

# MSNER: A Multilingual Speech Dataset for Named Entity Recognition

Quentin Meeus<sup>1,2</sup>, Marie-Francine Moens<sup>1</sup>, Hugo Van hamme<sup>2</sup>

20th Joint ACL-ISO Workshop on Interoperable Semantic Annotation

<sup>1</sup> LIIR Lab, Computer Science Dpt., KU Leuven    <sup>2</sup> PSI, Electrical Engineering Dpt., KU Leuven  
Quentin.Meeus@kuleuven.be

## Abstract

While extensively explored in text-based tasks, Named Entity Recognition (NER) remains largely neglected in spoken language understanding. Existing resources are limited to a single, English-only dataset. This paper addresses this gap by introducing MSNER, a freely available, multilingual speech corpus annotated with named entities. It provides annotations to the VoxPopuli dataset in four languages (Dutch, French, German, and Spanish). We have also releasing an efficient annotation tool that leverages automatic pre-annotations for faster manual refinement. This results in 590 and 15 hours of silver-annotated speech for training and validation, alongside a 17-hour, manually-annotated evaluation set. We further provide an analysis comparing silver and gold annotations. Finally, we present baseline NER models to stimulate further research on this newly available dataset.

**Keywords:** Spoken Named Entity Recognition, Spoken Language Understanding, Speech Dataset

## 1. Introduction

In an increasingly interconnected world where language knows no boundaries, the field of Speech Processing is undergoing a transformative shift towards multilingual applications. One such pivotal area is Spoken Named Entity Recognition (Spoken NER). Named Entity Recognition (NER) is a natural language processing (NLP) task that involves the identification and categorization of named entities within a text, typically into predefined categories such as names of persons, organizations, locations, dates, numerical values, and more. The primary objective of NER is to automatically recognize and extract specific pieces of information from unstructured text, making it easier to analyze and understand the content. NER plays a crucial role in various NLP applications, including information retrieval, question answering, sentiment analysis, and language understanding. In contrast, *Spoken NER* extracts named entities from audio documents, a task that is considerably more challenging. Indeed, aside from the inherent difficulties associated with speech processing, Spoken NER requires not only to identify and classify the entities, but also to transcribe them correctly. Variability in pronunciation, accents, and dialects can make the detection and especially the spelling of named entities very challenging. On the other hand, prosody, intonation and emphasis are cues that may be crucial for NER but are not readily available in written text. Recognizing the pressing need to facilitate cross-lingual research and to provide comprehensive evaluation resources for Spoken NER models, we have undertaken the task of manually annotating the popular speech dataset VoxPopuli's test sets in four

languages: Dutch, French, German, and Spanish. Additionally, we also provide machine-made annotations on the training and validation sets.

In the following sections, we provide a detailed overview of our efforts in the domain of Spoken NER. First, we give an overview of related works and datasets. Then, we introduce the newly annotated dataset and provide information about its size, multilingual coverage, and its potential significance in advancing Spoken NER technology. Additionally, we describe the methodology employed in the dataset's creation, breaking down the annotation process and data preparation. We also introduce the user-friendly annotation interface we've developed for this purpose. Furthermore, we present the results of various experiments and benchmarks conducted using this dataset. These experiments demonstrate its utility in evaluating Spoken NER models across the chosen languages, highlighting its role in advancing research and development in this field.

In summary, this article describes our contributions to the field of multilingual Spoken NER, including the dataset's creation, annotation methodology, and its role in advancing research in this domain.

## 2. Literature Review

In the field of NLP, there is not one unified label set. Both generic and specialized datasets exist with their own label sets defined. Specialized datasets might cover large amounts of topics with specific vocabulary and entities. For example, a NER system for doctors would include medications, dosages, medical reasons, etc. (Uzuner et al., 2010), and biomedical entities include names

of proteins, chemical, disease, or species (Crichton et al., 2017). Other datasets provide more generic entities that cover broader landscapes. One of the most widely used is CoNLL-2003 (Tjong Kim Sang and De Meulder, 2003), although it comes with only four entity types (LOC, ORG, PER and MISC). OntoNotes v5 enriches this set with 14 more classes (Table 2), to include things such as numbers, dates, and laws. Its high quality makes it one of the most widely used NER datasets, although it only covers three languages: English, Arabic and Chinese. Another notable mention is Tedeschi et al. (2021), which adds a few more generic classes to OntoNotes definitions to cover things such as animal names, diseases, food, and plants, and released a dataset derived from Wikipedia where named entities were annotated automatically with an annotation pipeline that effectively combined pretrained language models and knowledge-based approaches. A follow-up dataset was published covering more languages (Tedeschi and Navigli, 2022).

Currently, we know of only one Spoken NER dataset that is openly distributed as SLUE (Shon et al., 2021). This is an annotated subset of the larger VoxPopuli dataset (Wang et al., 2021), which comprises audio recordings and corresponding transcripts of sessions held in the European Parliament. The annotated portion of the dataset include approximately 25 hours of speech, divided into three subsets: 3/5 for training, 1/5 for validation, and 1/5 for testing purposes. While this initiative is a significant step forward, SLUE exclusively covers the English language. They used the same entities as OntoNotes (Weischedel et al., 2013) although in practice, they combine some types and remove rare ones to produce a new label set (Table 2, Column 2).

Another task in spoken language understanding is similar to Spoken NER: slot filling. This is the identification of information relevant to specific applications, such as flight booking (Hemphill et al., 1990). Although they share many grounds, there is a major difference: slot filling relates to a specific application, and in this regard, covers a much narrower domain than NER, often consisting of short commands for a computer interface (Lugosch et al., 2019; Saade et al., 2018; Bastianelli et al., 2020; Lugosch et al., 2021; Renkens and Van hamme, 2018) or a booking system (Hemphill et al., 1990). Since the vast majority of entity recognition datasets are text-based, the same goes for the applications. Consequently, NER is often framed as a token classification task, where each word or word piece must be assigned an entity type. Since an entity can cover many tokens, the entity classes are redefined in the BIO format, a widely used tagging scheme in NER tasks (Ramshaw and Marcus,

1995). This format provides a structured way to label and distinguish the boundaries of named entities within the text. Each word or token is tagged with one of three labels: “B” marks the beginning, or first word of an entity, “I” indicates the continuation of the named entity and always follows the “B” tag, and “O” is used for words that are not part of an entity. This marker, together with the entity type, makes the target for the classification task. Other annotation schemes are extensions of this (e.g. IO, IOBES, IOE, etc.). The major drawback of the BIO format is its inability to represent nested entities. The modern approach to NER is to add linear layers to a pretrained language model and fine-tune it on the chosen NER dataset. Sometimes, a conditional random field (CRF) (Lafferty et al., 2001) is added to learn the transition probabilities between the label classes (Ushio and Camacho-Collados, 2021). In Spoken NER, the two main approaches are pipeline and end-to-end models. As the name suggests, pipeline models first use automatic speech recognition to transcribe an audio recording, then use NER to predict the entities. In contrast, end-to-end models do not force the model to make hard decisions by choosing one token over another. Instead, it predicts entities directly from the hidden states. Finally, hybrid models or multitask models predict both the entities and the transcriptions simultaneously (Meeus et al., 2023).

subset	language	duration	size	entities
train	DE	224.5 h	86,410	97,492
	ES	141.5 h	47,611	66,482
	FR	186h	65,952	80,255
	NL	38.5 h	16,533	19,566
dev	DE	4h	1,610	1,880
	ES	4h48	1,529	2,094
	FR	4h22	1,527	1,884
	NL	2h16	963	1,074
test	DE	5h	1,966	2,061
	ES	5h	1,512	2,198
	FR	4h30	1,656	2,004
	NL	2h30	1,120	1,272

Table 1: MSNER Dataset statistics

### 3. Dataset description

The MSNER dataset is an annotated version of the VoxPopuli dataset (Wang et al., 2021) in four languages – Dutch, French, German, and Spanish. VoxPopuli is a collection of recorded sessions from the European Parliament, segmented to contain one or more sentence by one speaker. For each language in scope, we provide three annotated subsets (Table 1): a training and development set with machine-generated “silver” annotations, and a test set with manual “gold” annotations. The subsets

OntoNotes5	SLUE	DE	ES	FR	NL	Examples
date	WHEN	307	276	243	113	125 years ago, 15 maart, 1815—1830, 1997
time		12	21	10	8	24 hours, acht uur, de hele dag, mañana
cardinal number	QUANT	136	167	123	91	1, 10, 10 miljoen, 11, 11 billion
ordinal number		82	100	79	45	First, Ten derde, dritten
quantity		6	2	5	1	one and a half meter, two inches
money		26	16	18	8	200 million EUR, Dertig miljoen euro
percent		21	28	13	22	1 procent, 100%, 15 Prozent
geopolitical area	PLACE	259	285	283	176	Amsterdam, Australië, Barcelona, Belgium
location		128	139	214	110	Afrika, Balkanlanden, Europe
group	NORP	229	244	285	213	African, American, Christian
organization	ORG	621	638	527	362	Amnesty International, Charlie Hebdo
law	LAW	64	108	33	22	Paris Accords, US Constitution
person	PERSON	123	131	100	67	Angela Merkel, Barroso, Beyoncé
facility	-	6	2	8	12	Guantánamo, White House
event	-	23	25	21	8	Europees Semester, Rio conferentie
work of art	-	6	3	4	4	Green Book, Koran
product	-	4	1	2	8	2G, 4G, 5G, iPhone
language	-	3	12	6	2	Latin, Nederlands, Español

Table 2: Number of annotated entities per entity type in the test sets. Column SLUE correspond to the ‘combined’ entity set proposed by Shon et al. (2021).

of the four languages in scope were annotated according to OntoNotes’ 18 classes. The test sets were manually annotated by the authors following the methodology outlined in Section 4. Each example in the annotated dataset contains the VoxPopuli ID to identify the relevant audio recording in the original dataset, the transcribed sentence and the annotated named entities, that is, the list of entities, each composed of a text and a label component (Figure 1). For the silver label datasets, we also provide a probability score of each predicted entity. We discuss in Section 6 how this number is related to the uncertainty of the model.

We use the 18-classes OntoNotes label set (Weischedel et al., 2013). However, following the example from Shon et al. (2021), we provide annotations by using an alternative label set that combines entity types like places or numbers and discard the rarest classes like languages, events, and work of art (Table 2 Column 2).


ID	20090423-0900-PLENARY-26-fr_20 090423-21:55:26_4
Audio	
Text	200 milliards d’euros qu’il faut rapprocher aussi du niveau des déficits des pays européens.
Entities	(MONEY, 200 milliards d’euros) (NORP, européens)

Figure 1: Annotated example

## 4. Methodology

We provide two kinds of label quality: machine-generated “silver” labels and human-annotated “gold” labels. For obvious reasons, the silver labels are much cheaper and easier to produce. Therefore, we only provide human-made annotations for the test sets, and the training and validation sets annotations are entirely machine-generated. The methodology follows these four broad steps: (1) filtering out recordings without or with misaligned transcripts, (2) generate silver labels for all subsets, (3) manually annotate the test sets and (4) verify the human-made annotations to identify and rectify potential labelling errors. We detail each step in the following paragraphs.

### 4.1. Filtering

The VoxPopuli dataset contains a few alignment errors between the spoken content and its corresponding transcript. To address this issue, we employed an automatic speech recognition (ASR) system, initially transcribing the spoken utterances and subsequently calculating the word error rate by comparing the ASR-generated sentence to the provided transcript. For this task, we opted for the Whisper large v2 ASR model (Radford et al., 2022), because it showed near state-of-the-art performance across the selected languages. Notably, this model has been meticulously trained on extensive, well-curated data to perform both audio translation and transcription tasks.

For the training and development sets, we filter out examples with a WER larger than 20%, without verifying that the excluded examples were indeed problematic. This discards about 20% of the Ger-

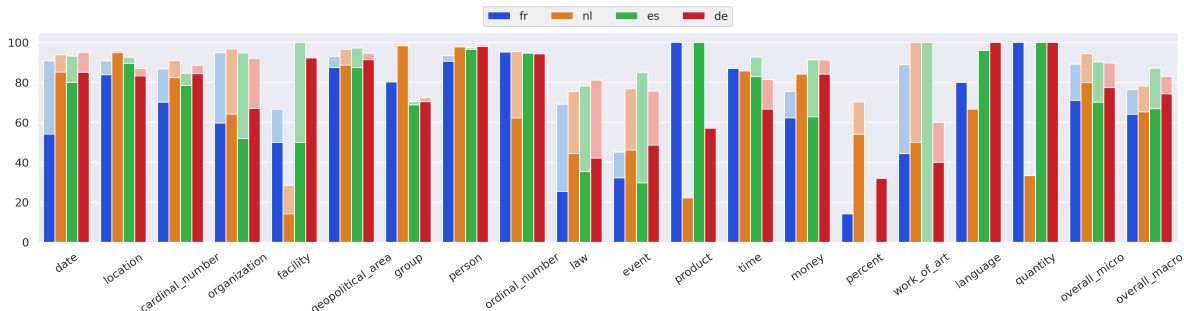


Figure 2: Evaluation of text-based pretrained NER model against our annotations. Bright colors correspond to the F1-score and faded colors correspond to the label-F1 score, a metric that ignores spelling mistakes and segmentation errors.

man and Dutch utterances, 10% of the French examples and 6% of the Spanish utterances.

For the test sets, instances where the word error rate (WER) between the machine-generated transcription and the original transcript exceeded 20%, we conducted a meticulous review process. This involved listening to the audio recording and cross-referencing it with the existing transcript. When feasible, we made necessary corrections to the transcript. However, in cases where multiple speakers were heard in the recording or no speech is present, we removed the problematic utterance from the dataset.

#### 4.2. Pseudo-annotations

We employed an established text-based Named Entity Recognition (NER) model to predict entities within the gold transcript. We chose to use the XLM-RoBERTa large pretrained model (Conneau et al., 2019), fine-tuned specifically on the OntoNotes v5 dataset (Weischedel et al., 2013). This model is readily accessible through the HuggingFace repository<sup>1</sup>.

While it’s important to note that this particular model’s fine-tuning was conducted solely on English data, its robustness and efficacy across multiple languages were remarkable. In our evaluation, we observed impressive performance, with most sentences annotated correctly.

#### 4.3. Annotation Tool

For each of the 6,254 pre-annotated sentences in the test sets, we corrected the annotations predicted by the model. For this purpose, we have developed a command line tool to quickly add, edit, merge or remove annotations in a sentence. This utility displays the pre-annotated sentence with a summary of the annotations below. Annotations appear as colored XML tags both in the text and in the summary. An annotated English translation can be displayed. The annotator then has access

<sup>1</sup><https://huggingface.co/asahi417/tner-xlm-roberta-base-ontonotes5>

to both the original sentence and the translation to make sure that the annotations are as accurate as possible. When presented with a sentence, the annotator has the choice to add a new annotation, delete an existing one, merge two annotations together or modify an annotation, either by changing the type or by adding or removing words. Once a sentence has been annotated, it is saved to a file in JSON format. Following this methodology and with the help of this tool, we were able to save a lot of time and effort without sacrificing accuracy. For this reason, we make the tool available online so that others will have the opportunity to contribute to this field of research by easily annotating more data in many more languages.

As mentioned in Section 3, we not only provide annotations according to OntoNotes 18 classes, but also the 7-classes combined set proposed in Shon et al. (2021). However, we chose to completely re-annotate the examples where entities are removed, instead of simply removing all the annotations of the same type from the dataset. To illustrate this, consider the following example:

```
<event> 15th conference on
speech of Toronto </event>
```

According to the combined set conversion rules (Table 2), all the entities of type `<event>` are to be discarded. Doing that would lead to two unannotated entities, ‘15th’ as a number and ‘Toronto’ as a place. Instead, we re-annotate the examples containing removed entities to make sure that we are not penalizing the models for correct assumptions.

#### 4.4. Verification

Finally, we verify the integrity of the test annotations by deriving a number of heuristics and rules that the annotations must abide. This involved grouping the annotations by category and verify each list one by one, comparing them to one another, searching in the text for frequent annotated terms to identify missing annotations, etc. In this last step, we also fix some remaining transcription issues. For example, we realized that VoxPopuli transcripts omitted





Figure 3: Distribution of predicted probability score per class given the target class for the text-based model’s predictions

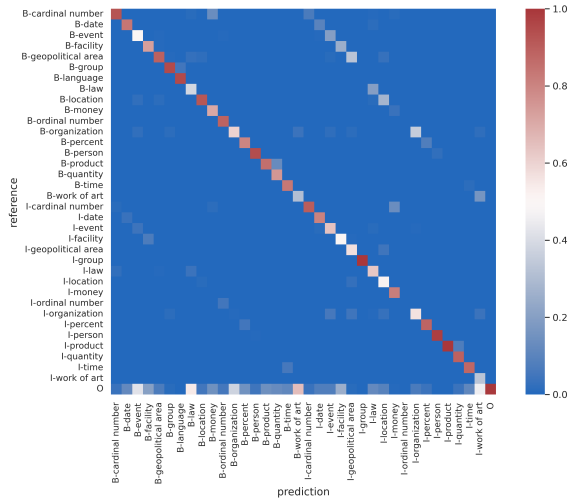


Figure 4: Confusion matrix, normalized to show the probability distribution of the tags predicted with the text-based model.

the symbol “%”, and sometimes the word “thousands” (in all languages). Consequently, for all entities marked as cardinal number, we added the missing tokens when necessary, following the rules specific to the language<sup>2</sup>. Another error often made by the text-based NER model is to predict the article as being part of the entity. As multiple sources advocate against doing so, we abided by the main guidelines (Maekawa, 2018; Benikova et al., 2014).

#### 4.5. Distribution

The annotated datasets are distributed in two formats: As JSON Lines files available on GitHub<sup>3</sup>, and on the HuggingFace repository (Wolf et al., 2020). There is one file per subset and per language, where each line is an annotated example. The audio files can be obtained by downloading VoxPopuli and matching the audio ID. The dataset version hosted on HuggingFace contains the audio

<sup>2</sup>In French and in Spanish, the symbol “%” is generally used, but in German and in Dutch, the word is more commonly spelled as Prozent or procent, respectively.

<sup>3</sup><https://github.com/qmeeus/MSNER>

recordings and the preprocessed annotations in BIO format, so that a researcher can already use the dataset after only two lines of code.

## 5. Evaluation Metrics

Following Shon et al. (2021), we recommend evaluating model predictions with the micro-averaged F1-score. The F1-score is the harmonic mean of precision and recall, calculated from an unordered list of named entities predicted for each utterance. Precision is the proportion of correctly predicted entities among all predicted entities, and recall is the proportion of ground truth entities that were correctly detected. An entity is considered to be predicted correctly if both the type and spelling are identical to the ground truth. To allow multiple entities with the same spelling and type in a sentence, we add a unique identifier to each entity/type pair. We recommend using the micro-averaged F1-score because the dataset is unbalanced. The label F1-score only considers the predicted type of the entity for correctness, leaving the transcribed entity out of the computations. This metric ignores spelling mistakes and segmentation errors. We provide an evaluation script<sup>4</sup> to compute these metrics and generate a breakdown of the prediction results per entity type.

## 6. Experiments

### 6.1. Setup

The first analysis compares the annotated test sets to the pseudo-annotations generated by the text-based NER model. Since the silver-label training and validation sets were generated with this model, this analysis is valuable for anyone intending to use these datasets for training. Indeed, it gives insights into the entities that are often confused with one another or remain undetected. It also gives some insights on the reliability of the model’s confidence score in assessing whether a prediction is correct.

<sup>4</sup><https://raw.githubusercontent.com/qmeeus/MSNER/main/src/evaluate.py>

We also consider two methods to predict named entities from speech, with a pipeline and an end-to-end model. The end-to-end model is a transformer encoder-decoder trained to perform both ASR and NER with a multitask objective (Meeus et al., 2023). This model is initialized from Whisper Large V2 (Radford et al., 2022), with an additional SLU module connected to the layers of the decoder with an adaptor. The end-to-end model was fine-tuned on English SLUE-VoxPopuli (Shon et al., 2021). The pipeline model transcribes the audio files and subsequently annotates the transcriptions. For the ASR model, we use Whisper Large V2 (Radford et al., 2022). For the pipeline model, we provide two options to allow for a better comparison. In Table 3, we use XML-RoBERTa fine-tuned on OntoNotes v5 (Weischedel et al., 2013) and compare it to the predictions generated by the text-based NER model from the gold transcripts. In Table 4, we fine-tuned the same XML-RoBERTa on SLUE-VoxPopuli (Shon et al., 2021), which provides a fair comparison to the end-to-end model. Although both models rely on multilingual pre-trained models, the fine-tuning dataset is entirely in English. Therefore, we evaluate the ability of these models to generalize from one language (English) to other languages (Dutch, French, German, and Spanish). Before computing the F1-scores, we normalize the text by putting it in lower case and removing symbols. It should be noted that the evaluation script does normalize the text further, which could have its importance depending on the model to be evaluated.

All results are presented on the human-annotated test sets proposed in this article.

## 6.2. Results

Figure 3 shows the distribution of calculated probabilities for predicted ‘B’ and ‘O’ tags conditional to whether they were predicted correctly or not. For each token position  $k$ , the probability of the most likely tag  $i^*$  is computed as follows:

$$P(y^k = i^*) = \max_i \frac{e^{z_i^k}}{\sum_j e^{z_j^k}}$$

where  $z_{1..N}^k$  are the logits predicted by the model for the token at position  $k$ . We observe that, on average, annotations for which there was no agreement between the annotator and the NER model were predicted with a lower probability than annotations that were correctly annotated from the start. However, we observe major differences between the class distributions. For the most frequent classes, like ‘O’, ‘organization’ or ‘date’, the probability distributions overlap considerably, and one should be careful if using this score as a proxy for the model’s uncertainty. This is not surprising, as transformers

are known to be overconfident (Ye et al., 2023). For rare quantitative classes like ‘percent’ and ‘quantity’, the model shows confidence when predictions are correct, and uncertain otherwise. This indicates that for those particular classes, the given probability could be relied upon when estimating the model’s uncertainty. The score breakdown by entity and language (Figure 2) indicates that in general, there are no major differences across languages, except for rare classes, where the variability increases significantly.

Figure 4 shows the confusion matrix of the NER model predictions against the manual annotations. Most errors are undetected entities (bottom row in Figure 4) and segmentation errors (I-tags predicted instead of B-tags and inversely, are visible on the lighter diagonals above and below the main diagonal). Some entities remain undetected more often than not, e.g. “work of art” and “event”, which is a sign that predictions are less reliable for these rare classes. Some other types are often confused with one another, like “money” and “cardinal number”. However, all types seem to have at most two confused types. We notice that “geopolitical area” is most often confused with “location” and “law”. In the latter case, this is because many laws are named after cities (e.g. the Paris Agreement, the Warsaw Treaty).

Table 3 compares the text-based NER predictions with the NER predictions obtained from the ASR transcript and generated by the same text-based NER model. The OntoNotes dataset, although in English, provides many well-curated annotations and the NER model trained on this dataset seem to generalize well to the other languages. However, this model was not trained to handle automatic transcripts and we observe a considerable drop in performance when it is asked to process ASR outputs. To make a fair comparison with the end-to-end model, we fine-tune XML-RoBERTa on SLUE-VoxPopuli and report the results in Table 4. The fine-tuning dataset being of much modest size (14.5 hours of training data), the models do not have many examples to learn from. The end-to-end model has a slight advantage because it learns simultaneously the ASR and NER tasks, and it is able to share part of its architecture between both tasks. For example, it seems well able to identify the presence of entities despite a lot of transcription and segmentation errors, as evidenced by the large label F1-score. In contrast, the pipeline suffers much more from the transcription errors because it was pretrained on curated texts and is not expecting noisy ASR transcriptions.

The text-based NER model performs best for Dutch, then German, French and finally Spanish. As the model was trained on English annotations, this ranking is not a surprise, although the ability of the

model to transfer to other languages is impressive. However, for the speech processing models, the same conclusion cannot be drawn. The entity F1-score seem to be correlated with the word error rate, which is influenced by the availability of the different languages in the pretraining set. In other words, for speech models, this is the model’s ability to transcribe foreign languages that will drive the quality of the predictions, rather than how similar the evaluation and the pretraining language are. The label-F1 indicates how accurate a model is at detecting the presence of entity types, disregarding of its ability to transcribe it correctly. Looking at those numbers, we observe again the same behavior as with the text-based entity predictions, namely that entities are more likely to be accurately detected when the evaluation language is more similar to the finetuning language.

Model	Metric	DE	ES	FR	NL
Gold	F1 (↑)	77.4	70.1	71.1	79.9
	Label-F1 (↑)	89.7	90.3	89.1	94.4
	F1 (↑)	52.4	50.6	44.7	52.7
ASR	Label-F1 (↑)	66.2	63.6	59.4	66.1
	WER (↓)	12.0	8.6	11.1	13.1

Table 3: Performance of text-based NER model trained on OntoNotes. Gold corresponds to the model’s predictions from the gold transcripts and ASR corresponds to the model’s predictions on the ASR transcripts.

Model	Metric	DE	ES	FR	NL
Pipeline	F1 (↑)	30.8	36.3	37.2	36.3
	Label-F1 (↑)	42.7	51.6	49.5	45.9
	WER (↓)	12.0	8.6	11.1	13.1
End2End	F1 (↑)	38.3	41.3	39.6	31.2
	Label-F1 (↑)	76.8	77.1	78.3	78.4
	WER (↓)	13.3	10.5	14.5	18.2

Table 4: Provided baselines on the annotated test sets for a pipeline ASR/NER model and an end-to-end multitask model. Both models were fine-tuned on SLUE-VoxPopuli (Shon et al., 2021)

## 7. Conclusion

In this manuscript, we have presented MSNER, a new dataset for evaluating multilingual Spoken NER systems. Although NER is a popular topic in NLP, this task has remained mostly unexplored in speech processing and spoken language understanding. To address this issue, we have used a pretrained model to annotate the VoxPopuli training and validation subsets in Dutch, French, German, and Spanish. Additionally, to provide researcher with a gold standard dataset for evaluating their

Spoken NER models, the authors have manually annotated the test sets for these subsets. By analyzing the predictions of a text-based NER model, and comparing them with our annotations, we were able to identify points of attentions for researchers who intend to train a model on silver annotations. For example, in some cases, the model confidence on the predictions can serve as a basis to estimate the correctness of the prediction, but this must be done carefully, since we have seen that transformers can be overconfident. Counter-intuitively, we have shown that most frequent classes are not always the ones where the model’s uncertainty is most reliable. We also looked at the classes that were often confused with one another, which gave us some ideas about which errors might be present in the training and validation sets.

We also provide baselines on the newly annotated evaluation subsets. We selected a pipeline and an end-to-end SLU model, both fine-tuned on English SLUE VoxPopuli (Shon et al., 2021), and we evaluate them on the manually annotated test sets. We saw that in a low resource scenario, the end-to-end model seems to benefit from learning simultaneously to transcribe and to annotate, which allows a better generalization across languages than the pipeline model fine-tuned on the same dataset. Finally, we found that the performance of text-based models on unseen languages is correlated with the similarity of the evaluation language with English. However, for speech models, this is the multilingual transcription accuracy that is the main driver for NER performance. Interestingly, we have seen that the end-to-end model was able to identify the presence of entities much better than the pipeline model, despite a similar overall performance, which illustrate the advantage of sharing parameters across tasks.

## 8. Bibliographical References

- Emanuele Bastianelli, Andrea Vanzo, Pawel Swietojanski, and Verena Rieser. 2020. [SLURP: A spoken language understanding resource package](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*. Association for Computational Linguistics.
- Darina Benikova, Chris Biemann, and Marc Reznicek. 2014. [NoSta-D named entity annotation for German: Guidelines and dataset](#). In *Proceedings of the Ninth International Conference on Language Resources and Evaluation (LREC'14)*, pages 2524–2531, Reykjavik, Iceland. European Language Resources Association (ELRA).
- Alexis Conneau, Kartikay Khandelwal, Naman Goyal, Vishrav Chaudhary, Guillaume Wenzek, Francisco Guzmán, Edouard Grave, Myle Ott, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Unsupervised cross-lingual representation learning at scale. *CoRR*, abs/1911.02116.
- Gamal Crichton, Sampo Pyysalo, Billy Chiu, and Anna Korhonen. 2017. [A neural network multi-task learning approach to biomedical named entity recognition](#). *BMC Bioinformatics*, 18.
- Charles T. Hemphill, John J. Godfrey, and George R. Doddington. 1990. The ATIS spoken language systems pilot corpus. In *Speech and Natural Language: Proceedings of a Workshop Held at Hidden Valley, Pennsylvania*.
- 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*. Morgan Kaufmann Publishers Inc.
- Loren Lugosch, Piyush Papreja, Mirco Ravanelli, Abdelwahab Heba, and Titouan Parcollet. 2021. Timers and such: A practical benchmark for spoken language understanding with numbers. *CoRR*, abs/2104.01604.
- Loren Lugosch, Mirco Ravanelli, Patrick Ignoto, Vikrant Singh Tomar, and Yoshua Bengio. 2019. Speech Model Pre-Training for End-to-End Spoken Language Understanding. In *Interspeech*.
- Emi Maekawa. 2018. [Annotation guidelines for named entities](#). online.
- Quentin Meeus, Marie-Francine Moens, and Hugo Van Hamme. 2023. [Whisper-slu: Extending a pretrained speech-to-text transformer for low resource spoken language understanding](#). In *2023 IEEE Automatic Speech Recognition and Understanding Workshop (ASRU)*.
- Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. 2022. Robust speech recognition via large-scale weak supervision. *CoRR*.
- Lance Ramshaw and Mitch Marcus. 1995. Text chunking using transformation-based learning. In *Third Workshop on Very Large Corpora*.
- Vincent Renkens and Hugo Van hamme. 2018. [Capsule networks for low resource spoken language understanding](#). In *Proc. Interspeech*. International Speech Communication Association.
- Alaa Saade, Alice Coucke, Alexandre Caulier, Joseph Dureau, Adrien Ball, Théodore Bluche, David Leroy, Clément Doumouro, Thibault Gisselbrecht, Francesco Caltagirone, Thibaut Lavril, and Mael Primet. 2018. Spoken language understanding on the edge. *CoRR*.
- Suwon Shon, Ankita Pasad, Felix Wu, Pablo Brusco, Yoav Artzi, Karen Livescu, and Kyu J. Han. 2021. SLUE: new benchmark tasks for spoken language understanding evaluation on natural speech. *CoRR*.
- Simone Tedeschi, Simone Conia, Francesco Cecconi, and Roberto Navigli. 2021. [Named entity recognition for entity linking: What works and what's next](#).
- Simone Tedeschi and Roberto Navigli. 2022. [Multi-NERD: A multilingual, multi-genre and fine-grained dataset for named entity recognition \(and disambiguation\)](#). In *Findings of the Association for Computational Linguistics: NAACL 2022*, Seattle, United States. Association for Computational Linguistics.
- Erik F. Tjong Kim Sang and Fien De Meulder. 2003. Introduction to the CoNLL-2003 shared task: Language-independent named entity recognition. In *Proceedings of the Seventh Conference on Natural Language Learning at HLT-NAACL 2003*.
- Asahi Ushio and Jose Camacho-Collados. 2021. T-NER: An all-round python library for transformer-based named entity recognition. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: System Demonstrations*. Association for Computational Linguistics.
- Özlem Uzuner, Imre Solti, Fei Xia, and Eithon Cadag. 2010. [Community annotation experiment](#)



for ground truth generation for the i2b2 medication challenge. *Journal of the American Medical Informatics Association*, 17(5).

Changhan Wang, Morgane Riviere, Ann Lee, Anne Wu, Chaitanya Talnikar, Daniel Haziza, Mary Williamson, Juan Pino, and Emmanuel Dupoux. 2021. VoxPopuli: A large-scale multilingual speech corpus for representation learning, semi-supervised learning and interpretation. Association for Computational Linguistics.

Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In *EMNLP. ACL*.

Wenqian Ye, Yunsheng Ma, Xu Cao, and Kun Tang. 2023. Mitigating transformer overconfidence via Lipschitz regularization. In *Proceedings of the Thirty-Ninth Conference on Uncertainty in Artificial Intelligence*, volume 216 of *Proceedings of Machine Learning Research*, pages 2422–2432. PMLR.

## 9. Language Resource References

Ralph Weischedel and Martha Palmer and Mitchell Marcus and Eduard Hovy and Sameer Pradhan and Lance Ramshaw and Nianwen Xue and Ann Taylor and Jeff Kaufman and Michelle Franchini and Mohammed El-Bachouti and Robert Belvin and Ann Houston. 2013. *OntoNotes Release 5.0*. Linguistic Data Consortium LDC2013T19, ISLRN 151-738-649-048-2.