

# For a Fistful of Puns: Evaluating a Puns in Multiword Expressions Identification Algorithm Without Dedicated Dataset

Anonymous ACL submission

## Abstract

Machine Translation systems has always faced challenges such as multiword expressions (MWEs) and wordplays, which impact their performances, being idiosyncratic and pervasive across different languages. In this context, we seek to explore the nature of puns created from multiword expressions (PMWEs), characterized by the creation of a wordplay from a source MWE to recontextualize it or to give it a humorous touch. Little work has been done on these entities in NLP. To address this challenge, we introduce ASMR, an alignment-based PMWE identification and tagging algorithm. We offer an in-depth analysis of three different approaches to ASMR, each created to identify different types of PMWEs. In the absence of PMWE-related datasets and resources, we proceed to a snowclone detection task in English. We also perform a MWE identification task in 26 languages to evaluate ASMR performances across different languages. We show that ASMR exhibits state-of-the-art results for the snowclone detection task and produces interesting results with the MWE identification task. These results may indicate that ASMR is suitable for a PMWE identification task.

## 1 Introduction

A lot of work has been done on multiword expressions (MWEs) in NLP since their introduction to the field by Sag et al. (2002); Choueka (1988). They are generally described as combinations of words with a certain degree of idiomaticity at the lexical, syntactic, semantic, pragmatic and/or statistical levels (Baldwin and Kim, 2010). MWEs are usually non-compositional or semi-compositional (Gross, 1982), diosyncratic, pervasive across different languages, and subject to varying degrees of variation. (Ramisch, 2023). Other phenomena, such as ambiguity and discontiguity, may also be an issue (Constant et al., 2017). Because of these features, they represent a particular

challenge in NLP, notably for Machine Translation systems, which need to take them into account (Zaninello and Birch, 2020).

Like any sequence of words, MWEs can serve as the basis for creating puns and other kinds of wordplays. Puns in multiword expressions (hereafter PMWEs) are characterized by the creation of a pun or a wordplay from a source multiword expression in order to recontextualize it or give it a humorous touch. By this process, MWEs such as (1) become (2) in the context of strikes in France in 2023.

1. "*l'heure est grave*"  
(FR, it's a **serious** time)
2. "*l'heure est grève*"  
(FR, it's a **strike** time)

Like MWEs, PMWEs can be translated from one language to another. For instance, (4) is a PMWE created from (3) working in both Italian and English. However, studies show that the translation of puns is not well handled by Machine Translation systems (Yu et al., 2018; Jiang et al., 2021).

3. "*l'alba dei morti viventi*"  
(IT, **Dawn** of the dead, 1978)
4. "*l'alba dei morti dementi*"  
(IT, **Shaun** of the dead, 2004)

To our knowledge, PMWEs have not been extensively studied in NLP, with very few resources available and almost no dedicated work on them. We find that PMWEs can be interesting due to their dual nature as MWE and wordplay. Machine Translation tasks as well as Automatic Humor Analysis could benefit from their study. Moreover, PMWEs might be useful to study the morphosyntactic and semantic evolutions of MWEs, since they tend to accept new forms and/or meaning over time (Fiala and Habert, 1989). In some cases, they may even be completely replaced by one of their own PMWE (Cusimano, 2015).

In addition to sharing the same difficulties as MWEs, PMWEs pose challenges of their own. Their identification in text can be even harder than that of MWEs, for several reasons: (i) they tend to be less frequent in texts than MWEs (ii) although their source MWE generally remains recognizable, several letters or words may be modified when creating a PMWE and (iii) their meaning can be altered, making the use of semantic-based approaches more challenging. Finally, differentiating a PMWE from a MWE can be a complex task, even for an individual with a certain expertise in these entities, as shown in [self reference \(1000\)](#).

In this paper, we introduce ASMR (Align, Segment, Match, Rank), an alignment-based algorithm whose goal is to identify and tag PMWE candidates in texts. We first present the architecture of ASMR. We then proceed to various experiments in two different datasets in order to evaluate the performances of this algorithm:

**Snowclone detection** For the first series of experiments, we evaluate how ASMR is able to detect snowclones (defined in Section 2) in a given set of sentences. To do so, we use the CATCHPHRASE dataset ([Sweed and Shahaf, 2021](#)).

**MWE identification** For the second series of experiments, we aim to evaluate ASMR’s ability to identify and tag MWEs in different languages by using the PARSEME 1.3 corpus ([Savary et al., 2023](#)), which consists of 26 languages.

With the help of an older prototype version of ASMR, we were able to identify PMWEs created from 216 MWEs in a corpus of French tweets ([self reference, 1000](#)). We were also able to identify a set of PMWEs created from formulas in Middle Arabic texts ([self reference, 1000](#)). Both approaches rely on qualitative evaluation carried out by experts on a selection of  $N$  PMWE candidates. In the absence of a PMWE annotated dataset, we have not yet been able to evaluate the performances of ASMR from a quantitative perspective. By combining a snowclone detection task with a MWE identification and tagging task, we hope to gain a better insight into ASMR’s functionalities.

## 2 Related Work

**MWE identification.** The main focus of MWE processing in NLP is the identification task, whose goal is to tag MWEs from a lexicon or a list in a text. Direct string matching and rule-based meth-

ods such as the ones proposed by [Stanković et al. \(2016\)](#); [Ramisch \(2015\)](#) were the first approaches used to address this task and are still used to this day. More recent approaches use Large Language Models (LLMs) such as BERT ([Devlin et al., 2019](#)). In fact, LLMs-based methodologies tend to outperform other approaches for the task of MWE identification ([Ramisch et al., 2020](#); [Bui and Savary, 2024](#)). For instance, [Tanner and Hoffman \(2023\)](#) use a rule-based pipeline along with a pretrained Bi-encoders for Word Sense Disambiguation ([Blevins and Zettlemoyer, 2020](#)). [Taslimipoor et al. \(2020\)](#) use a pretrained BERT model as well as a tree CRF architecture to tag verbal MWEs in the PARSEME 1.2 corpus. [Swaminathan and Cook \(2023\)](#) use multilingual LLMs to try to learn non-language-specific knowledge about MWEs and idiomaticity. Nevertheless, while pretrained LLMs seem to offer better results than more traditional approaches, they still have difficulties capturing their semantic aspect ([Tayyar Madabushi et al., 2021](#); [Zeng and Bhat, 2022](#)). [Wada et al. \(2023\)](#) paraphrase MWEs to address this problem, demonstrating that taking into account relevant semantic information can help to identify MWEs. Since there are very few resources on PMWEs, approaches using language models seem all the more costly to implement. We therefore drew inspiration from rule-based approaches to design ASMR, using known properties of PMWEs to characterize and identify them. We also plan to implement some semantic information in our methodology.

**Wordplays Detection.** Linguistic creativity, and therefore wordplays, are hard to deal with in NLP. As explained by [Netzer et al. \(2009\)](#); [Saussure et al. \(1949\)](#), humans tend to diversify their sets of relations between words, using cultural and emotional experiences for instance. As a result, the combinatorial possibilities for creating wordplays are almost infinite ([Knospe et al., 2016](#)). Few works report on wordplays detection. However, Since 2022, the JOKER-CLEF participative task challenge teams of scientists on several wordplay detection tasks ([Ermakova et al., 2022, 2023, 2024](#)).

**Snowclones Detection.** A snowclone is generally illustrated by a prototypical form of a MWE with flexible positions ("X be the new Y", X and Y being the flexible positions). It is derived from a reference sentence ("pink is the new black", allegedly said by Gloria Vanderbilt in India, 1960) and used to create new forms ("orange is the new black", Netflix

TV show, 2013). Snowclones have known a large set of definitions, often described as patterns that accept word substitutions (Lieberman, 2006), taking up known and institutionalized MWEs that remain identifiable in all circumstances (Hill, 2018; Traugott et al., 2016). Hartmann and Ungerer (2023) propose a quantitative study of two snowclones, "X be the new Y" and "the mother of all X", by extracting new forms of these snowclones. Sweed and Shahaf (2021) introduce the CATCHPHRASE dataset, constituted of 3,855 snowclone-sentence pairs, along with a snowclone detection methodology relying on an SVM-based approach and a RoBERTa-based approach (Liu et al., 2019) (Section 4). While snowclones tend to be PMWEs, there is no saying that all PMWEs are snowclones. Snowclones correspond to patterns with predefined word substitution positions, but we argue that PMWEs do not necessarily comply with this rule.

### 3 Introduction to ASMR

ASMR main purpose is the identification and tagging of PMWEs (Puns created from Multi Word Expressions). It can be described as an alignment-based, semi-supervised approach. ASMR takes a list of seeds, for instance prototypical forms of MWEs, as described in Pasquer (2019), and a list of sentences in which we want to identify PMWEs created from the seeds. As an output, ASMR create a ranking of PMWE candidates for each seed. It consists of a succession of 4 processes, which we describe here.

#### 3.1 Alignment

First, ASMR creates alignments between each seed-sentence pairs. An alignment can be defined as the superposition of the elements of two sequences in order to highlight their similarities and differences. We give an example of alignment between two sequences in Table 1.

May	the	-	beer	be	with	you
May	the	force	-	be	with	you

Table 1: Example of alignment at token level for the seed "May the force be with you" and the PMWE "May the beer be with you" (CATCHPHRASE dataset). In this example, the substitution of "force" by "beer" is highlighted by the misalignment between these tokens (in blue).

We use the BIOPYTHON package<sup>1</sup> to create these alignments. Initially dedicated to the alignment of DNA and RNA sequences, this package offers a fast token-level alignment. Furthermore, BIOPYTHON allows us to fetch multiple possible alignments for a given seed-sentence pair, as shown in Table 2.

there	s	no	place	like	long	island	no	place	like	home
there	s	-	-	-	-	-	no	place	like	home
there	s	no	place	like	-	-	-	-	-	home

Table 2: Two possible alignments between the seed "there's no place like home" and a sentence seen in the CATCHPHRASE dataset.

#### 3.2 Segmentation

Once the alignments made, we use them to find the longest common segment (hereafter LCS) between a seed and a sentence. This LCS will be our PMWE candidate. To find the LCS, we perform the following steps: (i) we retrieve each aligned token between a seed and a sentence and (ii) for each misalignment, we create a list containing all consecutive misaligned tokens, both for the seed and the sentence. We plan to use these misalignment lists in the next step in order to match unseen tokens from the seed with substitute tokens from the corresponding sentence.

#### 3.3 Matching

The matching process's goal is to isolate the LCS between a seed and a sentence. We provided 3 approaches to match tokens from the seed with tokens from the sentence, leading to the creation of 3 different approaches to ASMR:  $ASMR_{exact}$ ,  $ASMR_{fuzzy}$  and  $ASMR_{combined}$ .

**$ASMR_{exact}$**  Only identical tokens between the seed and the sentence are matched. In other word, only the aligned tokens are matched, while the misaligned ones are ignored.

**$ASMR_{fuzzy}$**  We match every single token between the first and the last common tokens between the aligned seed and sentence. If the  $X$  first tokens of the seed are unseen in the sentence, we match the first  $X$  tokens before the first common token in the sentence. We repeat this process with the  $Y$  last tokens of the seed: if they are unseen in the sentence, we match the  $Y$  first tokens after the last common token in the sentence.

<sup>1</sup><https://biopython.org/>

Seed Sentence	some men just want to watch the world burn some people really do just want to watch the world freeze													
Alignment	some some	men -	- people	- really	- do	just just	want want	to to	watch watch	the the	world world	burn -	- freeze	
Segmentation	some some	[men] [people,really,do]				just just	want want	to to	watch watch	the the	world world	[burn] [freeze]		
Match <sub>exact</sub>	some					just	want	to	watch	the	world			
Match <sub>fuzzy</sub>	some		people	really	do	just	want	to	watch	the	world		freeze	
Match <sub>combined</sub>	some		people			just	want	to	watch	the	world		freeze	
Cand. <sub>exact</sub>	some just want to watch the world													0.86
Cand. <sub>fuzzy</sub>	some <b>people really do</b> just want to watch the world <b>freeze</b>													0.70
Cand. <sub>combined</sub>	some <b>people</b> just want to watch the world <b>freeze</b>													0.80

Table 3: Alignment, segmentation, matching and resulting candidate (Cand.) for each approach for the seed "some men just want to watch the world burn" paired with the sentence "some people really just do want to watch the world freeze", found in the CATCHPHRASE dataset. For each seed-candidate pair, a cosine similarity is computed to rank candidates.

**ASMR<sub>combined</sub>** In addition to matching the aligned tokens between the seed and the sentence, we use misalignment lists to find the closest match for each unseen token from the seed. Let's take the following lists  $list_{seed}$  and  $list_{sent}$  from Table 3:

- $list_{seed} = [\text{men}]$
- $list_{sent} = [\text{people, really, do}]$

For each token  $tok_{seed}$  from  $list_{seed}$ , we compare its POS tag with the ones of each of the tokens in  $list_{sent}$ . The first token from  $list_{sent}$  with the same POS tag is matched with  $tok_{seed}$ . If no token possesses the same POS tag as  $tok_{seed}$ , we compute a Levenshtein score between  $tok_{seed}$  and each token in  $list_{sent}$  in order to find the best match. The only word from  $list_{seed}$ , "men", would therefore match with the first token of  $list_{sent}$ , "people", since they share the same POS tag.

Table 3 show the alignment, segmentation and matching steps. Each approach was designed to provide a solution to a specific problem. ASMR<sub>exact</sub> can help us identify puns-free MWEs and provide a minimal tagging of MWEs. In contrast, it should not be able to find substitutes to unseen tokens in the seed, and therefore is most likely not suitable for PMWEs identification. ASMR<sub>fuzzy</sub>, on the contrary, should be able to identify PMWEs, especially insertion and substitution based PMWEs, but will most probably produce a significant amount of noise, as it does not take discontinuity into consideration. Finally, ASMR<sub>combined</sub> will try to match the exact number of words seen in the seed by matching unseen tokens with substitutes. However, it should not be able to identify insertions.

### 3.4 Ranking

Prior to this step, we aligned, segmented and matched each seed with each sentence. As a result, we obtain a certain number of PMWE candidates for each seed. The final step of ASMR is to rank these candidates in order to sort them according to their probability of corresponding to a PMWE. We choose to use a cosine similarity score (illustrated in 1) using SCIKIT-LEARN to vectorize and to rank the candidates for each seed.

$$s_c(\vec{u}, \vec{v}) = \frac{\vec{u} \cdot \vec{v}}{\|\vec{u}\| \|\vec{v}\|} \quad (1)$$

We compute a cosine similarity matrix between each seed  $u$  and all the PMWE candidates  $v$  extracted with this seed, as shown in 2.

$$M = \begin{bmatrix} s_c(\vec{u}_1, \vec{v}_1) & s_c(\vec{u}_1, \vec{v}_2) & \cdots & s_c(\vec{u}_1, \vec{v}_n) \\ s_c(\vec{u}_2, \vec{v}_1) & s_c(\vec{u}_2, \vec{v}_2) & \cdots & s_c(\vec{u}_2, \vec{v}_n) \\ \vdots & \vdots & \ddots & \vdots \\ s_c(\vec{u}_m, \vec{v}_1) & s_c(\vec{u}_m, \vec{v}_2) & \cdots & s_c(\vec{u}_m, \vec{v}_n) \end{bmatrix} \quad (2)$$

The ranking step can be repeated for numerous linguistic information layers. For instance, if our seeds and sentences are POS tagged, we can compute another similarity matrix between the POS tags of the seeds and the ones of the candidates. We argue that such process allows us to take into account various information in order to adjust our ranking of the candidates for each seed. In order to take all the available linguistic information layers into account at the same time, we calculate the mean similarity score of all layers for each candidate. Finally, we ponder our scores by taking into account the difference of length  $N$  between the seed and the candidate: if the candidate has



	Recall	Precision	F-score	Accuracy
ASMR <sub>exact</sub>	<b>.89±.06</b>	.73±.14	.79±.09	<b>.89±.10</b>
ASMR <sub>fuzzy</sub>	.88±.03	.81±.09	<b>.84±.05</b>	.88±.03
ASMR <sub>combined</sub>	<b>.89±.02</b>	.80±.06	<b>.84±.03</b>	<b>.89±.02</b>
SVM (Sweed and Shahaf, 2021)	.78±.12	<b>.84±.13</b>	.81±NA	.85±.08
ROBERTA (Sweed and Shahaf, 2021)	.74±.18	.70±.15	.72±NA	.81±.94

Table 4: Results of ASMR for snowclone detection on the CATCHPHRASE test set. For the results of our approaches, the standard deviation is computed on 20 runs. Additionally, we manually computed F-scores for SVM and ROBERTA since (Sweed and Shahaf, 2021) did not report them.

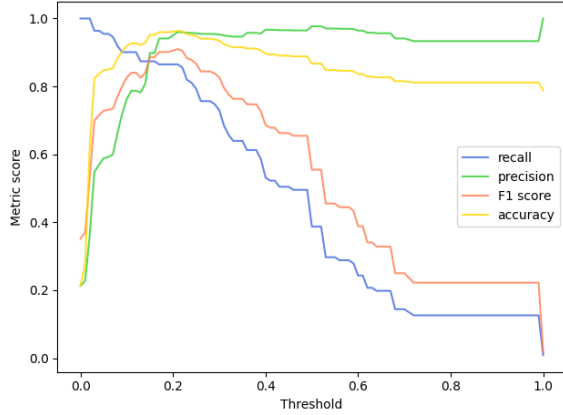


Figure 1: Impact of threshold on recall, precision, F-score and accuracy for the best run on the test partition of CATCHPHRASE with ASMR<sub>combined</sub>.

$X$  fewer tokens than the seed, we apply a rule of proportionality to its score  $S$ , as in 3.

$$S_{ponder} = \frac{S \cdot (N - X)}{N} \quad (3)$$

By applying this rule, we aim to discriminate candidates shorter than their seed, as a lot of them tend to be false positive. Additionally, shorter candidates that partially match the words of a seed tend to have better cosine similarity scores when compared with a seed, as seen for Cand.<sub>exact</sub> in Table 3.

## 4 Snowclone detection

We explained the features of ASMR. We now use the CATCHPHRASE dataset (Sweed and Shahaf, 2021) to evaluate ASMR capacity to detect if a sentence contains a snowclone.

**Dataset.** The CATCHPHRASE dataset consists of 3,855 snowclone-sentence pairs, of which 1,406 sentences allegedly contain the snowclone it was paired with. It proposes a binary classification task: for each snowclone-sentence pair, we must indicate whether the sentence contains the snowclone it was paired with. To achieve this classification

task, (Sweed and Shahaf, 2021) used a Feature-based SVM model as well as a ROBERTA-based model. We report the recall, precision and accuracy they obtained with these models in Table 4. Surprisingly, their SVM model performed better than their ROBERTA model.

**Parameters.** As ASMR does not learn from input data, we use the train and dev partitions of CATCHPHRASE to determine the best parameters to run our experiments. We run ASMR with 240 distinct sets of parameters on the train partition. These parameters include those of the vectorizer (number of ngrams and analyzer) and the threshold at which we consider a candidate to be a snowclone (according to its score). We only use token-level information during these runs. We select the 10 best sets of parameters for the train partition and the 10 best sets for the dev partition, for a total of 20 sets. We plan to run ASMR on the test partition with these 20 sets of parameters and to report standard deviation. We repeat this process for each approach, ASMR<sub>exact</sub>, ASMR<sub>fuzzy</sub> and ASMR<sub>combined</sub> for a total of 720 runs.

**Results.** We report the results we obtained with ASMR with each approach in Table 4. ASMR<sub>exact</sub> obtained the best recall and accuracy. ASMR<sub>fuzzy</sub> offer the best precision among the 3 tested algorithms, as well as the best F-score. ASMR<sub>combined</sub> achieve the best recall, F-score and accuracy. This matching algorithm might be slightly better than the other two due to the nature of snowclones, which are mainly created by substitution. Overall, ASMR performs better than the models used by Sweed and Shahaf (2021) for the task of snowclone detection, although we note that our accuracy is slightly behind that obtained by their SVM model. Figure 1 shows the impact of threshold on the metrics we used for the best run (with ASMR<sub>combined</sub>). As expected for a ranking system, the lower the threshold, the lower the precision and the higher the recall.

	ASMR <sub>exact</sub>			ASMR <sub>fuzzy</sub>			ASMR <sub>combined</sub>			s2s
	R	P	F	R	P	F	R	P	F	F
AR	32.2±02	<b>54.0±01</b>	<b>40.3±01</b>	25.4±08	40.3±12	30.4±09	<b>34.0±06</b>	40.6±11	35.3±03	50.9
BG	<b>72.1±01</b>	55.3±00	<b>62.5±00</b>	61.9±11	56.8±04	58.4±04	63.3±10	<b>57.4±03</b>	59.5±03	65.7
CS	59.4±00	<b>64.9±00</b>	<b>62.0±00</b>	46.4±11	57.5±08	51.0±09	<b>60.0±03</b>	59.0±05	59.3±02	74.1
DE	<b>20.7±00</b>	<b>67.3±03</b>	<b>31.6±00</b>	16.4±03	38.8±19	22.2±06	18.6±01	43.3±15	25.5±03	71.4
EL	<b>57.9±03</b>	57.3±01	<b>57.5±01</b>	44.4±13	55.4±05	48.4±08	55.4±07	<b>59.5±02</b>	56.9±03	66.3
EN	<b>44.4±01</b>	<b>78.0±00</b>	<b>56.5±00</b>	32.4±08	66.7±15	42.6±08	42.8±03	72.1±08	53.6±03	59.9
ES	<b>53.8±00</b>	<b>54.7±00</b>	<b>54.2±00</b>	45.3±08	49.9±05	47.3±06	50.2±05	51.5±04	50.6±03	55.6
EU	<b>72.7±01</b>	<b>76.4±03</b>	<b>74.4±01</b>	62.3±10	74.5±07	67.2±07	71.3±05	69.1±08	69.8±04	<u>82.1</u>
FA	61.8±00	77.8±01	68.8±00	64.0±03	<b>78.0±05</b>	70.1±01	<b>66.4±04</b>	76.5±02	<b>71.0±02</b>	71.9
FR	<b>66.2±04</b>	<b>73.6±01</b>	<b>69.6±02</b>	50.2±13	57.5±13	53.5±13	65.9±04	65.1±07	65.2±03	78.7
GA	<b>19.4±00</b>	<b>52.0±00</b>	<b>28.2±00</b>	17.2±06	49.3±13	23.9±05	<b>19.4±01</b>	51.6±07	28.0±00	26.6
HE	35.8±01	<b>64.1±01</b>	<b>45.9±00</b>	33.9±02	53.6±10	41.3±04	<b>36.3±01</b>	57.4±06	44.4±02	<u>46.9</u>
HI	45.2±00	<b>80.6±01</b>	57.9±00	<b>51.2±08</b>	75.4±12	<b>59.6±02</b>	46.4±01	70.7±07	55.9±02	58.7
HR	<b>64.1±01</b>	<b>91.9±00</b>	<b>75.5±01</b>	49.7±13	77.9±14	60.1±13	61.9±03	79.8±09	69.5±04	75.3
HU	<b>18.5±02</b>	<b>81.8±21</b>	<b>29.4±01</b>	15.8±03	69.3±16	25.2±03	18.4±02	76.0±18	28.9±01	32.0
IT	<b>59.0±01</b>	<b>64.0±01</b>	<b>61.4±00</b>	50.0±07	55.2±09	52.2±07	58.6±02	61.1±03	59.8±01	<u>65.0</u>
LT	27.5±00	<b>83.2±00</b>	<b>41.3±00</b>	20.2±05	65.1±15	30.7±07	<b>27.7±02</b>	78.1±05	40.9±02	48.9
MT	14.2±02	<b>19.2±01</b>	<b>16.3±01</b>	<b>16.0±04</b>	16.4±03	15.7±02	10.4±04	15.2±04	12.1±04	<u>16.5</u>
PL	<b>62.4±05</b>	<b>90.1±01</b>	<b>73.6±03</b>	52.3±11	80.6±11	62.9±11	60.1±06	77.8±10	67.4±05	<u>82.5</u>
PT	51.4±07	<b>70.0±07</b>	<b>58.5±04</b>	34.6±05	47.3±07	39.3±03	<b>53.4±10</b>	59.0±07	54.8±05	<u>74.0</u>
RO	<b>88.4±00</b>	<b>61.1±00</b>	<b>72.3±00</b>	69.3±17	53.8±07	60.0±10	83.8±07	54.7±05	66.0±05	<u>74.8</u>
SL	<b>51.2±04</b>	<b>33.2±01</b>	<b>40.2±01</b>	33.7±16	29.2±04	30.0±09	49.7±04	30.2±03	37.4±02	<u>41.8</u>
SR	37.8±01	<b>87.1±00</b>	<b>52.7±01</b>	34.5±05	74.9±14	46.9±06	<b>38.8±02</b>	79.0±10	51.6±01	62.0
SV	<b>29.2±01</b>	<b>80.8±03</b>	<b>42.8±01</b>	25.1±04	70.1±13	36.7±05	28.4±02	74.2±03	41.0±02	<u>82.2</u>
TR	<b>71.8±03</b>	<b>58.4±02</b>	<b>64.4±01</b>	67.8±06	57.8±04	62.2±03	65.8±08	53.4±04	58.5±05	65.0
ZH	22.0±00	40.5±00	28.4±00	20.0±01	<b>42.0±01</b>	27.1±00	<b>23.5±01</b>	39.2±01	<b>29.2±01</b>	35.0
M	<b>47.7</b>	<b>66</b>	<b>52.6</b>	40	57.4	44.8	46.6	59.7	49.7	60.1

Table 5: Global MWE-based results on the test set of PARSEME 1.3 for 26 languages using ASMR. We report recall (R), precision (P), F-score (F) and mean (M) for all languages. Since we performed 10 runs for each language for our approaches, we also report the standard deviation. For the sake of comparison, we add SEEN2SEEN (s2s) system results. We underline state-of-the-art results.

## 5 MWE identification

We measured the performance of ASMR for the task of snowclone detection in sentences with CATCHPHRASE. We now want to evaluate its ability to identify tokens belonging to MWEs in a given set of sentences.

**Dataset.** We use the version 1.3 of the PARSEME corpus (Savary et al., 2023), composed of 26 languages and mainly containing verbal MWEs. This corpus proposes a MWEs tagging task. So far, only 2 systems have been tested on the PARSEME 1.3 corpus: SEEN2SEEN (Pasquer et al., 2020) and MTLB-STRUCT (Taslimipoor et al., 2020). Savary et al. (2023) report the results for these 2 systems on PARSEME 1.3.

**Parameters.** The following steps are repeated for each language: (i) We retrieve a list of every MWEs seen in the train partition (lemmas and POS tags included). Since we collected lemmas for each word, we use them to align each MWE with each sentence. (ii) We run ASMR with 256 sets of parameters on the dev partition. Those parameters consist of cosine similarity thresholds for the token

layer, the morphosyntactic layer and the lemma layer. The possible thresholds were 0.1, 0.3, 0.7 and 1. We also compute a semantic score between each candidate and MWE using the SENTENCE-TRANSFORMERS package. This addition will enable us to assess the impact of semantic information on a MWE identification task using ASMR. Additionally, we remove candidates with discontinuities of more than 4 words. As shown in Pasquer (2019), the vast majority of discontinuous MWEs tend to have shorter discontinuities. (iii) We select the 10 best sets of parameters for the dev partition to run them on the test partition. We repeat this process for each approach with ASMR, totaling 768 runs per language. In the end, we performed 19,968 runs on PARSEME 1.3.

**Results.** Table 5 shows the global MWE-based results we obtained on the test set of PARSEME 1.3 for each language. Overall, ASMR<sub>exact</sub> obtained the best results among all the ASMR approaches, with a mean F-score of 52.6. Since ASMR<sub>exact</sub> only tag aligned words between a seed and a sentence, this result does not come as a surprise. ASMR<sub>fuzzy</sub> offer the best F-score for Hindi

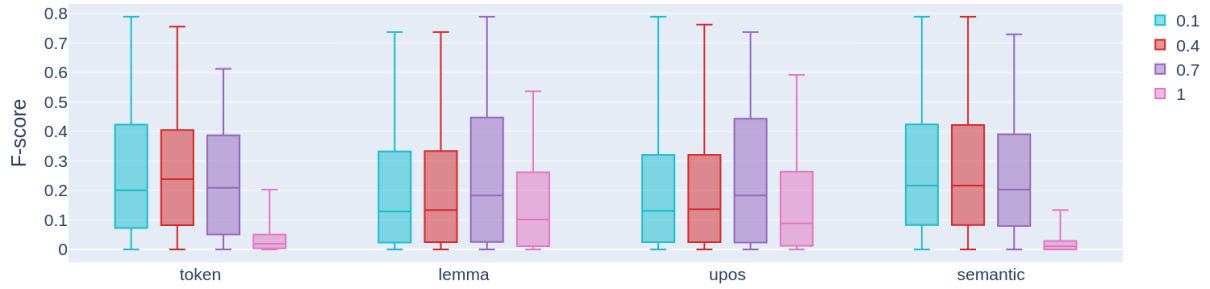


Figure 2: Boxplots of F-scores obtained on the dev partition of PARSEME 1.3 for each linguistic features (token, lemma, upos and semantic similarity) for different thresholds (1, 0.7, 0.4 and 0.1) and each language. We used the F-scores obtained with  $ASMR_{exact}$  since it has the best mean F-score among the 3 approaches we used.

(HI), while  $ASMR_{combined}$  obtained the best F-scores with Persian (FA) and Chinese (ZH). We report state-of-the-art results on the PARSEME 1.3 corpus with ASMR for Irish (GA) and Croatian (HR).

In order to analyze the impact of each feature used to identify MWEs with ASMR, we generate boxplots for each feature and each threshold used with these features in Figure 2. These boxplots consist of F-scores obtained with every run made with  $ASMR_{exact}$  on all languages on the dev partition. For instance, the first boxplot represents all the F-scores obtained with a threshold of 0.1 for the token feature. We observe that (i) regardless of the feature, a threshold of 1 seems to be too restrictive, as F-scores tend to be much lower, (ii) for the token and semantic features, we observe almost no variation with different thresholds, which can indicate that those features are not the most determinant for MWE identification with ASMR and (iii) the lemma and upos features show better F-scores with a threshold of 0.7, meaning that those features are probably the most helpful to identify MWEs with ASMR.

## 6 Error analysis with PARSEME

Since PARSEME 1.3 offers several metrics on different subsets of MWEs, such as discontinuous and unseen ones, we can perform a more refined analysis of ASMR capabilities. Table 6 shows the mean F-scores including all languages obtained with each approach on different subsets of MWEs. We observe that for two subsets (discontinuous and unseen-in-train) we achieve lower F-scores. Additionally, since ASMR was designed to identify

PMWEs, we could argue that the Variant-of-train score is lower than expected.

	Exact	Fuzzy	Combined
Tok-based	55.0	48.1	52.6
Continuous	57.1	52.3	54.8
<b>Discontinuous</b>	<b>41.9</b>	<b>14.8</b>	<b>38.4</b>
Seen-in-train	68.0	61.0	66.9
<b>Unseen-in-train</b>	<b>00.9</b>	<b>06.8</b>	<b>05.2</b>
<b>Variant-of-train</b>	<b>60.1</b>	<b>50.3</b>	<b>59.4</b>
Identical-to-train	78.6	72.8	76.4

Table 6: mean F-scores including all languages obtained with each approach on different subsets of MWEs. We highlight the most interesting subsets (in **bold**).

**Discontinuous** Discontinuous MWEs are a recurring challenge for MWE identification task (Constant et al., 2017). As  $ASMR_{fuzzy}$  and  $ASMR_{combined}$  try to match misaligned words between a MWE and a candidate, their low F-scores are expected. This is especially the case for  $ASMR_{fuzzy}$ , which match every word between the first and the last common words between a MWE and a candidate (as previously seen in Table 3). We observe that  $ASMR_{exact}$ , by tagging only aligned tokens, manage to obtain the highest mean F-score among the 3 approaches.

**Unseen-in-train** One could argue that ASMR should be able to see a minimal number of unseen-in-train MWEs, especially with the  $ASMR_{fuzzy}$  and  $ASMR_{combined}$  approaches. We argue that this can be the case, notably with shorter, more generic MWEs, such as "break up". Table 7 shows 10 candidates found with  $ASMR_{combined}$  for the MWE "break up". We observe the presence of other seen-in-train MWEs as well as 2 unseen-in-train MWEs. We also report erroneous candidates, which

does not correspond to a MWE. While ASMR is capable of capturing both closely related MWEs and unseen MWEs, it might be difficult for it to distinguish good candidates from bad ones. This is highlighted by the ranking in Table 7, where seen, unseen and erroned MWEs tend to blend together in the ranking.

Candidate	Cat	Tok	Upo	Lem	Sem	M
broke up	see	0.10	1.00	1.00	0.84	0.73
<b>speak up</b>	<b>uns</b>	0.55	1.00	0.48	0.45	0.62
<b>fuck up</b>	<b>uns</b>	0.20	1.00	0.21	0.44	0.46
look up	see	0.22	1.00	0.18	0.29	0.42
make up	see	0.12	1.00	0.11	0.41	0.41
ensure up	err	0.07	1.00	0.07	0.49	0.41
end up	see	0.03	1.00	0.03	0.54	0.40
jangle up	err	0.01	1.00	0.02	0.51	0.38
grow up	see	0.02	1.00	0.02	0.49	0.38
have up	err	0.02	1.00	0.03	0.46	0.38

Table 7: 10 ranked candidates for the MWE "break up". For each candidate, we report its score for each feature as well as its mean score (M, used for the ranking) and its subset (Cat). Possible categories are seen (see), unseen (uns) and erroned (err).

**Variant-of-train** Variants of MWEs can correspond to several instances in the PARSEME 1.3 (see guidelines<sup>2</sup>). Among these instances, we find (i) syntactic variants, such as conjugated verb, change of tense or number and (ii) MWEs with some open slots (to make/take a decision). The former should be handled by morphosyntactic and lemmas analysis in most case, but the latter may have a direct impact on MWE identification, especially with ASMR. Table 8 shows 10 candidates found with ASMR<sub>fuzzy</sub> for the MWE "получа помощ" (BG, get help). We observe possible variations for both words of this MWE. "получа" can be conjugated and/or replaced by "взе" and "помощ" can be replaced by "подкрепата". Once again, the possible variations of this MWE blend with erroned candidates in the ranking, making it hard to distinguish them. However, we observe that simple syntactic variants seems to obtain a higher score than other candidates, making it to the top of the ranking and therefore easier to identify.

## 7 Discussion

In this work, we introduced ASMR, a PMWE identification and tagging algorithm relying on sentence-level alignments and similarity scores in

<sup>2</sup>[https://parsemefr.lis-lab.fr/parseme-st-guidelines/1.3/?page=010\\_Definitions\\_and\\_scope/030\\_Syntactic\\_variants\\_of\\_VMWEs](https://parsemefr.lis-lab.fr/parseme-st-guidelines/1.3/?page=010_Definitions_and_scope/030_Syntactic_variants_of_VMWEs)

Candidate	EN	Cat	M
получат помощ	get help	var	0.90
получиха помощ	get help	var	0.84
получават помощ	get help	var	0.73
каза помощ	say help	err	0.65
каза помощта	say help	err	0.63
поеха помощ	ask help	err	0.63
взе помощ	take help	var	0.61
стана помощ	become help	err	0.61
получи подкрепа	receive support	var	0.54
получи подкрепата	receive support	var	0.54

Table 8: 10 ranked candidates for the MWE "получа помощ" (BG, get help). For each candidate, we propose a minimal translation in english (EN) as well as its mean score (M) and its category (Cat). Possible categories are identical (idt), variant (var) and erroned (err).

order to propose a ranking of PMWE candidates in a given set of sentences. While earlier studies show that ASMR can be used to extract good PMWE candidates in both French and Arabic (self reference, 1000), no quantitative evaluation was yet performed, due to a lack of a PMWE annotated dataset. To get around this issue, we proceeded to 2 experiments in order to evaluate ASMR functionalities. We first used a snowclone detection task on the CATCHPHRASE dataset in order to evaluate ASMR's capacity to assert the presence of a PMWE candidate in a sentence. We then used the PARSEME 1.3 corpus to evaluate ASMR identification and tagging performances on MWEs for 26 languages. We show that ASMR obtains state-of-the-art results on the snowclone detection task and for two languages with the MWE identification task (Irish and Croatian).

We performed an in-depth analysis of the limitations we encountered with some subsets of MWEs within PARSEME, which allowed us to get a better understanding of ASMR performances. This analysis has shown that, while true positive and false positive candidates tend to blend together in the ranking, the top  $N$  candidates seem to be pertinent in most cases. This observation is highlighted by both the MWE identification task and the snowclone detection task, where higher thresholds lead to higher precision and lower recall. We also note that, while we performed multiple run for each task, our standard deviations are low, which can account for the robustness of ASMR.

We plan to create a PMWE dataset through participative sciences to further evaluate the performances of ASMR. Such dataset would also be useful to test the performances of other systems, either created for MWE or PMWE identification.



## Limitations

**CATCHPHRASE experiment.** We take into account several limitations, due to either the CATCHPHRASE dataset or the methodology we used: (i) the dataset itself is imbalanced. As stated by its authors, 64 % of the sentences do not contain the snowclone they were paired with (Sweed and Shahaf, 2021). (ii) the task doesn’t evaluate the capacity of a system to tag tokens belonging to a snowclone. (iii) since CATCHPHRASE does not come with POS tag nor lemmas, we only tokenized both the snowclones and the sentences. (iv) the threshold itself can be seen as a limitation: the ideal threshold found for the train and dev partitions of the dataset might not always be the same for the test partition. Nevertheless, we find that for CATCHPHRASE, the ideal threshold seems to be roughly the same for all partitions (between 0.1 and 0.3).

**PARSEME 1.3 experiment.** To avoid overloading our calculation server, we had to limit the number of runs we made on the PARSEME 1.3 corpus. To limit this number, we did not manipulate the features of the vectorizer used to compute cosine similarity scores, which remained the same among all languages. We also limited to 4 the number of thresholds we used for each feature (using only thresholds of 0.1, 0.4, 0.7 and 1). Moreover, since ASMR was not initially designed to strictly identify MWEs, we added a rule to limit the size of possible discontinuities to 4. While this rule is also found in other systems, such as the one of Pasquer et al. (2020), we did not evaluate its impact on the MWE identification task with ASMR. Finally, ASMR does not account for phenomena such as permutation yet, which might have an impact on the results we obtained, since some MWEs allow word permutations.

## Ethical considerations

We ran ASMR on an AMD EPYC MODEL 7543P MILAN 32 CORE CPU with 32GB of memory. We ran it on every language in parallel threads, for a cumulated time of 58 hours and a maximum time of 13 hours. We use this information along with the carbon intensity in France in 2024<sup>3</sup> to estimate our carbon footprint, which amounts to 120.45g estimated CO2 emission (or 0.12 estimated

CO2 emission kilogram). This estimation remains approximate, as we couldn’t take every parameter into account.

In comparison, Large Language Models such as BERT usually have a much higher carbon footprint (Wang et al., 2023).

## Acknowledgments

We thank SACADO for allowing us to use their calculation server to run our experiments on the PARSEME 1.3 corpus.

## References

- Timothy Baldwin and Su Nam Kim. 2010. *Multiword Expressions*, 2 edition. Chapman and Hall/CRC.
- Terra Blevins and Luke Zettlemoyer. 2020. [Moving down the long tail of word sense disambiguation with gloss informed bi-encoders](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 1006–1017, Online. Association for Computational Linguistics.
- Van-Tuan Bui and Agata Savary. 2024. [Cross-type French multiword expression identification with pre-trained masked language models](#). In *Proceedings of the 2024 Joint International Conference on Computational Linguistics, Language Resources and Evaluation (LREC-COLING 2024)*, pages 4198–4204, Torino, Italia. ELRA and ICCL.
- Yaacov Choueka. 1988. Looking for needles in a haystack or locating interesting collocational expressions in large textual databases. In *User-Oriented Content-Based Text and Image Handling*, RIAO ’88, pages 609–623, Paris, FRA.
- Mathieu Constant, Gülşen Eryiğit, Johanna Monti, Lonneke van der Plas, Carlos Ramisch, Michael Rosner, and Amalia Todirascu. 2017. [Survey: Multiword Expression Processing: A Survey](#). *Computational Linguistics*, 43(4):837–892.
- Christophe Cusimano. 2015. [Figement de séquences défigées](#). *Pratiques*, (159-160):69 78.
- 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, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Liana Ermakova, Anne-Gwenn Bosser, Tristan Miller, Victor Preciado, Grigori Sidorov, and Adam Jatowt. 2024. [Overview of the CLEF 2024 JOKER Track: Automatic Humour Analysis](#), pages 165–182.

<sup>3</sup><https://www.sfen.org/rgn/2024-record-production-electricite/>

657	Liana Ermakova, Tristan Miller, Julien Boccou, Albin	Caroline Pasquer, Agata Savary, Carlos Ramisch,	711
658	Digue, Aurianne Damoy, and Paul Campen. 2022.	and Jean-Yves Antoine. 2020. <a href="#">Seen2Unseen at</a>	712
659	Overview of the clef 2022 joker task 2: translate	<a href="#">PARSEME shared task 2020: All roads do not lead</a>	713
660	wordplay in named entities. <i>Proceedings of the Work-</i>	<a href="#">to unseen verb-noun VMWEs</a> . In <i>Proceedings of the</i>	714
661	<i>ing Notes of CLEF</i> , pages 1666–1680.	<i>Joint Workshop on Multiword Expressions and Elec-</i>	715
662	Liana Ermakova, Tristan Miller, Anne-Gwenn Bosser,	<i>tronic Lexicons</i> , pages 124–129, online. Association	716
663	Victor Manuel Palma Preciado, Grigori Sidorov, and	for Computational Linguistics.	717
664	Adam Jatowt. 2023. Overview of joker-clef-2023	Carlos Ramisch. 2015. <i>Multiword Expressions Acqui-</i>	718
665	track on automatic wordplay analysis. In <i>Interna-</i>	<i>sition: A Generic and Open Framework</i> . Theory	719
666	<i>tional Conference of the Cross-Language Evalua-</i>	and Applications of Natural Language Processing.	720
667	<i>tion Forum for European Languages</i> , pages 397–415.	Springer International Publishing, Cham.	721
668	Springer.	Carlos Ramisch. 2023. <i>Multiword expressions in com-</i>	722
669	Pierre Fiala and Benoît Habert. 1989. <a href="#">La langue de bois</a>	<i>putational linguistics</i> . thesis, Aix Marseille Univer-	723
670	<a href="#">en éclat : les défigements dans les titres de presse quo-</a>	sité (AMU).	724
671	<a href="#">tidienne française</a> . <i>Mots. Les langages du politique</i> ,	Carlos Ramisch, Agata Savary, Bruno Guillaume,	725
672	21(1):83–99. 25 citations (Semantic Scholar/DOI)	Jakub Waszczuk, Marie Candito, Ashwini Vaidya,	726
673	[2022-11-15] Publisher: Persée - Portail des revues	Verginica Barbu Mititelu, Archana Bhatia, Uxoá Iñur-	727
674	scientifiques en SHS.	rieta, Voula Giouli, Tunga Güngör, Menghan Jiang,	728
675	Maurice Gross. 1982. <a href="#">Une classification des phrases</a>	Timm Lichte, Chaya Liebeskind, Johanna Monti,	729
676	<a href="#">« figées » du français</a> . <i>Revue québécoise de linguis-</i>	Renata Ramisch, Sara Stymne, Abigail Walsh, and	730
677	<i>tique</i> , 11(2):151.	Hongzhi Xu. 2020. <a href="#">Edition 1.2 of the PARSEME</a>	731
678	Stefan Hartmann and Tobias Ungerer. 2023. <a href="#">Attack of</a>	<a href="#">shared task on semi-supervised identification of ver-</a>	732
679	<a href="#">the snowclones: A corpus-based analysis of extrava-</a>	<a href="#">bal multiword expressions</a> . In <i>Proceedings of the</i>	733
680	<a href="#">gant formulaic patterns</a> . <i>Journal of Linguistics</i> , pages	<i>Joint Workshop on Multiword Expressions and Elec-</i>	734
681	1–36.	<i>tronic Lexicons</i> , pages 107–118, online. Association	735
682	Ian E. J. Hill. 2018. <a href="#">Memes, munitions, and collective</a>	for Computational Linguistics.	736
683	<a href="#">copia: The durability of the perpetual peace weapons</a>	Ivan A. Sag, Timothy Baldwin, Francis Bond, Ann	737
684	<a href="#">snowclone</a> . <i>Quarterly Journal of Speech</i> , 104(4):422–	Copestake, and Dan Flickinger. 2002. <a href="#">Multiword</a>	738
685	443.	<a href="#">Expressions: A Pain in the Neck for NLP</a> . In <i>Com-</i>	739
686	Siyu Jiang, Zhiheng Zhang, Qiong Zhong, Jin Xie, and	<i>putational Linguistics and Intelligent Text Process-</i>	740
687	Weilin Wu. 2021. <a href="#">The system analysis and research</a>	<i>ing</i> , Lecture Notes in Computer Science, pages 1–15,	741
688	<a href="#">based on pun recognition</a> . <i>Journal of Physics: Con-</i>	Berlin, Heidelberg. Springer.	742
689	<i>ference Series</i> , 2044(1):012190.	F. de Saussure, C. Bally, A. Riedlinger, and	743
690	Sebastian Knospe, Alexander Onysko, and Maik Goth.	A. Sechehayé. 1949. <i>Cours de linguistique générale</i> .	744
691	2016. <i>Crossing Languages to Play with Words: Mul-</i>	Payot, Paris.	745
692	<i>tidisciplinary Perspectives</i> . Walter de Gruyter GmbH	Agata Savary, Cherifa Ben Khelil, Carlos Ramisch,	746
693	& Co KG.	Voula Giouli, Verginica Barbu Mititelu, Najet	747
694	Mark Liberman. 2006. <a href="#">The proper treatment of snow-</a>	Hadj Mohamed, Cvetana Krstev, Chaya Liebeskind,	748
695	<a href="#">clones in ordinary english</a> .	Hongzhi Xu, Sara Stymne, Tunga Güngör, Thomas	749
696	Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Man-	Pickard, Bruno Guillaume, Eduard Bejček, Archana	750
697	dar Joshi, Danqi Chen, Omer Levy, M. Lewis, Luke	Bhatia, Marie Candito, Polona Gantar, Uxoá Iñurri-	751
698	Zettlemoyer, and Veselin Stoyanov. 2019. Roberta:	eta, Albert Gatt, Jolanta Kovalevskaite, Timm Lichte,	752
699	A robustly optimized bert pretraining approach.	Nikola Ljubešić, Johanna Monti, Carla Parra Escartín,	753
700	<i>ArXiv</i> , abs/1907.11692.	Mehnoush Shamsfard, Ivelina Stoyanova, Veronika	754
701	Yael Netzer, David Gabay, Yoav Goldberg, and Michael	Vincze, and Abigail Walsh. 2023. <a href="#">PARSEME corpus</a>	755
702	Elhadad. 2009. <a href="#">Gaiku : Generating haiku with word</a>	<a href="#">release 1.3</a> . In <i>Proceedings of the 19th Workshop on</i>	756
703	<a href="#">associations norms</a> . In <i>Proceedings of the Workshop</i>	<i>Multiword Expressions (MWE 2023)</i> , pages 24–35,	757
704	<i>on Computational Approaches to Linguistic Creativ-</i>	Dubrovnik, Croatia. Association for Computational	758
705	<i>ity</i> , pages 32–39, Boulder, Colorado. Association for	Linguistics.	759
706	Computational Linguistics.	self reference. 1000. self-reference paper title. <i>self-</i>	760
707	Caroline Pasquer. 2019. <a href="#">Garder la trace, mettre de</a>	<i>reference paper conference</i> .	761
708	<a href="#">l'ordre et relier les points : modéliser la variation et</a>	Ranka Stanković, Cvetana Krstev, Ivan Obradović, Bil-	762
709	<a href="#">l'ambiguïté des expressions polylexicales</a> . Phd thesis,	jana Lazić, and Aleksandra Trtovac. 2016. <a href="#">Rule-</a>	763
710	Tours, France.	<a href="#">based Automatic Multi-Word Term Extraction and</a>	764
		<a href="#">Lemmatization</a> . In <i>LREC</i> , pages 507–514, Portorož,	765
		Slovenia.	766

Raghuraman Swaminathan and Paul Cook. 2023. [Token-level identification of multiword expressions using pre-trained multilingual language models](#). In *Proceedings of the 19th Workshop on Multiword Expressions (MWE 2023)*, pages 1–6, Dubrovnik, Croatia. Association for Computational Linguistics.

Nir Sweed and Dafna Shahaf. 2021. [Catchphrase: Automatic Detection of Cultural References](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 2: Short Papers)*, pages 1–7, Online. Association for Computational Linguistics.

Joshua Tanner and Jacob Hoffman. 2023. [MWE as WSD: Solving Multiword Expression Identification with Word Sense Disambiguation](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 181–193, Singapore. Association for Computational Linguistics.

Shiva Taslimipour, Sara Bahaadini, and Ekaterina Kochmar. 2020. [MTLB-STRUCT @parseme 2020: Capturing unseen multiword expressions using multi-task learning and pre-trained masked language models](#). In *Proceedings of the Joint Workshop on Multiword Expressions and Electronic Lexicons*, pages 142–148, online. Association for Computational Linguistics.

Harish Tayyar Madabushi, Edward Gow-Smith, Carolina Scarton, and Aline Villavicencio. 2021. [AStitchInLanguageModels: Dataset and methods for the exploration of idiomaticity in pre-trained language models](#). In *Findings of the Association for Computational Linguistics: EMNLP 2021*, pages 3464–3477, Punta Cana, Dominican Republic. Association for Computational Linguistics.

Elizabeth Closs Traugott, Graeme Trousdale, Elizabeth Closs Traugott, and Graeme Trousdale. 2016. *Constructionalization and Constructional Changes*. Oxford Studies in Diachronic and Historical Linguistics. Oxford University Press, Oxford, New York.

Takashi Wada, Yuji Matsumoto, Timothy Baldwin, and Jey Han Lau. 2023. [Unsupervised paraphrasing of multiword expressions](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 4732–4746, Toronto, Canada. Association for Computational Linguistics.

Xiaorong Wang, Clara Na, Emma Strubell, Sorelle Friedler, and Sasha Luccioni. 2023. [Energy and carbon considerations of fine-tuning BERT](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 9058–9069, Singapore. Association for Computational Linguistics.

Zhiwei Yu, Jiwei Tan, and Xiaojun Wan. 2018. [A neural approach to pun generation](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1650–1660, Melbourne, Australia. Association for Computational Linguistics.

Andrea Zaninello and Alexandra Birch. 2020. [Multiword expression aware neural machine translation](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 3816–3825, Marseille, France. European Language Resources Association.

Ziheng Zeng and Suma Bhat. 2022. [Getting bart to ride the idiomatic train: Learning to represent idiomatic expressions](#). *Transactions of the Association for Computational Linguistics*, 10:1120–1137.

## A Snowclone detection task details

In this Section, we report some basic metadata on the CATCHPHRASE dataset as well as the parameters used during the training session and the best sets of parameters we used during the test session. Finally, we show an example of the ranking obtained for the snowclone "may the force be with you".

### A.1 CATCHPHRASE metadata

Table 9 shows some statistics on the CATCHPHRASE dataset. Table 10 shows a sample of the CATCHPHRASE dataset for two snowclones. For each snowclone, we report an identical match, a partial match and a mismatch.

	#Token	#Sentence	#Snowclone
train	50,292	2,974	1,235
dev	11,068	682	60
test	10,389	520	111
total	58,785	3,855	1,406

Table 9: Number of tokens, sentences and sentences containing a snowclone in CATCHPHRASE.

### A.2 Run parameters

The tested parameters include those of the vectorizer and the threshold at which we consider a candidate to correspond to a snowclone. The possible parameters were as follows:

- ngram: 1,2 | 1,3 | 2,3 | 2,4 | 3,4 | 3,5 | 4,5 | 4,6;
- analyzer: word | char | char\_wb;
- threshold: 1, 0.9, 0.8, ... 0.2, 0.1, 0.

The best runs on the test partition of the CATCHPHRASE dataset were the following:

- ASMR<sub>exact</sub>: ngram = 3,4 | analyzer = char | threshold = 0.3;
- ASMR<sub>fuzzy</sub>: ngram = 2,4 | analyzer = word | threshold = 0.3;
- ASMR<sub>combined</sub>: ngram = 1,2 | analyzer = word | threshold = 0.2.

Snowclone	Sentence	Label
may the force be with you	thank you and <b>may the force be with you</b>	1
may the force be with you	<b>may the gods be with you</b>	1
may the force be with you	the ache in my chest from not being able to be with you	0
i love the smell of napalm in the morning	<b>i love the smell of napalm in the morning</b>	1
i love the smell of napalm in the morning	<b>they love the smell of racism in the morning</b>	1
i love the smell of napalm in the morning	i love the smell of christmas	0

Table 10: Some entries of the CATCHPHRASE dataset. We highlight in **bold** the snowclones in each sentence. A label of 1 indicates that the snowclone is seen in the sentence, while a label of 0 indicates that the snowclone is not present in it.

### A.3 Resulting ranking

We report some ranked candidates for the snowclone "may the force be with you" in Figure 3 for the 3 approaches to ASMR. For each approach, we use the best parameters found on the train and dev set for the vectorizer with this approach, which is why some candidate's scores may vary. We also report the impact of threshold on the best run with each approach in Figure 4.

## B MWE identification task details

In this Section, we report basic metadata for the PARSEME 1.3 corpus, the parameters we used during the training session and the best parameters for each language, for each approach. We also show some instances of ranking.

### B.1 PARSEME 1.3 metadata

Figure 5 show the number of sentences and MWE for each language in the PARSEME 1.3 corpus. We notice that some languages are much more represented than others. This is especially the case for Portuguese (PT), Romanian (RO), Chinese (ZH) and Czech (CS), which all contain more than 30,000 sentences.

### B.2 Run parameters

The tested parameters all correspond to a threshold for each linguistic information layer we used during our experiments on the PARSEME 1.3 corpus (token level, morphosyntactic level, lemmas and a semantic similarity score). The possible thresholds were 0.1, 0.4, 0.7 and 1. We limited them in order to avoid overloading our calculation server with longer runs. We report in Table 11 the best parameters for each language and for each approach to ASMR. For the semantic scores, we used the PARAPHRASE-XLM-R-MULTILINGUAL-V1 model from the SENTENCE-TRANSFORMERS python package. This model covers all of the 26 languages of ASMR.

### B.3 Resulting ranking

For 21 language, we show the top 3 candidates of our ranking system obtained with  $ASMR_{combined}$  for a random MWE in Table 12.

## C Error analysis details

We show the the F-scores obtained for each subset of MWE for each language in the PARSEME 1.3 corpus for each approach in Table 13, Table 14 and Table 15.



	ASMR <sub>exact</sub>				ASMR <sub>fuzzy</sub>				ASMR <sub>combined</sub>			
	tok	upos	lem	sem	tok	upos	lem	sem	tok	upos	lem	sem
AR	0.1	0.4	0.7	0.1	0.1	0.4	0.7	0.1	0.1	0.7	0.1	0.1
BG	0.4	0.4	0.7	0.1	0.4	0.4	0.7	0.1	0.4	0.4	0.7	0.1
CS	0.4	0.4	0.7	0.4	0.4	0.4	0.7	0.4	0.4	0.4	0.7	0.1
DE	0.1	0.7	0.7	0.4	0.1	0.4	0.7	0.4	0.1	0.7	0.4	0.4
EL	0.1	0.7	0.7	0.1	0.1	0.7	0.4	0.1	0.1	0.7	0.7	0.1
EN	0.1	0.7	0.7	0.4	0.1	0.7	0.1	0.4	0.1	0.7	0.7	0.4
ES	0.1	0.1	0.7	0.4	0.1	0.1	0.7	0.4	0.1	0.1	0.7	0.4
EU	0.1	0.7	0.7	0.4	0.1	0.7	0.4	0.4	0.1	0.7	0.4	0.4
FA	0.1	0.7	0.4	0.1	0.1	0.7	0.7	0.1	0.1	0.7	0.7	0.4
FR	0.1	0.7	0.7	0.1	0.1	0.7	0.4	0.1	0.1	0.7	0.4	0.1
GA	0.1	0.4	0.7	0.4	0.1	0.7	0.1	0.4	0.1	0.7	0.7	0.4
HE	0.4	0.1	0.7	0.4	0.4	0.1	0.7	0.4	0.4	0.1	0.7	0.4
HI	0.1	0.7	0.7	0.1	0.1	0.7	0.7	0.4	0.1	0.7	0.1	0.1
HR	0.1	0.1	0.7	0.1	0.1	0.1	0.7	0.1	0.1	0.7	0.4	0.1
HU	0.1	0.1	0.4	0.7	0.1	0.1	0.4	0.7	0.1	0.1	0.4	0.7
IT	0.1	0.7	0.7	0.4	0.1	0.4	0.7	0.7	0.1	0.7	0.4	0.4
LT	0.1	0.7	0.7	0.4	0.1	0.7	0.4	0.4	0.1	0.7	0.7	0.4
MT	0.1	0.7	0.7	0.4	0.1	0.7	0.4	0.4	0.1	0.7	0.7	0.4
PL	0.1	0.1	0.7	0.4	0.1	0.1	0.7	0.4	0.1	0.1	0.7	0.1
PT	0.1	0.7	0.7	0.7	0.1	0.7	0.7	0.4	0.1	0.7	0.7	0.7
RO	0.1	0.7	0.7	0.4	0.1	0.7	0.4	0.4	0.1	0.7	0.4	0.4
SL	0.1	0.7	0.7	0.1	0.1	0.7	0.4	0.1	0.1	0.7	0.4	0.1
SR	0.1	0.1	0.7	0.1	0.1	0.1	0.7	0.1	0.1	0.1	0.7	0.1
SV	0.1	0.4	0.7	0.4	0.1	0.4	0.7	0.4	0.1	0.7	0.4	0.4
TR	0.4	0.1	0.7	0.1	0.4	0.1	0.7	0.1	0.4	0.1	0.7	0.1
ZH	0.7	0.7	0.7	0.7	0.7	0.7	0.1	0.1	0.7	0.4	0.7	0.7

Table 11: Best run parameters for each language for each approach for each linguistic information layer: token (tok), morphosyntactic (upos), lemmas (lem) and for the semantic similarity (sem).

Language	Candidate	mean	tok	upos	lem	sem
BG	решаване на проблеми	0.99	0.99	1.0	1.0	0.96
	решаване на проблема	0.98	0.93	1.0	1.0	0.99
	решаване на проблемите	0.97	0.91	1.0	1.0	0.95
CS	mít problém	1.0	1.0	1.0	1.0	1.0
	mít problémů	0.97	0.91	1.0	1.0	0.98
	má problém	0.91	0.69	1.0	1.0	0.95
DE	der entscheiden	1.0	1.0	1.0	1.0	1.0
	Der entscheiden	0.97	1.0	1.0	1.0	0.87
	den entscheiden	0.97	0.89	1.0	1.0	0.98
EL	το παίρνει	0.92	0.83	1.0	1.0	0.85
	Το παίρνει	0.9	0.83	1.0	1.0	0.76
	Ο παίρνει	0.9	0.84	1.0	1.0	0.76
EN	Look forward	1.0	1.0	1.0	1.0	0.99
	look forward	1.0	1.0	1.0	1.0	1.0
	looking forward	0.91	0.7	1.0	1.0	0.94
ES	informar de	1.0	1.0	1.0	1.0	1.0
	informa de	0.93	0.81	1.0	1.0	0.93
	informaron de	0.92	0.81	1.0	1.0	0.88
EU	aintzat hartu	1.0	1.0	1.0	1.0	1.0
	aintzat har	0.99	0.98	1.0	1.0	0.97
	aintzat hartuz	0.95	0.88	1.0	1.0	0.93
FR	se rendre compte	1.0	1.0	1.0	1.0	1.0
	s' rendre compte	0.95	0.8	1.0	1.0	0.99
	se rendant compte	0.88	0.55	1.0	1.0	0.96
GA	baint le	1.0	1.0	1.0	1.0	1.0
	baint leis	0.96	0.94	1.0	1.0	0.92
	bhaint leo	0.87	0.67	1.0	1.0	0.79
HR	nastaviti s	0.99	0.98	1.0	1.0	0.99
	Nastaviti s	0.98	0.98	1.0	1.0	0.95
	nastavi s	0.92	0.71	1.0	1.0	0.98
HU	kötött szerződés	1.0	1.0	1.0	1.0	1.0
	kötött szerződést	0.98	0.94	1.0	1.0	0.99
	kötött szerződésben	0.97	0.91	1.0	1.0	0.98
IT	si prestare	1.0	1.0	1.0	1.0	1.0
	Si prestare	0.98	1.0	1.0	1.0	0.9
	Si prestata	0.92	0.74	1.0	1.0	0.92
LT	sprendimas priimtas	0.92	0.81	1.0	1.0	0.89
	Sprendimas priimtas	0.91	0.81	1.0	1.0	0.84
	sprendimą priimti	0.91	0.65	1.0	1.0	0.98
MT	Il- indusirija	1.0	1.0	1.0	1.0	0.98
	I- indusirija	0.98	0.99	1.0	1.0	0.94
	Iż- indusirija	0.96	0.91	1.0	1.0	0.94
PL	spodziewać się	1.0	1.0	1.0	1.0	1.0
	spodziewają się	0.89	0.69	1.0	1.0	0.88
	spodziewał się	0.89	0.77	1.0	1.0	0.78
PT	ter qualidade	1.0	1.0	1.0	1.0	1.0
	tem qualidade	0.94	0.83	1.0	1.0	0.94
	teve qualidade	0.91	0.75	1.0	1.0	0.9
RO	beneficia de	1.0	1.0	1.0	1.0	1.0
	beneficiat de	0.96	0.87	1.0	1.0	0.97
	beneficiază de	0.94	0.8	1.0	1.0	0.97
SL	se privoščiti	1.0	1.0	1.0	1.0	1.0
	se privoščite	0.97	0.9	1.0	1.0	0.98
	si privoščiti	0.97	0.93	1.0	1.0	0.95
SR	biti u problema	0.98	0.94	1.0	1.0	1.0
	je u problem	0.93	0.79	1.0	1.0	0.93
	sam od problem	0.77	0.5	1.0	0.73	0.83
SV	ta reda på	1.0	1.0	1.0	1.0	1.0
	Ta reda på	0.99	1.0	1.0	1.0	0.97
	får reda på	0.81	0.6	1.0	0.65	0.97
TR	teşekkür etti	0.98	0.97	1.0	1.0	0.97
	teşekkür ederim	0.97	0.89	1.0	1.0	1.0
	teşekkür eden	0.97	0.95	1.0	1.0	0.94

Table 12: top 3 results obtained for a random MWE for 21 languages in PARSEME 1.3 with ASMR<sub>combined</sub>.

ASMR <sub>exact</sub> Candidate	Score	Freq	ASMR <sub>fuzzy</sub> Candidate	Score	Freq
may the force be with you	1.00	51	may the force be with you	1.00	51
may the force be with	0.81	3	may the force be with your	0.69	1
the force be with you	0.74	6	let the force be with you	0.51	6
the force be with	0.58	1	may the force be good to you	0.29	1
may the force be you	0.50	1	may some of the force be with you	0.22	1
force be with you	0.44	3	may the gravity force be with you	0.20	3

(a)

ASMR <sub>combined</sub> Candidate	Score	Freq
may the force be with you	1.00	51
may the force be with your	0.81	1
may some force be with you	0.23	2
may the peace be with you	0.15	1
may the god be with you	0.14	3
may the boop be with you	0.14	1

(c)

Figure 3: Some ranked candidates for the snowclone "may the force be with you" (Star Wars franchise), found in the CATCHPHRASE dataset. We report candidates for each approach.

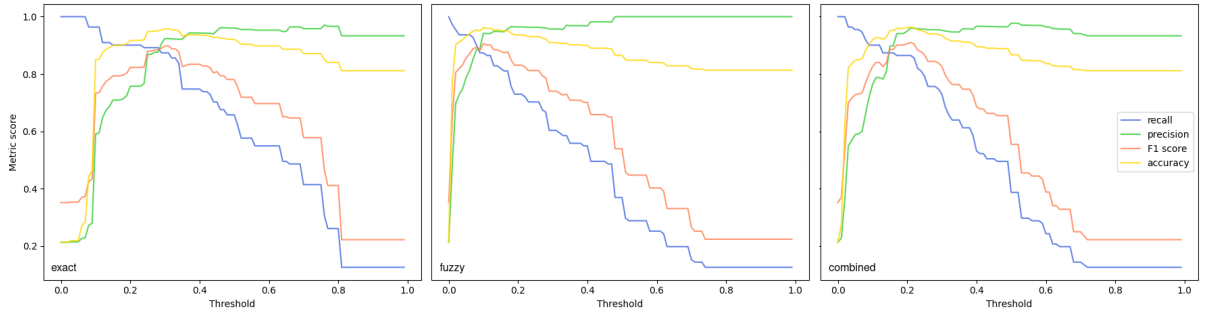


Figure 4: Impact of threshold on recall, precision, F-score and accuracy for the best run on the test partition of CATCHPHRASE with each approach.

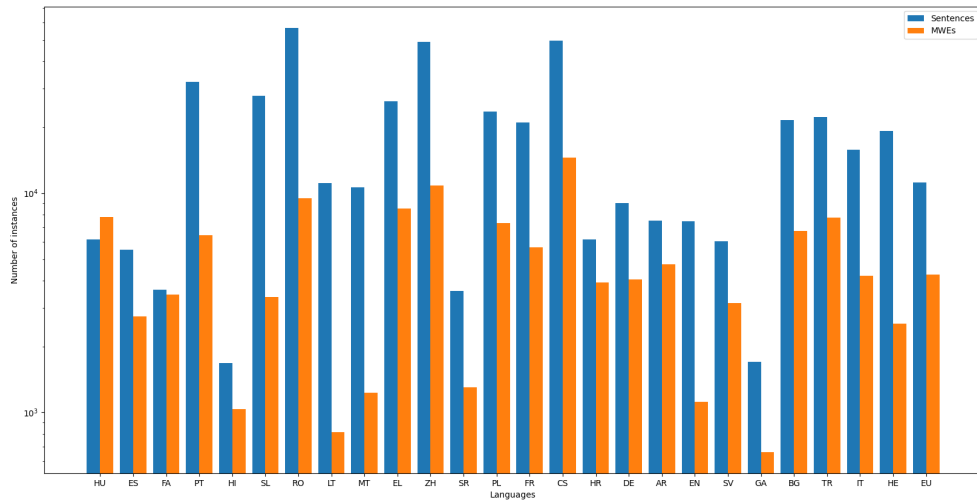


Figure 5: Number of sentences and MWEs for each language in the PARSEME 1.3 corpus.

	Tok-based	Continuous	Discontinuous	Seen	Unseen	Variant	Identical
AR	43.6	48.6	26.5	55.8	00.7	46.7	69.6
BG	62.9	66.7	47.0	70.0	00.0	51.1	80.4
CS	67.6	73.6	49.7	69.1	01.7	61.1	86.6
DE	37.8	38.7	22.0	44.1	00.0	38.3	51.8
EL	59.8	63.6	50.5	72.1	00.0	62.2	85.6
EN	55.5	62.1	47.1	83.2	00.0	73.1	93.3
ES	56.1	59.6	38.2	68.8	00.0	57.8	84.8
EU	75.6	83.7	45.8	81.6	00.0	69.0	95.8
FA	71.3	75.5	39.4	84.5	00.8	75.8	93.2
FR	71.2	75.2	61.0	80.2	00.0	73.1	86.6
GA	30.7	40.4	16.3	62.7	00.0	53.9	92.9
HE	46.4	48.0	38.5	74.2	00.0	54.8	92.1
HI	59.8	62.1	23.9	87.3	00.0	77.4	96.3
HR	75.1	83.3	62.9	87.1	00.0	74.9	93.3
HU	43.0	26.1	62.8	33.1	00.0	48.5	28.9
IT	61.5	67.0	47.8	76.9	00.0	68.1	88.9
LT	38.7	41.1	41.4	76.8	00.0	72.6	98.3
MT	19.5	18.0	11.3	32.1	00.0	31.8	32.5
PL	73.5	80.6	54.4	85.5	00.0	79.5	92.3
PT	59.0	61.3	55.2	82.8	20.3	79.4	90.0
RO	73.8	76.4	63.9	74.9	00.0	46.7	87.9
SL	40.2	43.3	37.4	44.8	00.0	40.4	58.1
SR	53.4	56.4	44.6	80.5	00.0	75.0	94.1
SV	54.1	39.2	56.1	51.9	00.0	58.5	45.9
TR	64.7	69.5	13.5	71.1	00.7	64.5	84.3
ZH	35.5	27.4	32.8	37.6	00.0	31.8	38.7
Mean	55.0	57.1	41.9	68.0	00.8	60.1	78.6

Table 13: F-score obtained for each subset of MWE in each language with the PARSEME 1.3 corpus, using  $ASMR_{exact}$ .



	Tok-based	Continuous	Discontinuous	Seen	Unseen	Variant	Identical
AR	34.6	38.5	11.0	45.3	05.8	33.9	60.6
BG	58.5	63.7	21.0	67.2	06.3	42.0	77.9
CS	56.0	66.4	20.5	60.6	05.5	46.6	84.7
DE	32.6	29.5	06.2	36.5	01.7	28.1	46.3
EL	51.9	60.5	22.4	61.5	08.7	50.2	75.0
EN	41.8	54.6	10.8	69.3	03.8	56.4	80.2
ES	49.8	53.5	17.0	65.2	02.6	51.7	81.5
EU	68.8	76.2	17.1	78.1	03.3	62.8	92.2
FA	72.3	77.1	13.1	84.6	25.8	75.4	92.6
FR	57.8	63.9	22.9	68.4	02.8	55.1	78.1
GA	26.8	36.7	04.8	57.5	06.3	48.1	80.8
HE	43.1	45.7	16.1	69.6	06.9	46.5	88.6
HI	61.1	62.6	05.3	90.8	14.0	84.8	95.8
HR	61.8	73.5	25.9	73.9	03.6	62.1	79.5
HU	39.5	25.5	17.9	28.4	05.0	39.5	25.6
IT	54.6	59.6	17.9	72.1	04.1	61.4	84.8
LT	29.5	39.3	11.7	61.0	02.4	54.3	91.8
MT	18.8	17.5	04.9	30.4	05.7	29.9	31.0
PL	63.1	72.5	25.8	76.5	04.1	66.2	86.8
PT	42.1	52.5	00.0	60.3	20.7	49.2	78.4
RO	64.1	66.9	26.3	68.2	01.9	43.2	77.7
SL	31.0	37.4	15.4	36.4	01.3	29.7	51.9
SR	48.1	54.5	22.3	73.0	06.9	65.7	89.7
SV	47.5	37.8	23.1	44.9	07.1	45.5	43.8
TR	62.2	65.0	05.6	71.6	09.0	66.0	81.8
ZH	35.3	29.9	03.0	34.2	10.5	13.3	36.7
Mean	48.1	52.3	14.8	61.0	06.8	50.3	72.8

Table 14: F-score obtained for each subset of MWE in each language with the PARSEME 1.3 corpus, using  $ASMR_{fuzzy}$ .

	Tok-based	Continuous	Discontinuous	Seen	Unseen	Variant	Identical
AR	40.2	44.4	21.0	54.1	06.2	46.2	66.1
BG	59.5	64.7	38.4	67.6	06.6	45.7	78.1
CS	64.8	70.8	47.3	68.3	04.9	60.0	86.8
DE	35.3	29.4	20.5	39.7	02.8	39.1	40.3
EL	59.1	64.8	48.5	70.0	10.0	62.2	80.7
EN	52.9	58.6	45.3	81.6	00.8	72.2	91.2
ES	52.4	56.0	35.5	67.6	01.3	56.8	82.6
EU	71.5	79.6	41.8	81.0	03.4	68.7	94.8
FA	73.2	78.1	39.4	84.8	23.0	76.4	92.9
FR	67.2	72.0	55.6	80.5	01.0	73.8	86.4
GA	31.7	40.2	16.2	69.0	02.3	61.1	91.2
HE	45.3	47.1	34.9	72.8	05.0	53.1	90.9
HI	58.5	61.7	19.0	86.9	04.1	77.1	95.8
HR	70.0	77.2	57.4	84.8	03.5	73.4	90.7
HU	42.6	25.8	58.5	32.6	03.6	48.1	28.4
IT	60.1	66.5	44.5	77.4	02.1	69.1	88.6
LT	38.5	40.3	41.6	75.6	03.0	71.6	96.2
MT	15.2	13.1	09.4	25.0	01.6	24.7	25.4
PL	67.8	75.0	47.5	83.3	04.2	76.8	90.8
PT	55.9	60.4	48.3	81.3	23.2	77.6	88.9
RO	67.6	71.1	56.4	73.6	02.6	47.5	85.2
SL	37.9	40.2	35.0	44.3	01.7	40.2	56.8
SR	52.9	56.3	42.1	79.8	04.6	74.1	93.8
SV	52.3	37.8	52.6	50.6	02.7	56.9	45.0
TR	59.4	65.0	08.6	68.1	06.2	61.8	80.4
ZH	36.5	28.4	32.9	37.8	05.0	30.5	39.2
Mean	52.6	54.8	38.4	66.9	05.2	59.4	76.4

Table 15: F-score obtained for each subset of MWE in each language with the PARSEME 1.3 corpus, using  $ASMR_{combined}$ .