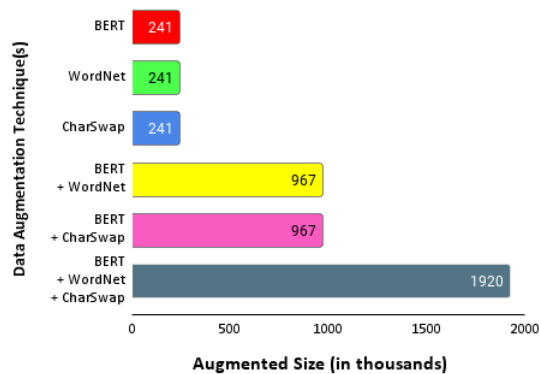


Figure 4: Effective size of the target monolingual corpora with different combinations of augmentation.

## 5 Experimental Setup

All our experiments have been implemented in PyTorch, except for the Bi-LSTM and CRF which were done in Tensorflow. In addition to this, we used FairSeq (Ott et al., 2019), FLAIR (Akbik et al., 2019a), HuggingFace Transformers (Wolf et al., 2019), and Open-NMT (Klein et al., 2017) frameworks in Python. Nvidia Apex was used for memory optimisation using fp-16 training. Experiments related to Bi-LSTM, CRF, vanilla transformers, and FT were performed on a single $8GB$ Nvidia

| Technique | S | US | SS | HE |
|---|---|---|---|---|
| *Vanilla Transformer* | | | | |
| UrIIISeg | 36.32 | | | 2.202 |
| UrIIIComp | 33.45 | | | 2.242 |
| AllSeg | 37.01 | | | 2.360 |
| AllComp | 42.23 | | | 2.431 |
| +3×FT* | | | 41.98 | 2.358 |
| **+5×FT** | | | **44.14** | **2.504** |
| +7×FT | | | 42.95 | 2.367 |
| *XLM* | | | | |
| MLM, `Orig` | | 4.49 | 15.04 | |
| MLM + TLM, `WMT` | | 0.94 | – | |
| `Mixed` | | 13.08 | 21.23 | 1.104, – |
| `Orig` | | 12.73 | 24.64 | 1.294, – |
| *XLM + Data Augmentation* | | | | |
| `BERT` | | 13.06 | 29.50 | 1.320, 1.704 |
| `WordNet` | | 13.08 | 28.57 | 1.269, 1.690 |
| `CharSwap` | | 12.92 | 29.04 | |
| `BERT+WordNet` | | 13.34 | 26.57 | 1.460, 1.666 |
| `BERT+CharSwap` `+WordNet` | | 13.23 | 30.10 | – , 1.757 |

Table 1: Sumerian-English Machine Translation. Here, S: Supervised, US: Unsupervised, SS: Semi-Supervised and HE: Human Evaluation. Each of the available values for the first three columns (BLEU) is compared with a value under HE (out of 3). *Number of shards used for FT.

| | F1-Score |
|---|---|
| HMM | 0.815 |
| **Rules + CRF** | **0.991** |
| Bi-LSTM + CRF | 0.763 |
| FLAIR | 0.499 |
| RoBERTa | 0.949 |

Table 2: POS Tagging for Sumerian. CRF with rules outperform large models like FLAIR and RoBERTa.

| | F1-Score |
|---|---|
| HMM | 0.656 |
| Rules + CRF | 0.913 |
| Bi-LSTM + CRF | 0.775 |
| FLAIR | 0.187 |
| **RoBERTa** | **0.953** |

Table 3: NER for Sumerian. RoBERTa performs best among others. Due to high character-level noise, FLAIR fails to generalise well.

GeForce RTX 2070 GPU, while the pre-training and fine-tuning of FLAIR, RoBERTa, and XLM on various data configurations were performed on 2 16 GB Nvidia V100 GPUs. We used development sets to tune the hyper-parameters for all our models, especially those for POS and NER. For RoBERTa and vanilla transformer, $N = 6$ encoder layers with $h = 16$ attention heads were used, while $N = 4$ and $h = 12$ was used for XLM. A beam-size of 5 was used for our FT experiments. Adam (Kingma and Ba, 2015) optimiser with a learning rate of 0.001, $\beta_1 = 0.90$, $\beta_2 = 0.98$ and a decay factor of 0.5 was used. Additional regularisation was done via Dropout and Attention Dropout (wherever applicable) layers with $p_{drop} = 0.1$. We used a batch size of 32 or 64 and an early stopping criteria based on the validation loss.

# 6 Results and Analysis

**Sequence Labeling** Tables 2 and 3 represent the metric scores of our different models for POS and NER tasks, respectively. CRF with domain-specific rules gives the best F1-score for the POS tagging task, even better than the complex RoBERTa and FLAIR language models which are the current state-of-the-art techniques for most languages. The prevalence of distorted words and short phrases in the corpora makes context learning difficult, although the domain-specific rules help learn short-term dependencies by learning feature weights.

RoBERTa performs well for both of the tasks, while being the best among others for NER (95.37 F1 score). To make the most out of the limited vocabulary and noisy text, we used Byte-Level BPE (Radford et al., 2019) to train the language model and further fine-tuned it on our POS and NER dataset with a batch size of 128. We also tried FLAIR language model across various word embeddings (character, Word2vec, FastText, GloVe) along with an additional CRF layer for both of the tasks. Although a high precision is observed using this approach, the F1 scores is seen to be significantly low due to low recall. In addition to the F1 metric we also

conducted human evaluation by language expert for the best performing models, out of randomly selected 76 (496 words) phrases, only 8 and 6 words were misclassified by NER and POS models, giving an error of 1.20 and 1.61%, respectively.

**Machine Translation** Table 1 summarises our results for all supervised, semi-supervised, and unsupervised techniques. Forward translation on vanilla transformer outperforms all other techniques by at least 2 BLEU. The variation of its performance with more monolingual source text is shown. The superior performance of `AllComp` over the other configurations in vanilla transformer signifies the value of both context and out-of-domain data together. Even though the XLM-based models show lower performance, it could be attributed to the lesser number of encoder layers and attention heads used for them. What is interesting to note, though, is the variation of its performance across various training strategies. We experiment across MLM and TLM (+ MLM) initialization for XLM, where the latter comfortably outperforms the former. We do not test with random initialization and CLM, following up from the conclusions made for NMT in Conneau and Lample (2019). Pre-training the XLM on augmented target-side text works surprisingly well. We note that using pre-training on `BERT` and `WordNet` augmentations results in better Unsupervised performance while introducing `CharSwap` improves the semi-supervised models. The human evaluation presented in the table was made by three Assyriologists, who rated 100 output examples for each model, on a scale of 3. A pairwise inter-annotator agreement of 0.673 (Cohen's Kappa) was observed.[6]

## 7 Interpretability Analysis

Oftentimes in case of Deep Learning Architectures, metric scores like Accuracy, F1 and BLEU are unable to portray the true behavior of the models. For languages like Sumerian, the human-understanding itself is scarce. Visualizing the representations and correlations made by the model could provide insights into which elements of the context can give additional information to support semantic analysis of the terms. Thus, we herein introduce a generalisable interpretability toolkit, *InterpretLR*, to interpret algorithms for **L**ow-**R**esource NLP and

---

[6]Elaborate evaluation criteria mentioned in the Appendix.

further apply it for the aforementioned tasks and models.

*InterpretLR* is primarily aimed at fabricating attribution saliency maps, i.e., tracing back the model output so as to assign an importance score to each input token, based on its 'influence' on that output. We do this using two kinds of interpretability techniques– gradient-based (Sundararajan et al., 2017; Simonyan et al., 2014; Shrikumar et al., 2017), and perturbation-based (Zeiler and Fergus, 2014; Castro et al., 2009).

Due to the inherently discrete nature of natural language text, the starting point for all our approaches is the embedding of the input sentence across the model to interpret. Most of our analysis is done for the encoder of the network architecture, thus analyzing the effect of different pre-training and fine-tuning techniques on how the model eventually represents the language attributes. We use the word 'Attribution' as a better-defined substitute for the 'Influence' measure of an input span of text on the output.

A part of our visual analysis is shown and elaborated here, while a complete analysis with all our models and layer-wise heat-maps is presented in the Appendix.

In Table 4a, we apply *InterpretLR* on 3 different configurations of XLM for a randomly chosen sentence from NMT's evaluation set. A human expert was asked to annotate the source sentence in accordance with the expected reference for each output token in the actual English translation, as shown in the first column. The highlighted visualizations for each of the 3 models were obtained using Integrated Gradients (Sundararajan et al., 2017) across the three input embeddings- token, position, and language. A lot of interesting observations could be made from these attributions.

Firstly, the named entity in the sentence *ur-{d}asznan* (*UrAnan*) has been wrongly translated by all the three models. Although this behavior is expected (learning the context of a named entity is extremely difficult without excessive supervision around the same, which is largely absent our training text) the models even largely fail to attend to the right words in the input.

Secondly, words like *rations*, *weavers* and *seal* which appear frequently in the parallel Ur III Admin corpora and have a contextual meaning attached to them, are translated perfectly by the models, this property is observed among these models

| Actual | Human Expert | Model-1 | Semi-Supervised DataAug XLM | Model-2 | Unsupervised DataAug XLM | Model-3 | Unsupervised `Orig` TLM XLM |
|---|---|---|---|---|---|---|---|
| Output Word | Annotations | Output Word | Visualisations | Output Word | Visualisations | Output Word | Visualisations |
| barley | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | barley | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | Monthy | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | Basketoftablets | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e |
| rations | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | rations | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | rations | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | rations | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e |
| weavers | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | weavers | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | weavers | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | weavers | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e |
| under | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | under | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | from | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | 255 | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e |
| seal | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | seal | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | seal | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | seal | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e |
| of | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | of | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | of | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | of | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e |
| UrAnan | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | Lugalniglagare | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | Ninlil | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | weavers | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e |
| foreman | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | foreman | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | foreman | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e | female | #s sze-ba geme2 usz-bar kiszib3 ur-{d}asznan ugula #e |

(a) MT- Selected output tokens for Sumerian Input text of "sze-ba geme2 usz-bar kiszib3 ur-dasznan ugula", which translates to "barley rations of the female weavers under seal of UrAnan the foreman".[7]

| Actual | Human Expert | Model | RoBERTa |
|---|---|---|---|
| N | 5(disz) gin2 ku3-babbar | N | 5 ( disz ) gin 2 ku- 3 - babbar |

(b) POS- With tagged word "ku3-babbar"

| Actual | Human Expert | Model | RoBERTa |
|---|---|---|---|
| GN | mu ur-bi2-lum{ki} ba-hul | GN | mu ur - bi 2- lum { ki } ba - hul |

(c) NER- With tagged word "ur-bi2-lumki"

Table 4: Highlighted attributions for randomly selected examples. **Green** and **Red** represent correct and wrong predictions, respectively, while Green and Red highlights represent positive and negative attributions, respectively.

in general. Even the unsupervised models that do not have access to the one-to-one mapping of the translation during training manage to infer these words from the appropriate context. It can be assumed that they learn the right representations of such tokens. But at the same time, there are instances like *sze-ba* (*barley*), which the two unsupervised models rightly refer to but do not give the right translations, which thus is a direct result of the absence of supervision.

Lastly, English words like *under*, *of* and *from* do not have any direct translations in Sumerian and are mostly inferred from the context, even by the human annotators. At such places, again, supervision might play a critical role as in the $4^{th}$ row of Table 4a. There are also instances like the $6^{th}$ row where the supervised model fails to attend to the right words, and the correct output word could very well be out of memorisation.

Tables 4b and 4c represent visualizations for two randomly selected phrases for our sequence labeling tasks, indicating the attributions for each sub-word for tagging the corresponding target word with their predicted labels. It can be observed from Table 4b that word *gin* (*unit*) and sub-word *ku*, are contributing to the attribution score positively, depicting positive model attribution to tag *ku3-babbar*

as a Noun (N), whereas in Table 4c the sub-words *ur*, *hul* and *ki* are contributing *ur-bi2-lum{ki}* to be tagged as the label GN (Geographical Name). As observed from the corresponding human annotation, *ur* and *ki* are the most associated for Geographical names and GNs are mostly followed by a verb part, which is *hul* (*destroy*) in this case. It can thus be inferred that RoBERTa identifies this correspondence well and makes the decision accordingly.

## 8 Conclusion

In this work, we introduced the first information extraction and translation pipeline for Sumerian cuneiform. We first undertook the tasks of POS Tagging and NER, where we observed that *deeper is not necessarily better*. A simple CRF model with well-defined rules outperformed the large language model RoBERTa for POS Tagging. Further, for machine translation we overcame unprecedented challenges pertaining to lack of in-domain text, sparse sentence formation, and incoherence. We found that using out-of-domain text along with specific data-augmentation can have huge impacts in a low-resource setting. All components of this work are generalisable to other low-resource languages, including *InterpretLR*, and we open way to future research in this direction.

---

[7]The left-out tokens were rightly predicted by all the three models, with almost the same attributions.

## References

Alan Akbik, Tanja Bergmann, Duncan Blythe, Kashif Rasul, Stefan Schweter, and Roland Vollgraf. 2019a. FLAIR: An easy-to-use framework for state-of-the-art NLP. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics (Demonstrations)*, pages 54–59, Minneapolis, Minnesota. Association for Computational Linguistics.

Alan Akbik, Tanja Bergmann, and Roland Vollgraf. 2019b. Pooled contextualized embeddings for named entity recognition. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 724–728. Association for Computational Linguistics.

Alan Akbik, Duncan Blythe, and Roland Vollgraf. 2018. Contextual string embeddings for sequence labeling. In *Proceedings of the 27th International Conference on Computational Linguistics, COLING 2018, Santa Fe, New Mexico, USA, August 20-26, 2018*, pages 1638–1649. Association for Computational Linguistics.

Franck Burlot and François Yvon. 2018. Using monolingual data in neural machine translation: a systematic study. In *Proceedings of the Third Conference on Machine Translation: Research Papers*, pages 144–155, Brussels, Belgium. Association for Computational Linguistics.

Javier Castro, Daniel Gómez, and Juan Tejada. 2009. Polynomial calculation of the shapley value based on sampling. *Comput. Oper. Res.*, 36(5):1726–1730.

Alexis Conneau and Guillaume Lample. 2019. Cross-lingual language model pretraining. In *Advances in Neural Information Processing Systems 32: Annual Conference on Neural Information Processing Systems 2019, NeurIPS 2019, December 8-14, 2019, Vancouver, BC, Canada*, pages 7057–7067.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. 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, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics.

Sergey Edunov, Myle Ott, Michael Auli, and David Grangier. 2018. Understanding back-translation at scale. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, pages 489–500. Association for Computational Linguistics.

Robert K. Englund. 2009. The smell of the cage. https://cdli.ucla.edu/pubs/cdlj/2009/cdlj2009_004.html. Online; accessed 2009.

Di He, Yingce Xia, Tao Qin, Liwei Wang, Nenghai Yu, Tie-Yan Liu, and Wei-Ying Ma. 2016. Dual learning for machine translation. In *Advances in Neural Information Processing Systems 29: Annual Conference on Neural Information Processing Systems 2016, December 5-10, 2016, Barcelona, Spain*, pages 820–828.

Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural Comput.*, 9(8):1735–1780.

Ari Holtzman, Jan Buys, Maxwell Forbes, Antoine Bosselut, David Golub, and Yejin Choi. 2018. Learning to write with cooperative discriminators. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers*, pages 1638–1649. Association for Computational Linguistics.

Zhiheng Huang, Wei Xu, and Kai Yu. 2015. Bidirectional lstm-crf models for sequence tagging. *arXiv preprint arXiv:1508.01991*.

Diederik P. Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*.

Guillaume Klein, Yoon Kim, Yuntian Deng, Jean Senellart, and Alexander Rush. 2017. OpenNMT: Open-source toolkit for neural machine translation. In *Proceedings of ACL 2017, System Demonstrations*, pages 67–72, Vancouver, Canada. Association for Computational Linguistics.

John D. Lafferty, Andrew McCallum, and Fernando C. N. Pereira. 2001. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of the Eighteenth International Conference on Machine Learning (ICML 2001), Williams College, Williamstown, MA, USA, June 28 - July 1, 2001*, pages 282–289. Morgan Kaufmann.

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. *CoRR*, abs/1907.11692.

Yudong Liu, Clinton Burkhart, James Hearne, and Liang Luo. 2015. Enhancing sumerian lemmatization by unsupervised named-entity recognition. In *NAACL HLT 2015, The 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Denver, Colorado, USA, May 31 - June 5, 2015*, pages 1446–1451. The Association for Computational Linguistics.

Liang Luo, Yudong Liu, James Hearne, and Clinton Burkhart. 2015. Unsupervised sumerian personal name recognition. In *Proceedings of the Twenty-Eighth International Florida Artificial Intelligence Research Society Conference, FLAIRS 2015, Hollywood, Florida, USA, May 18-20, 2015*, pages 193–198. AAAI Press.

George A Miller, Richard Beckwith, Christiane Fellbaum, Derek Gross, and Katherine J Miller. 1990. Introduction to wordNet: An on-line lexical database. *International journal of lexicography*, 3(4):235–244.

Myle Ott, Sergey Edunov, Alexei Baevski, Angela Fan, Sam Gross, Nathan Ng, David Grangier, and Michael Auli. 2019. fairseq: A fast, extensible toolkit for sequence modeling. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Demonstrations*, pages 48–53. Association for Computational Linguistics.

Émilie Pagé-Perron, Maria Sukhareva, Ilya Khait, and Christian Chiarcos. 2017. Machine translation and automated analysis of the Sumerian language. In *Proceedings of the Joint SIGHUM Workshop on Computational Linguistics for Cultural Heritage, Social Sciences, Humanities and Literature*, pages 10–16, Vancouver, Canada. Association for Computational Linguistics.

Ravneet Punia, Niko Schenk, Christian Chiarcos, and Émilie Pagé-Perron. 2020. Towards the first machine translation system for sumerian transliterations. In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 3454–3460.

Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. 2019. Language models are unsupervised multitask learners.

Mike Schuster and Kuldip K. Paliwal. 1997. Bidirectional recurrent neural networks. *IEEE Trans. Signal Process.*, 45(11):2673–2681.

Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving neural machine translation models with monolingual data. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL 2016, August 7-12, 2016, Berlin, Germany, Volume 1: Long Papers*. The Association for Computer Linguistics.

Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning important features through propagating activation differences. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 3145–3153. PMLR.

Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2014. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Workshop Track Proceedings*.

Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. MASS: masked sequence to sequence pre-training for language generation. In *Proceedings of the 36th International Conference on Machine Learning, ICML 2019, 9-15 June 2019, Long Beach, California, USA*, volume 97 of *Proceedings of Machine Learning Research*, pages 5926–5936. PMLR.

Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*, volume 70 of *Proceedings of Machine Learning Research*, pages 3319–3328. PMLR.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. Attention is all

you need. In *Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA*, pages 5998–6008.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, and Jamie Brew. 2019. Huggingface's transformers: State-of-the-art natural language processing. *CoRR*, abs/1910.03771.

Lijun Wu, Yiren Wang, Yingce Xia, Tao Qin, Jianhuang Lai, and Tie-Yan Liu. 2019. Exploiting monolingual data at scale for neural machine translation. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4207–4216, Hong Kong, China. Association for Computational Linguistics.

Matthew D. Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *Computer Vision - ECCV 2014 - 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part I*, volume 8689 of *Lecture Notes in Computer Science*, pages 818–833. Springer.

Jiajun Zhang and Chengqing Zong. 2016. Exploiting source-side monolingual data in neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1535–1545, Austin, Texas. Association for Computational Linguistics.

Barret Zoph, Deniz Yuret, Jonathan May, and Kevin Knight. 2016. Transfer learning for low-resource neural machine translation. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing, EMNLP 2016, Austin, Texas, USA, November 1-4, 2016*, pages 1568–1575. The Association for Computational Linguistics.

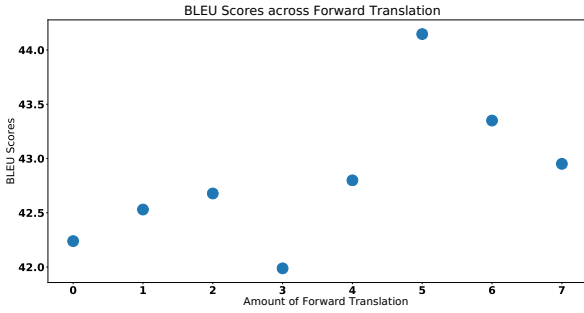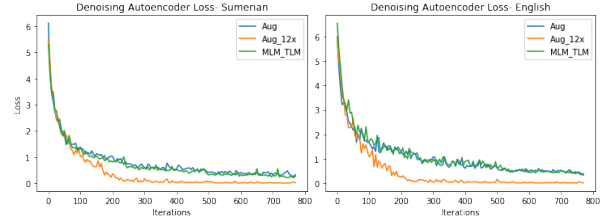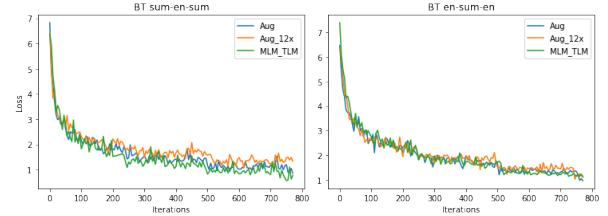## A Detailed Evaluation and Analysis

Figure 5

Forward Translation with Vanilla Transformer gave the best results for Sumerian-English Neural Machine Translation. Figure 5 shows the variation of the BLEU score with the amount of source monolingual data used. Here, the X-Axis represents the number of shards used, with each shard consisting of 80K sentences. It can be observed that the translation accuracy is not linear with the amount of text used.

Figure 6 shows the variation of several performance metrics during the Unsupervised fine-tuning of various XLM configurations. The comparison is made between XLM pre-training without any data augmentation (MLM_TLM), with one augmentation (Aug) and with all three augmentations (Aug_12x). It can be seen from Figure 6a that an XLM pre-trained on the Aug_12x configuration converges the fastest among the others, in terms of the main Denoising Auto-encoding Loss. It can also be observed that the curve corresponding to this configuration is much smoother than the others, which shows a positive regularizing effect of a better weight initialisation (through appropriate pre-training). A similar pattern is observed for the validation accuracy across the epochs as shown in Figure 6c, although, the trend of Back Translation loss remains mostly inseparable for the three configurations.
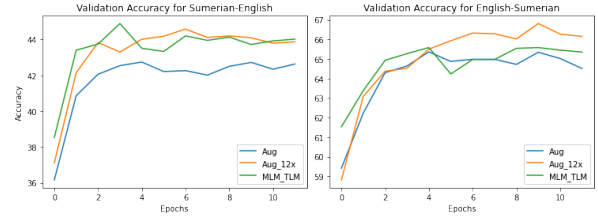
Table 5 depicts the net percentage error found by an human expert on the POS and NER results for the entire evaluation set across the best performing model. Table 6 and 7 represents the detailed results of POS and NER models. It can be observed from the tables, that although CRF and RoBERTa models gave the best results, FLAIR language model along with character embeddings also gave high precision for both of the tasks.

(a) Denoising Auto-encoder Loss (AE Loss) variation across the 1st Epoch

(b) Back Translation Loss variation in XLM across the 1st Epoch

(c) MT accuracy across a number of training epochs

Figure 6: Quantitative comparison of various models during Unsupervised MT fine-tuning

|  | POS error (in %) | NER error (in %) |
|---|---|---|
| Human Evaluation | 1.61 | 1.20 |

Table 5: Human Evaluation for POS and NER

## B Extended Interpretations

Here we present the interpretability analysis across a larger set of models and visualisations. We use and compare the different algorithms across layer-level, gradient-based, and perturbation-based techniques to obtain the attributions.

Figure 7 visualises the Multi-head Self Attention (MHSA) using Layer Conductance Dhamdhere, Sundararajan, and Yan 2018) across the 4 encoder layers we employ in XLMs[8]. The first two output tokens *barley* and *female* are known to be one-on-one mapping between the input words of *sze-ba* and *geme2* respectively. While the third output token *barley* is not a direct translation and

---

[8]The supervised version of the augmented pre-training is used here.

| | Part of Speech Tagging | | |
|---|---|---|---|
| | Precision | Recall | F1-Score |
| HMM | 0.857 | 0.794 | 0.815 |
| Rules + CRF | **0.994** | **0.989** | **0.991** |
| BBi-LSTM + CRF | 0.852 | 0.710 | 0.7631 |
| FLAIR | 0.9323 | 0.4766 | 0.4999 |
| RoBERTa | 0.9500 | 0.9489 | 0.9495 |

Table 6: POS Tagging Models for Ur III Sumerian Text

| | Named Entity Recognition | | |
|---|---|---|---|
| | Precision | Recall | F1-Score |
| HMM | 0.810 | 0.599 | 0.656 |
| Rules + CRF | 0.916 | 0.910 | 0.913 |
| Bi-LSTM + CRF | 0.864 | 0.704 | 0.775 |
| FLAIR | **0.9562** | 0.1817 | 0.1873 |
| RoBERTa | 0.9540 | **0.9534** | **0.9537** |

Table 7: NER Models for Ur III Sumerian Text



Figure 7: Layer Conductance across MHSA Layers

is needed to be inferred from context.

Figure 9a represents the attribution heat-map when gradient-normalisation saliency (Simonyan, Vedaldi, and Zisserman 2013) is used. Being one of the most conventional techniques for finding attribution, it is more prone to inconsistent interpretations. Whereas, the attribution heat-map in Figure 9b represents the Integrated Gradients (IG) (Sundararajan, Taly, and Yan512017) approach. Being a path-based technique, which measures the gradient attribution relation using a straight-line path from a baseline (usually all-zeros), to the given input, it is much more robust and stable.

Even though the gradient-based methods are much faster than perturbation-based methods, we observe that the heavy dependency of IG on hyper-parameters like the number of input steps to be considered when going from a baseline to the actual input, $n\_steps$, to be a major setback. The final attribution is generally found out after integrating (or summing) over the attributions of these sub-steps. We found that the attributions do not change when going beyond $n\_steps = 250$, thus, we experiment by varying it between 10 to 250. We observe that there is no ideal value of $n\_steps$, IG's faithfulness to the model varies largely over this range. For some inputs, t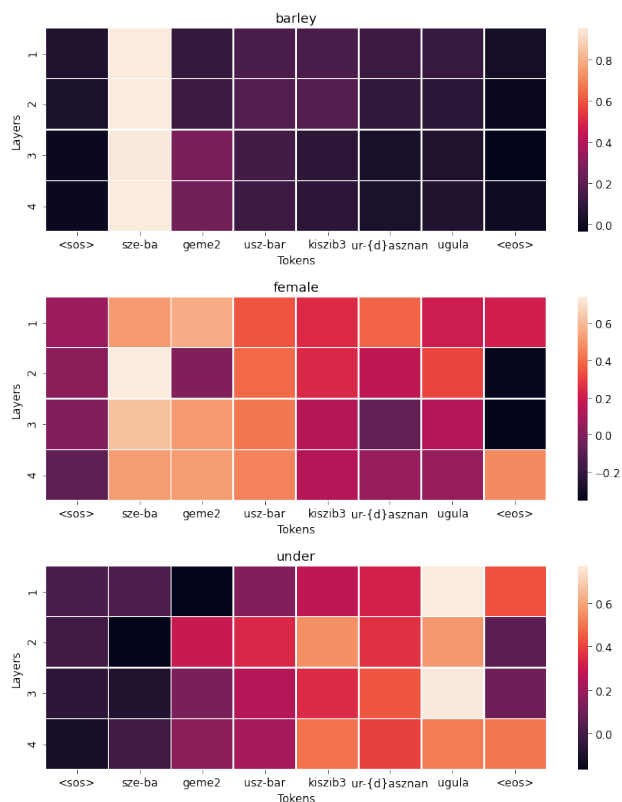he best value is $n\_steps = 50$ while for others $n\_steps = 250$ is the most ideal. We judge this by considering how much the attribution is given to *sos* and *eos* tokens for each output token. Thus, based on both *plausibility* and *faithfulness*. We use $n\_steps = 50$ for obtaining the heat-maps in Figure 9b.

Figure 10 represents the visualization for our sequence labeling tasks. It indicates two major things, 1) the effect of words, sub-words (depends on tokenization) on tagging the target word and 2) the effect of 6 transformer encoder layers. We created the hook on embeddings of RoBERTa with layer IG and obtained the visualizations for how each sub-word is contributing to tag the target word. Similarly, to obtain the heat-map we created the hook on RoBERTa embeddings and used the Layer Conductance.

From Figure 10a it can be observed that *ku* and *du* contribute the most to the attribution scores for tagging *ku3-babbar* and *ba-du3* as a Noun (N) and Verb (V), respectively. From the heat-maps it is also noted that *ku* shows the effect on all 6 layers whereas in second example effects are majorly due to the initial transformer layers. Similarly in the Figure 10b *ur* and *lugal* are the most effective sub-words to tag *ur-bi2-lumki* and *lugal-tesz2-mu* as
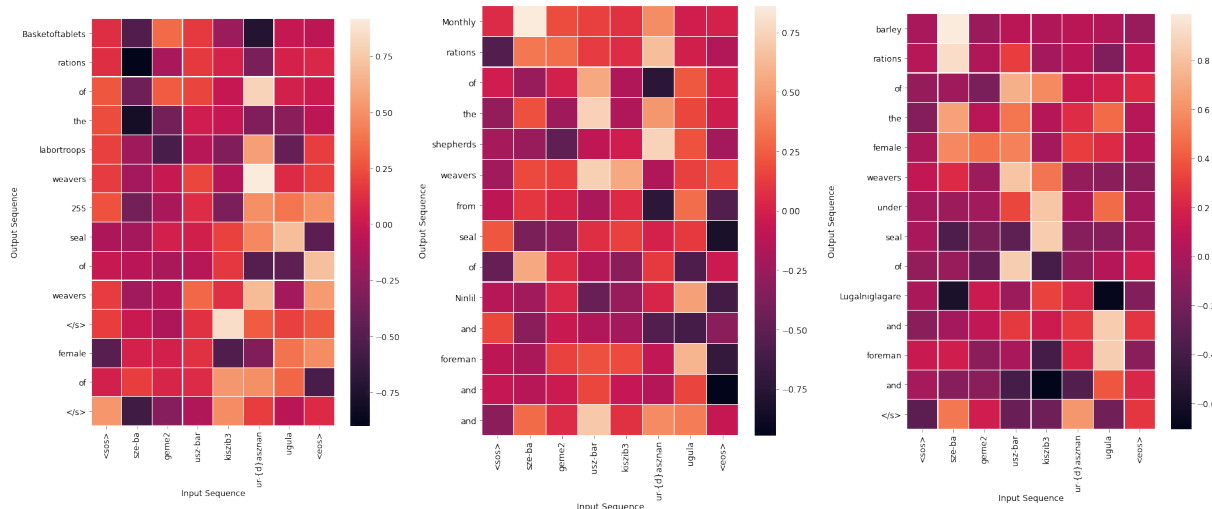
Figure 8: Feature Ablation in *InterpretLR*

GN (Geographical Name) and PN (Personal Name) respectively. It is also interesting to note that both of these sub-words have a very positive impact in the initial layers but are contributing oppositely in the last layer.

### B.1 Human Evaluation

The scoring by human experts was done independently for each result according to the following criterion:
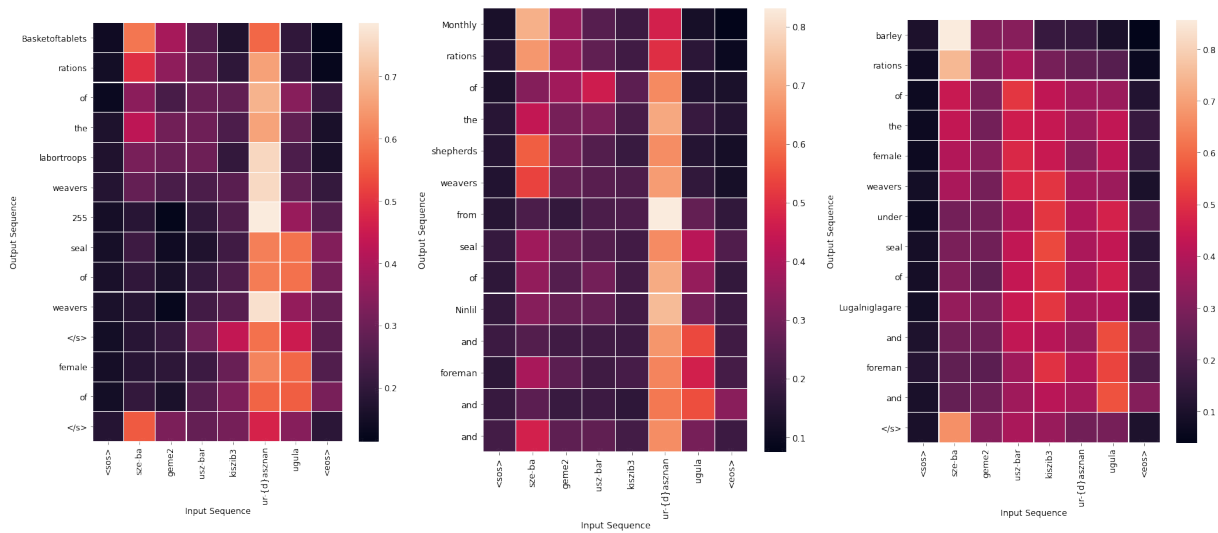
- **3 (good)**: interpretable in the correct meaning by a native speaker of English; (almost) no incorrectly translated content word (e.g., tolerant against some errors in word order, but not in incorrect words).

- **2 (helpful)**: partially distorted, but interpretable with some context information (tolerant against errors in word order and against incorrect function words).

- **1 (incorrect)**: contains incorrectly translated content words and/or is un-interpretable.
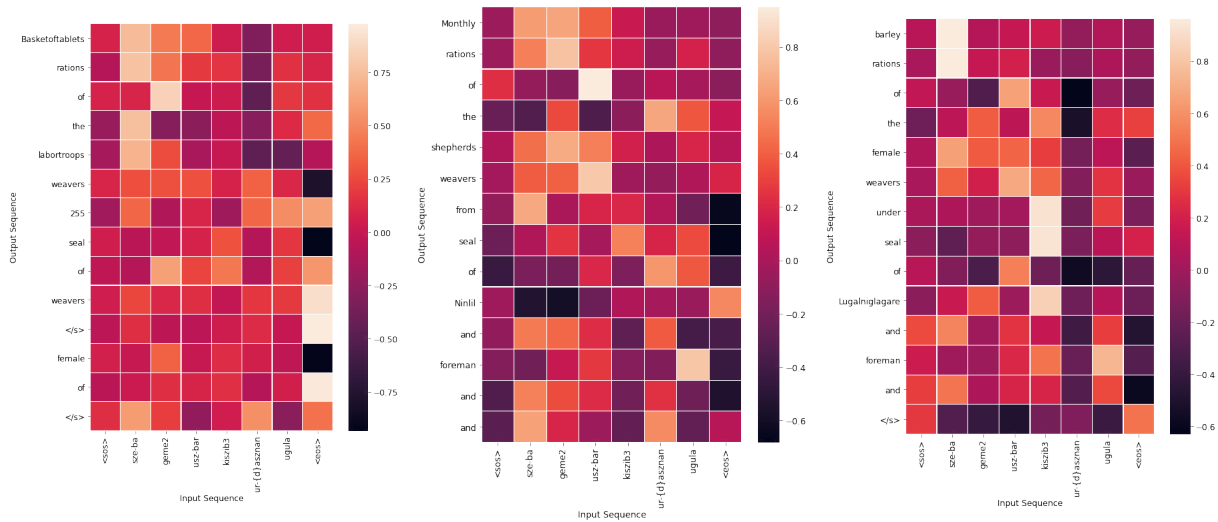
### C Rules for POS Tagging and NER

We used certain language-specific rules to assist CRF for the sequence labeling tasks. The rules were identified by human experts and some of them are as mentioned here:

- A word starting with "ur-", "lu2-", or "dumu" is most likely to be a personal name.

- If a word is followed by "mu", then the next phrase denotes a year name.

- If a word is followed by "iti", it denotes a month name.

- Words containing "ki" are mostly associated with geographical names (GN).

- Words ending with part "-hul" majorly denotes verbs.

- Words containing "{d}" denotes either personal name (PN) or divine name (DN).
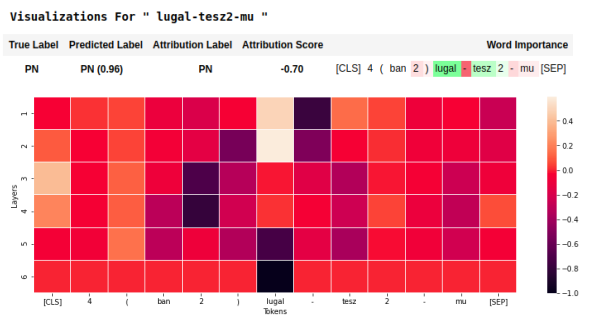
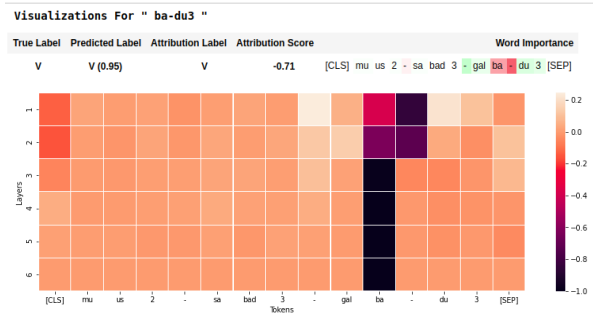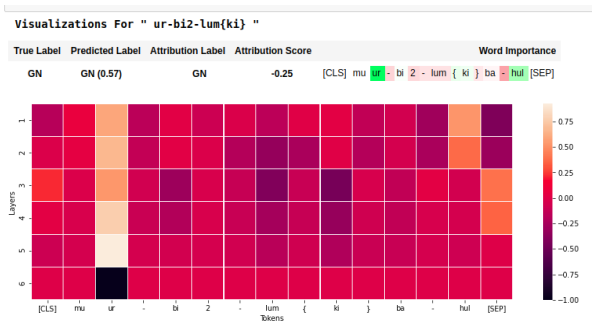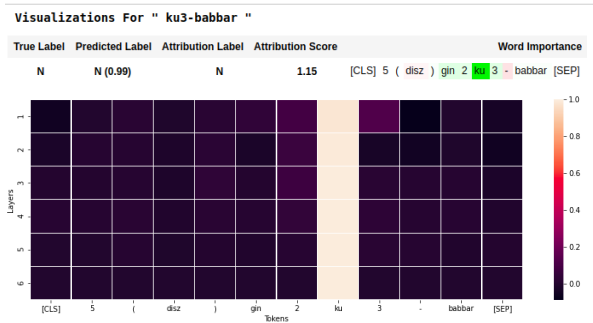- A word followed by "gin" (*unit*) majorly replicate a noun.

(a) Grad-Norm Saliency



(b) Integrated Gradients

Figure 9: Comparing different gradient-based approaches used in *InterpretLR*

(a) POS Tagging

(b) NER

Figure 10: *InterpretLR* on RoBERTa for Sequence Labeling