# Cross-Linguistic Intelligibility of Non-Compositional Expressions in Spoken Context

*Iuliia Zaitova[1], Irina Stenger[1], Wei Xue[1], Tania Avgustinova[1], Bernd Möbius[1], Dietrich Klakow[1]*

[1]Saarland University, Germany

izaitova@lsv.uni-saarland.de, ira.stenger@mx.uni-saarland.de,
weixue@lst.uni-saarland.de, avgustinova@coli.uni-saarland.de,
moebius@lst.uni-saarland.de, Dietrich.Klakow@lsv.uni-saarland.de

## Abstract

This study investigates intelligibility of non-compositional expressions in spoken context for five closely related Slavic languages (Belarusian, Bulgarian, Czech, Polish, and Ukrainian) by native Russian speakers. Our investigation employs a web-based experiment involving free-response and multiple-choice translation tasks. Drawing on prior research, two factors were examined: (1) linguistic similarities (orthographic and phonological distances), and (2) surprisal scores obtained from two multilingual speech representation (SR) models fine-tuned for Russian (Wav2Vec2-Large-Ru-Golos-With-LM and Whisper Medium Russian). According to the results of Pearson correlation and regression analyses, phonological distance appears to be a better predictor of intelligibility scores than SR surprisal.

**Index Terms**: language intelligibility, non-compositional expressions, speech recognition

## 1. Introduction

Non-compositional expressions, which include idioms, metaphors, and fixed phrases, diverge from the linguistic principle of compositionality, which asserts that the meaning of a whole expression comes from the meanings of its parts [1]. Understanding non-compositional expressions extends beyond compositional interpretation [2, 3] and often depends on culture, context, or common understandings. Examples of non-compositional expressions include idioms (e.g., English: "to beat around the bush" meaning: to avoid answering a question; to stall; to waste time), metaphors (Czech: "Život je jako jízda na horské dráze", meaning "Life is like a ride on a roller coaster"), and certain fixed phrases (Bulgarian: "не веднъж" transliterated as "ne vednž"[1], meaning "not once"; Russian: "в конце" transliterated as "v konce", meaning "at the end of").

The intelligibility of non-compositional expressions in an unfamiliar but closely related language is a challenging task, which mainly depends on the mutual intelligibility of languages, i.e., the degree to which a speaker of one language understands the speaker of another language [4]. The degree of success in intelligibility also differs between spoken and written modalities. In the spoken modality, the time available for auditory input processing is limited, whereas in the written modality, one can jump back at will during visual input processing [5]. Moreover, this difference depends on various linguistic and non-linguistic factors varying between languages [6]. For example, a recent study by [7] investigated the auditory intelligibility of idiomatic phrases, which is also a type of non-compositional expression, in two closely related

Slavic languages, i.e., Polish and Russian. The study built on measures of word adaptation surprisal, coupled with syntactic distances (a measure of linguistic similarity) between non-compositional expressions, to predict lay translators' intelligibility scores. However, the study did not include sentential context, which could be an important factor in successful understanding of non-compositional language.

In this paper, we study the intelligibility of non-compositional expressions in their spoken sentential context by native Russian (RU) speakers across five unfamiliar but closely related Slavic languages, i.e., Belarusian (BE), Ukrainian (UK), Bulgarian (BG), Czech (CS), and Polish (PL). The languages RU, BE, and UK belong to the sub-group of East Slavic languages, whereas CS and PL are West Slavic languages and BG is a South Slavic language [8]. Notably, RU, BE, UK, and BG use the Cyrillic script while CS and PL use the Latin script. We conducted a web-based experiment with audio fragments containing target non-compositional expressions in their sentential context along with the written form of the non-compositional expressions. While presenting written stimuli together with their contextual sentences in audio may seem unconventional, this scenario provides a controlled environment to isolate the influence of spoken context. For instance, it enables us to compare intelligibility between written and spoken contexts, offering insights into the dynamics of mutual intelligibility across modalities. Studies with written modality and a combination of the two modalities can be found in [9] and [4], respectively.

Our experiment includes two tasks: a free translation task and a multiple-choice question task (MCQ) task. We analyze the correspondences between the obtained intelligibility scores (percentage of correct responses out of total responses) from the two tasks and two predictive factors: 1) linguistic distances, and 2) surprisal scores from speech representation models. Surprisal scores measure unpredictability by quantifying the negative log-likelihood of encountering a unit given its preceding context [10]. As the probability decreases, the surprisal increases, indicating higher unexpectedness. Through these analyses, we aim to explore which of the two factors affect and better predict intelligibility[2].

## 2. Methodology

### 2.1. Experimental setup

#### 2.1.1. Stimulus preparation

The materials for the audio fragments used in our experiments were based on an existing dataset designed for analyzing non-

---

[1]Here and further, we used ISO 9:1995 transliteration from Cyrillic.

[2]The code and the data for this paper are available at the following link: https://github.com/IuliiaZaitova/spoken-non-compositional-expressions-slavic

compositional expressions, encompassing 227 Russian expressions with their translational correlates and two parallel bilingual context sentences across the five target Slavic languages [11]. For each Slavic language in the dataset, we selected a total of 60 expressions with their context sentences. The mean number of tokens per sentence is as follows: BE: 15.3, BG: 14.9, CS: 11.3, PL: 13.6, and UK: 14.8.

In preparing the assumed incorrect translations, we relied on word-by-word translations from the online bilingual Glosbe Dictionary (https://glosbe.com). Additionally, for identifying cognates, we consulted the etymological online dictionary of the Russian language by Max Vasmer[3]. Including literal translations as incorrect options aims to offer insights into participants' ability to move beyond surface-level comprehension and engage with the deeper (non-compositional) meanings of the expressions.

### 2.1.2. Audio preparation

The audio recordings of our non-compositional expressions with their sentential context were obtained through self-paced reading sessions with five speakers, each of whom was a native speaker of a target language. The three speakers for BG, CS, and UK were female (Bulgarian, Czech, and Ukrainian), and the other two for BE and PL were male. The speakers' age ranges between 21 and 29, with a mean age of 25. All audio recordings were collected in an acoustically controlled environment, at a 44.1 kHz sampling rate in an uncompressed format. The mean duration of each sentence recording across languages is 6.88 sec for BE, 6.82 sec for BG, 4.59 sec for CS, 6.93 sec for UK, and 6.10 sec for PL.

### 2.1.3. Participants

Overall, 118 subjects participated in our experiment, including 92 females, 41 males, and 1 person who identified as another gender. The age of the participants ranged from 18 to 59, with a mean age of 32. The participants were untrained in translation and were recruited for participation in the experiment through Prolific[4], an online platform specializing in participant recruitment for research purposes. We excluded the participants if they had any knowledge of the target foreign language. Since the Prolific platform is in English, we expected the participants to be familiar with the Latin script used by CS and PL languages.

### 2.1.4. Experiment implementation

The experiment was conducted using a custom-built application available online[5]. Participants were presented with instructions in Russian about the tasks and procedures to follow. After familiarizing themselves with the task, participants registered on the website hosting our web application and completed a questionnaire about their background and language skills. Only native Russian speakers without any knowledge of the five foreign languages were included in the analysis.

We used 60 non-compositional expressions with sentential context per language and split them into five subsets, each containing 12 expressions. During the experiment, each participant was exposed to five randomly selected sets, each of which being a subset for one of the five target languages (i.e., BE, UK, BG, CS, and PL). This means each participant received 60 expressions in total. Each subset was presented to each participant only once to avoid repetition effects.

The selected 60 non-compositional expressions were presented in 60 separate trials. In each trial, participants were first asked to listen to the audio containing a non-compositional expression and its sentential context by pressing the play button and then to type a free translation for the target expression highlighted in the audio display bar as shown in Figure 1. The time allocated for translating the highlighted non-compositional expression was based on a formula of 10 seconds per token plus an additional 3 seconds per stimulus. Participants were allowed to play each audio fragment of the sentence containing the target non-compositional expression up to 3 times, which simulates a real-life scenario of a listener asking to repeat what the speaker said. After participants finished the free translation task for an expression presented in its sentential context, they were immediately asked to complete the multiple-choice question (MCQ) task for the same expression, as shown in Figure 2. The task contained two options for participants to choose from: (i) the original non-compositional translation, and (ii) an alternative word-by-word literal translation, which is an inaccurate translation of the expression in terms of semantics. Both options were in Russian, and participants were asked to choose the most suitable one. The goal of the MCQ task was to test participants' preferences for either the non-compositional (correct) option or the literal (incorrect) option.
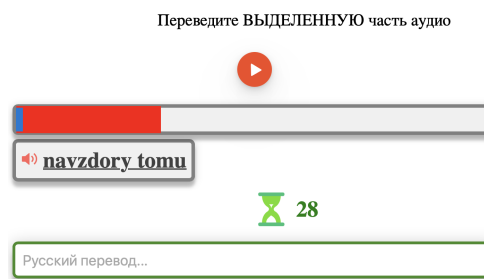


Figure 1: *Experimental screenshot of the free translation task as seen by Russian participants. The instruction on top is: 'Translate the highlighted words without using a dictionary'. Participants' translation is to be typed in the white box, which says 'Russian translation'. The Czech test expression is 'in spite of'.*

## 2.2. Linguistic similarities

Previous research suggests that orthographic and phonological distances are reliable predictors of cross-lingual intelligibility [12, 5, 13]. It is expected that greater linguistic distances correlate negatively with intelligibility scores since greater phonological dissimilarity poses a challenge to mutual intelligibility.

**Orthographic distance.** Measuring the orthographic distance between modern Slavic languages could be challenging due to the use of two writing scripts – Latin and Cyrillic. Normalized Word Adaptation Surprisal (nWAS) quantifies the degree of unexpectedness of a word form given a possibly related word form and a set of transformation probabilities [14]. To utilize the metric, we adapted the code and the orthographic substitution costs computed for Slavic languages used in [15].

**Phonological distance.** Phonologically Weighted Levenshtein Distance (PWLD) quantifies the phonological similarity between different phonemic sequences or word forms [16].
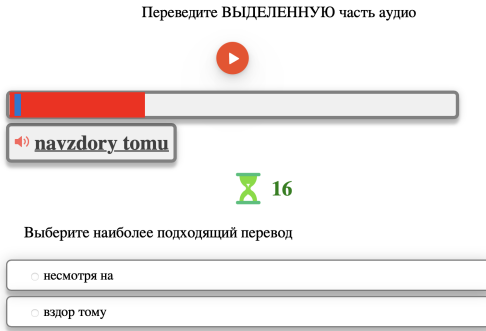
Figure 2: *Experimental screenshot of the MCQ task as seen by Russian participants. The instruction on top says: 'Translate the highlighted words without using a dictionary'. Here, the line below the sandclock says 'Choose the most suitable translation'. The Czech expressions translated to 'in spite of'. The two Russian options are 'in spite of' (non-compositional) and 'nonsense to that' (literal).*

This metric extends the string-based Levenshtein distance by considering the cost of each phoneme substitution based on their phonological features. We employ the same adaptation of the original PWLD as the one proposed in [17]. The phonemic transcriptions for all non-compositional expressions in the target languages and RU were obtained using CharsiuG2P, a transformer-based tool for grapheme-to-phoneme conversion [18].

### 2.3. Surprisal scores from speech representation models

We calculated sentence-level surprisal for all the audio fragments (entire sentences containing non-compositional expressions) presented in the experiment using two widely used speech representation (SR) models, both accessible through the HuggingFace Model Hub. The first one is Wav2Vec2-Large-Ru-Golos-With-LM (Wav2vec) [19], a large-sized model fine-tuned in Russian using Sberdevices Golos [20] with audio augmentations. The second one is Whisper Medium Russian (Whisper), a medium-sized model fine-tuned using audio data from the Open STT Russian Dataset [21]. Surprisal from SR models serves as a proxy for the difficulty of processing the audio fragment. The process of extracting surprisal using Python involved 1) preprocessing the audio, 2) generating predictions, and 3) calculating surprisal by obtaining the negative log-likelihood of predicted probabilities for each unit in the fragment.

## 3. Results and Discussion

### 3.1. Intelligibility scores

Figure 3 illustrates the intelligibility scores (percentage of correct responses out of total responses) for the free translation and MCQ tasks, represented as the percentage of correct responses out of total responses. Overall, participants performed worse in the free translation task compared to the MCQ task. In both tasks, the highest scores can be observed for BE and UK, which aligns with previous studies on the intelligibility between Slavic languages [22] and the fact that BE, UK, and RU belong to the same sub-group of Slavic languages (see Section 1). When it comes to the MCQ task, aside from BE and UK, BG also shows a relatively high intelligibility score and outperforms CS and

Table 1: *Pearson correlation coefficients between experiment scores and four predictive variables*

| | nWAS | PWLD | Wav2vec | Whisper |
|---|---|---|---|---|
| **Intelligibility scores in free translation task** | | | | |
| All | -0.178*** | -0.403*** | **0.404**\*** | 0.223*** |
| BE | -0.210 (NS) | -0.300* | 0.054 (NS) | **0.427**\*** |
| UK | -0.054 (NS) | **-0.558**\*** | 0.355** | 0.120 (NS) |
| BG | -0.339** | **-0.445**\*** | 0.243 (NS) | 0.099 (NS) |
| PL | -0.012 (NS) | **-0.284*** | 0.030 (NS) | 0.039 (NS) |
| CS | -0.006 (NS) | -0.093 (NS) | 0.009 (NS) | 0.232 (NS) |
| **Intelligibility scores in MCQ task** | | | | |
| All | -0.259*** | **-0.415**\*** | 0.399*** | 0.193*** |
| BE | -0.200 (NS) | -0.225 (NS) | 0.167 (NS) | **0.393*** |
| UK | -0.395** | **-0.623**\*** | 0.465*** | 0.142 (NS) |
| BG | -0.323* | **-0.330*** | 0.318* | 0.056 (NS) |
| PL | 0.183 (NS) | **-0.354*** | -0.058 (NS) | 0.194 (NS) |
| CS | -0.219 (NS) | **-0.423**\*** | -0.136 (NS) | 0.064 (NS) |

\*=$p < .05$, \*\*=$p < .01$, \*\*\*=$p < .001$, NS=Not Significant

PL. This could be attributed to the use of the Cyrillic script in both BG and RU, which likely improves the scores in the MCQ task, where participants have to map the target expression to two written options. However, as BG and RU are not from the same sub-group of Slavic languages, participants may not be able to use the auditory information to achieve better performance in the free translation task. Interestingly, the intelligibility scores of free translation task in this study (in the spoken context) are higher in comparison to the written experiment described in [9]. We hypothesize that the presence of both spoken and written forms of the expression may aid in better comprehension.
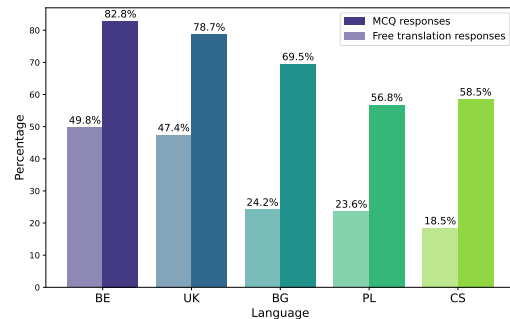


Figure 3: *Intelligibility scores of Free Translation and MCQ Responses.*

### 3.2. Correlation of intelligibility scores with linguistic factors and surprisal scores

For each experimental task, we correlated the intelligibility scores with the two linguistic distance variables 1) Phonologically Weighted Levenshtein Distance (PWLD) and 2) Normalized Word Adaptation Surprisal (nWAS). We also correlated the scores with the audio surprisal from Wav2vec and Whisper, as described in Section 2.3. The results are presented in Table 1. The results were calculated on the basis of all target languages jointly as well as for each language separately.

Regarding the results of linguistic similarities, the analysis of all languages jointly shows that both PWLD and nWAS have a significant Pearson correlation with intelligibility scores

in both tasks. The highest correlation for intelligibility score in the MCQ task is the correlation with PWLD. When taking languages individually, we can again observe a significant negative correlation with PWLD in both tasks for most of the languages. However, for BE, there is only a non-significant correlation with PWLD, which potentially indicates that phonological distance may not be a dominant factor influencing intelligibility in this specific language. As for nWAS, a significant correlation is only observed for BG in both free translation and MCQ tasks, as well as for UK but only in the MCQ task. Notably, these two languages and RU use the Cyrillic script, which makes orthographic mapping more straightforward. The negative correlations of intelligibility scores with linguistic similarities align with the expectation that greater phonological and orthographic dissimilarity poses challenges for intelligibility, as mentioned in Section 2.2.

Regarding the results of SR surprisal, when all languages are combined, the highest absolute correlation of intelligibility scores can be observed with Wav2vec surprisal in the free translation task. For BE and UK, languages that belong to the same sub-group of Slavic languages as RU, we can observe a positive correlation with SR model surprisal in both tasks. For BE, the correlation with the Whisper model surprisal outperforms all the other variables. The positive association between surprisal and improved intelligibility seems counterintuitive at first. A possible explanation is the transferability of linguistic features in closely related languages. Languages like BE and UK share a considerable amount of linguistic similarity with RU. The positive correlation with surprisal scores suggests that the models are capturing linguistic structures common to these languages, and surprisal scores only function when these linguistic structures are similar enough. When a language is rather different from RU, the model might perceive the audio input of that language as noise continuously and thus produce lower surprisal.

### 3.3. Stepwise regression for predicting intelligibility scores using linguistic factors and surprisal scores

We further analyzed the intelligibility score by conducting stepwise regression analyses with intelligibility scores as the dependent variable for the two tasks separately. The analyses were performed to identify the effect of joint predictors. The results for the best performing models of the stepwise analyses are presented collectively for all languages and for each language individually for each task. We present only significant results of these best-performing models in Table 2. As shown in Table 2, the amount of explained variance is generally low. When it comes to the model for all languages, all predictors play a significant role. However, the languages with the lowest intelligibility scores (CS and PL) do not incorporate surprisal from the models as predictors. Also, for CS, we did not observe any significant predictors in the best-performing model in the free translation task. These results potentially indicate that spoken context is less helpful for more distant languages.

### 3.4. Limitations

While this study contributes valuable insights into cross-lingual intelligibility of non-compositional expressions in spoken context, certain limitations should be acknowledged. First of all, the study relies on a specific group of participants, namely native Russian speakers, which might limit the generalizability of the findings. Secondly, for our analyses, we use SR models fine-tuned specifically for Russian. Generalizing the findings to other languages should be done cautiously. Moreover, the stim-

Table 2: *Stepwise regression analysis*

**Intelligibility scores in free translation task**

|     | Predictor | $R^2$ coeff | t-value | p-value |
|-----|-----------|-------------|---------|---------|
| All | nWAS      |             | -5.157  | 0.000   |
|     | +PWLD     |             | -9.583  | 0.000   |
|     | +Wav2vec  |             | 11.828  | 0.000   |
|     | +Whisper  | 0.3         | 5.606   | 0.000   |
| BE  | Whisper   | 0.24        | 3.165   | 0.003   |
| UK  | nWAS      |             | 2.083   | 0.042   |
| BG  | nWAS      |             | -2.015  | 0.049   |
|     | +PWLD     | 0.29        | -2.888  | 0.006   |
| PL  | PWLD      | 0.09        | -2.227  | 0.030   |

**Intelligibility scores in MCQ**

|     | Predictor | $R^2$ coeff | t-value | p-value |
|-----|-----------|-------------|---------|---------|
| All | nWAS      |             | -3.991  | 0.000   |
|     | +PWLD     |             | -9.126  | 0.000   |
|     | +Wav2vec  |             | 11.68   | 0.000   |
|     | +Whisper  | 0.31        | 4.376   | 0.000   |
| BE  | Whisper   | 0.2         | 2.800   | 0.007   |
| UK  | PWLD      |             | -4.367  | 0.000   |
| BG  | nWAS      |             | -2.128  | 0.038   |
|     | +Wav2vec  | 0.25        | 2.472   | 0.017   |
| PL  | PWLD      | 0.21        | -2.912  | 0.005   |
| CS  | PWLD      | 0.19        | -2.939  | 0.005   |

uli used in the experiment may introduce a gender bias. The three speakers for Bulgarian, Czech, and Ukrainian were female, while the two for Belarusian and Polish were male. This mixture of genders over speakers may potentially influence participants' perceptions. Acknowledging these limitations is crucial for interpreting the study's results.

## 4. Conclusion

In this study, we presented the results of a web-based experiment on the intelligibility of non-compositional expressions from Belarusian, Bulgarian, Czech, Polish, and Ukrainian in spoken context by native Russian speakers. The experiment consisted of two tasks, i.e., a free translation task and a multiple-choice task.

We observed that the intelligibility scores in both tasks are highest for Belarusian and Ukrainian, languages within the same (East Slavic) sub-group as Russian, followed by Bulgarian (South Slavic), and lowest for Czech and Polish (West Slavic). The intelligibility scores can be better explained by the phonological distance between the target non-compositional expressions and the correct Russian non-compositional expressions. Additionally, for Belarusian and Ukrainian, languages that are closer to Russian compared to the other three languages, surprisal extracted from speech representational models fine-tuned on Russian was found to be a significant predicting factor. However, surprisal is not a significant predictor of intelligibility for languages that are more distant from Russian (i.e., Czech and Polish). Future work will include further exploration of intelligibility for these non-compositional expressions in different language groups and across various modalities.

## 5. Acknowledgements

# 6. References

[1] B. Partee, *Compositionality in Formal Semantics: Selected Papers*, ser. Explorations in Semantics. Wiley, 2008.

[2] T. Baldwin and S. N. Kim, "Multiword expressions," in *Handbook of Natural Language Processing, Second Edition*, N. Indurkhya and F. J. Damerau, Eds. Chapman and Hall/CRC, 2010, pp. 267–292.

[3] R. Jackendoff, *Foundations of Language: Brain, Meaning, Grammar, Evolution*. Oxford University Press UK, 2002.

[4] C. Gooskens and F. Swarte, "Linguistic and extra-linguistic predictors of mutual intelligibility between germanic languages," *Nordic Journal of Linguistics*, vol. 40, pp. 123–147, 10 2017.

[5] R. Möller and L. Zeevaert, "Investigating word recognition in intercomprehension: Methods and findings," *Linguistics*, vol. 53, 03 2015.

[6] C. Gooskens, *8. Receptive Multilingualism*. Berlin, Boston: De Gruyter Mouton, 2019, pp. 149–174.

[7] J. Kudera, I. Stenger, P. Georgis, B. Möbius, T. Avgustinova, and D. Klakow, "Cross-linguistic intelligibility of idiomatic phrases in polish-russian translation tasks," in *Phraseology, Constructions and Translation: Corpus-based, Computational and Cultural Aspects*, J.-P. Colson, Ed. Presses Universitaires de Louvain, 2023, pp. 237–249.

[8] R. Sussex and P. Cubberley, *The Slavic Languages*. Cambridge: Cambridge University Press, 2006.

[9] I. Zaitova, I. Stenger, M. U. Butt, and T. Avgustinova, "Cross-linguistic processing of non-compositional expressions in Slavic languages," in *Proceedings of the Workshop on Cognitive Aspects of the Lexicon @ LREC-COLING 2024*, M. Zock, E. Chersoni, Y.-Y. Hsu, and S. de Deyne, Eds. Torino, Italia: ELRA and ICCL, May 2024, pp. 86–97. [Online]. Available: https://aclanthology.org/2024.cogalex-1.10

[10] M. Crocker, V. Demberg, and E. Teich, "Information density and linguistic encoding (ideal)," *Künstliche Intelligenz*, vol. 30, pp. 77–81, 2016.

[11] I. Zaitova, I. Stenger, and T. Avgustinova, "Microsyntactic unit detection using word embedding models: Experiments on slavic languages," in *Proceedings of the International Conference on Recent Advances in Natural Language Processing*. INCOMA Ltd., 2023, pp. 1265–1273.

[12] J. Vanhove and R. Berthele, "Item-related determinants of cognate guessing in multilinguals," *Crosslinguistic Influence and Crosslinguistic Interaction in Multilingual Language Learning*, vol. 95, p. 118, 2015.

[13] C. Gooskens and F. Swarte, "Linguistic and extra-linguistic predictors of mutual intelligibility between germanic languages," *Nordic Journal of Linguistics*, vol. 40, pp. 123–147, 10 2017.

[14] I. Stenger, K. Jágrová, A. Fischer, T. Avgustinova, D. Klakow, and R. Marti, "Modeling the impact of orthographic coding on czech–polish and bulgarian–russian reading intercomprehension," *Nordic Journal of Linguistics*, vol. 40, no. 2, p. 175–199, 2017.

[15] I. Stenger, P. Georgis, T. Avgustinova, B. Möbius, and D. Klakow, "Modeling the impact of syntactic distance and surprisal on cross-Slavic text comprehension," in *Proceedings of the Thirteenth Language Resources and Evaluation Conference*. Marseille, France: European Language Resources Association, Jun. 2022, pp. 7368–7376.

[16] L. Fontan, I. Ferrané, J. Farinas, J. Pinquier, and X. Aumont, "Using phonologically weighted levenshtein distances for the prediction of microscopic intelligibility," in *Annual conference Interspeech (INTERSPEECH 2016)*, 2016, p. 650.

[17] B. M. Abdullah, M. Mosbach, I. Zaitova, B. Möbius, and D. Klakow, "Do Acoustic Word Embeddings Capture Phonological Similarity? An Empirical Study," in *Proceedings of Interspeech 2021*, 2021, pp. 4194–4198.

[18] J. Zhu, C. Zhang, and D. Jurgens, "Byt5 model for massively multilingual grapheme-to-phoneme conversion," 2022.

[19] I. Bondarenko, "Xlsr wav2vec2 russian with 2-gram language model by ivan bondarenko," https://huggingface.co/bond005/wav2vec2-large-ru-golos-with-lm, 2022.

[20] N. Karpov, A. Denisenko, and F. Minkin, "Golos: Russian dataset for speech research," 2021. [Online]. Available: https://arxiv.org/abs/2106.10161

[21] A. Slizhikova, A. Veysov, D. Nurtdinova, D. Voronin, and Y. Baburov, "Russian open speech to text (stt/asr) dataset v1.0," https://github.com/snakers4/open_stt/, 2019.

[22] I. Stenger and T. Avgustinova, "On slavic cognate recognition in context," in *Computational Linguistics and Intellectual Technologies: Papers from the Annual International Conference 'Dialogue'*, vol. 20, Moscow, Russia, June 2021, pp. 660–668.