

Павел Малецкий 

АГГ Науково-Технічний Університет/AGH University of Krakow

Магдалена Піотровска 

АГГ Науково-Технічний Університет/AGH University of Krakow

# Аналіза і класифікація русиньської бесіди язовим модельом штучной інтелігенції OpenAI Whisper

## Abstrakt

**Analiza i klasyfikacja języka rusińskiego przy użyciu modelu sztucznej sieci neuronowej ASR OpenAI Whisper**

Artykuł przedstawia analizę lingwistyczną języka rusińskiego, koncentrując się na jego złożonych i zmieniających się aspektach, takich jak wymowa oraz różnice indywidualne, regionalne i historyczne. Do przeprowadzenia badania wykorzystano sztuczną sieć neuronową opartą na modelu OpenAI Whisper. Model ten, choć szkolony na danych z większości państwowych języków urzędowych, nie był bezpośrednio trenowany na bazach próbek języka rusińskiego ze względu na jego lokalny i mniejszościowy/etniczny charakter. Stąd próbki mowy tego języka klasyfikowane były przy użyciu najbardziej zbliżonych dostępnych etykiet, co pozwoliło na wyznaczenie podobieństwa języka rusińskiego do innych słowiańskich języków. Badanie objęło użytkowników zróżnicowanych pod względem płci, wieku i lokalizacji (Polska, Ukraina, Słowacja, Serbia), wykazując znaczące podobieństwa do języków dominujących w tych krajach oraz zależności między wyznaczonym podobieństwem językowym a wiekiem mówców.

**Słowa kluczowe:** język rusiński, fonetyka, klasyfikacja, asymilacja, AI, sztuczne sieci neuronowe, ASR

**Abstract****Analysis and Classification of the Rusyn Language Using the OpenAI Whisper ASR Model**

The paper presents a linguistic analysis of the Rusyn language, focusing on its complex and dynamic aspects, such as pronunciation and individual, regional, and historical variations. The study employed a neural network based on the OpenAI Whisper automatic speech recognition (ASR) model. While trained on data from a majority of official state languages, the model lacked direct training on Rusyn language samples due to its localized and minority/ethnic nature. Consequently, speech samples were classified using the closest available language labels, enabling the identification of similarities between Rusyn and other Slavic languages. The study encompassed a diverse range of speakers across gender, age, and location (Poland, Ukraine, Slovakia, Serbia), revealing significant similarities to the dominant languages in these respective countries. Furthermore, the research highlights correlations between the identified linguistic similarities and the age of the speakers.

**Keywords:** Rusyn language, Phonetics, Classification, Assimilation, AI, ANN, ASR

**Ключовы слова:** русиньскій язык, фонетыка, класифікація, асиміляція, ШІ, штучны нейронны сіті, ASR

## 1. Вступ

Русиньскій язык то східньославянскій язык хоснуваний через русиньску етнічну групу, головні в карпатскым регіоні Серединовой Європы, котрый обнимат част України, Словації, Польщы, Мадяр і Румунії. Має місцевы одміны, што оддзеркаляють зложены історичны і культуровы вплиня в поєдных регіонах. Класифікація русиньского языка в шыршым контексті східньославянських і інчых славянських языків є зложеном і дискусийном темом з огляду на його унікальне положыня і історичны вплиня. Русиньскій язык має спільны приметы зо східньославянськыма языками (як російскій, українскій і білорускій), а тіж выказує вплиня західньославянських языків (польскій, словацькый) і полудньовославянських (сербскій, хорватскій). Тота вельоаспектова мішанина оддзеркалят його положыня на скривуваню тых языковых груп (Kushko 2007).

Історичні русинський язык был залежный од языковой політыкы паноучых держав ци цисарств, медже інчыма Австро-Мадярской Імперіи ци Совітского Союзу. Языковой проблем є кісно повязаны з достоменністю, а дискусії на тему того ци русинський то окремый язык, ци диалект українського, тырвають іщы аж і днес. Прихыльниками той другой теорії сут мало не выключні дослідники і політыкы, які ідентифікують Руснів як Українців. В Україні Русины не сут офіційно узнаны як окрема етнічна група, што іщы веце комплікує охорону языка (Nikitin і ін. 2009). Русинський язык был скодифікованы в ХХ ст., што оддзеркаляло одроджыня народовой і языковой достоменности. Кодифікаційны діяння довели до выокремліня регіональных літературных стандартів опертых на місцевых приметах – в Словації, Польщы, Україні і в Сербії (Plišková 2008).

Русинський єст переважні класифікованы як східньославянський язык з огляду на його історичне і языкове коріння. Його класифікацію комплікує єднак значуче вплиня західньо- і полудньовославянських языків (Moser 2016).

Змінінст языка, што обнимат диасхронны, диастратычны, диалектычны і діятотічны аспекты, то вызваня для обробкы натурального языка (NLP, en. Natural Language Processing). Диалектычны ріжніці медже одмінами того самого языка вплиют на ефективність аплікацій, таких як машынове тлумачыня ци розпознаваня бесіды (Zampieri і ін. 2020). Прото тіж росне заінтерсуваня досліджынями над обробком споріднених языків, языковых одмін і диалектів, на што доказом сут чысленны публікації і науковы події, такы як варштаты VarDial (Zampieri і ін. 2020).

В примірі языків з низкыма ресурсами, як нп. русинський, придбати відповідні языковы корпусы єст особливі тяжко (Scherrer, Rabus 2019). Часто мож того розвязати транскрибуючы бесіду, як в примірі корпусу ArchiMob для німецкых диалектів або хоснуючы тлумачыня, як в примірі корпусу MADAR для арабскых диалектів (Bouamor і ін. 2019). Ахім Rabus і Ів Scherrer описують вызваня звязаны з творіньом ресурсів NLP для русинського языка і пропонують індукцію морфосинтактычного лексикону зо схоснуваньом ресурсів славянських языків (Rabus, Scherrer 2017). Досліджають они морфосинтактычне етыкетуваня (анг. tagging) для русинського языка, хоснуючы трансфер вчыня зо споріднених языків. Вказаны выбраны однесія сут доказом того, што тема є інтенсивні досліджана і експлоатувана. Істніють спеціалістычны програмы і аплікації до

дослідження языковых примет, але схоснування моделю так барз зложеного і чутливого, як увзгляднений в тым дописі, ест інновацийным підходом.

Цілю працы є дослідити ступін розпознавальности і клясифікації русиньского языка через модель автоматычного розпознавання бесіды OpenAI Whisper і оцінити вплиня домінуючых языків держав, де жытют Русины на результат клясифікації. Обсяг працы обнимат аналізу звуковых даных зобраных з русиньскоязычных радиювых авдиций зо штырьох держав: Польщы, Украіны, Словаціі і Сербії, а тіж другы (до-датковы) аналізы з огляду на групы бесідників подля іх років. В статі поміщено основны інформациі на тему моделю OpenAI Whisper, в тым його архітектуру і приметы, якы чынят можливом вельоязычну транскрипцію і дрібницьовый опис дослідничой методології, што обнимат зміст бази звуковых даных, альгоритмів сегментации і аналізы звуку та параметрів схоснуваных в модели.

Як результаты представлено клясифікацію русиньского языка при увзгляднію держав, одкале походят бесідуючы особы і ріжниц медже групамы подля років. Результаты омовлено в контексті вплиня домінуючых языків на язык Русинів і проходячых асиміляцийных процесів. Представлено тіж квестії звязаны з ограничынями моделю і специфікы языковой клясифікації. Працу завершат підсумуваня, што вказуе на основны высновкы і можливы напрямы дальшых досліджынь в обшыри автоматычного розпознавання бесіды для меншыньовых языків.

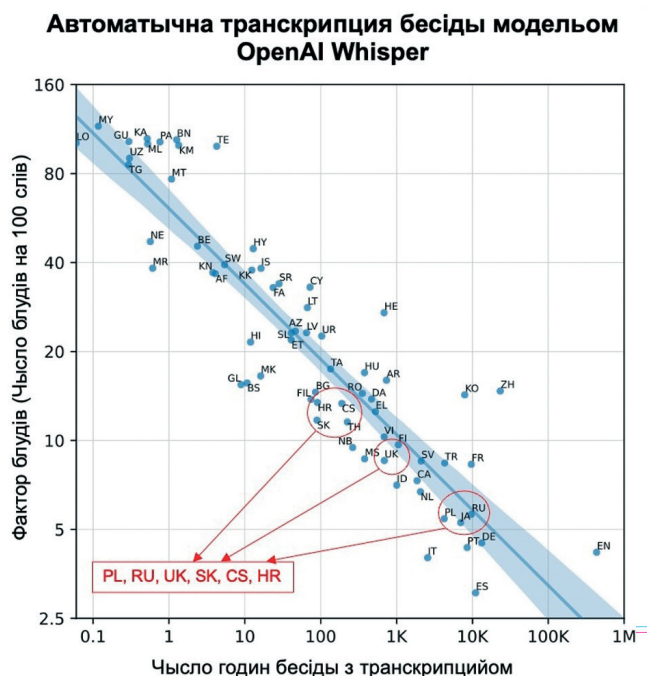
## 2. Языковой модель OpenAI Whisper

OpenAI Whisper то заавансуваный модель автоматычного розпознавання бесіды ASR (анг. Automatic Speech Recognition) (Radford і ін. 2022). Модель і програмовый код сут одкрыто доступны на ліценції тыпу *open source*, што уможливят динамічне розвита аплікації і повязаны досліджыня в обшыри заавансуваной переробкы бесіды. Сес модель вытренувано на основі веце як 680 тысячи годин вельоязычных награнь з транскрипциейом. Обшырна і ріжнорідна база даных дозволила осягнути велику одпорніст системы на зріжницюваны акценты, галас на фоні ци спеціалістычну термінологію, спомогла при тым транскрипцію в вельох языках.

Архітектура OpenAI Whisper ест приміром реализации системы ASR зо схоснуваньом невруновой сіти тыпу *Transformer* в конфігурації

енкодер-декодер (анг. encoder-decoder). Звукові дані на вході моделю ділені сут на 30-секундові частки, пак переформлювані на *log-Mel spectrogram* (є то дигітальна, графічна репрезентація звуку, яка бере до уваги приметы перцепції чловека), а дале перерабляны через екодер. Декодер передвидує зміст, што одповідат звукам і вводит додаткові інформаційны токены, такы як ідентифікація языка бесідуючої особы і часовы значныкы з докладністю до рівня фраз бесіды.

Во вчыню моделю OpenAI Whisper схоснувано барз велику і зріжницювану базу даных, што ілюструє Рыв. 1, на котрым вказано залежніст процента невластиві розпознаных слів од кількості даных до тренуваня. В базі была менше-вече третя част награнь в інчым як англійській языку. І хоц в базі єт о вельо менше тестовых награнь в інчых языках, то ефектывніст розпознаваня языка, при базі велькості більше як сто годин, єт на рівні 80 проц. і высшым. До того зараховуют ся вшыткы значучы з перспектывы дослїджынь славянськы языкы, што тіж зазначено на Рыв. 1.



Рыв. 1. Залежніст ефектывності моделю ASR, в залежності од велькості даных до тренуваня (Radford і ін. 2022). На выкресі зазначено выбраны славянськы языкы (польській, українській, словацкій і російській)

Модель OpenAI Whisper розпознає язык wypowiedi через вельоетаповый процес, што ест зінтегрований з його архітектуром енкодер-декодер. Логарытмічны Мелы спектрограмы сут передаваны до енкодера, што выокремлят зложены звуковы приметы. Енкодер выокремлят взірці в звуковых даных на барз высокым рівни зложености, якы характеризуют ріжны языки. Прогноза основана ест на высокоплощыновых приметах, што обнимают фонетычны і фонологічны взірці особливы для каждого языка. В час вчыня модель был експонуваний на звуковы даны в барз вельох языках, ведно з відповідніма для них текстовима етикетами, в котрых сут токены до ідентифікації языка. В процесі декодування звуку декодер хоснує сесы вивчены токены, жебы окрислити язык на основі введеных даных/звуків. Завдякы так обшырному вчыню модель аналізує сигналы в комплексовый спосіб в ріжных языках і акцентах, што поправлят його спосібніст до прецизийного розпознавання языка. OpenAI Whisper хоснує тіж здібности *zero-shot learning*, што означат, же на основі даных з вчыня може робити узагальніня, жебы розпознати і транскрибувати языки, на яких не был конкретні вчений. Осігат того хоснуючы вельоязычний характер даных до тренування і заавансуваны можливости выокремляня примет через енкодер.

Схоснування моделю OpenAI Whisper до клясифікації русиньского языка, хоц модель тот не был вчений на русиньській базі, мож пояснити парома чынниками. Модель OpenAI Whisper был вытренуваний на великій, вельоязычній базі даных, в якій были ріжны славянськы языки (такы як польській, українській, словацкій, сербській, хорватській ци російській), што сут зближены до русиньского языка під оглядом ґраматычной структуры, фонетыкы ци лексики. Завдякы тому модель годен розпознати деякы взірці і приметы характерны для группы славянських языків, што уможливят му зближену клясифікацію русиньского языка, хоц не было непосреднього вчыня на його базі. В примірі русиньского языка, котрый не ест релятивні масовый, схоснування моделю вытренованого на шыроком спектрі языків дозвоят на його приближену клясифікацію на основі примет спільных з інчыма языками.

### 3. Методологія

#### а. База даних

До мовної аналізи схоснувано архівальні награвня лемківської радіоїмової стації Лем.фм. В базі звукових даних є барз вельо зріжницюваних примірів русиньського бесідуваного мови, што уможливило комплексове досліджыня з шырокой перспектывы. З каждой авдиції усунено фрагменты музыкы, джінглі і реклямовы споты, а з інтервю лишено голос лем едной особы.

Так приготовлена база складала ся з 470 ріжных авдицій. Час тырваня цілой базы то даде 121,8 годин. Поділ походжыня бесідующих осіб з кількістю авдицій і часом тырваня досліджаних пробок представлено в Табелі 1.

Скорочыня	Держава	Кількіст авдицій	Час тырваня [годин]	Чысло лекторів
PL	Польша	209	45,4	веще як 100
UK	Україна	99	37,5	веще як 40
SK	Словачия	131	23,2	веще як 50
CS	Сербія	31	11,5	10

Таб. 1.

Додатковы метадамы приписаны до каждой авдиції:

- Група подля років: дві категорії – бесідующы особы што мают вещь або менше як 70 років (лем для осіб з Польщы).
- Ідентифікація бесідующей особы: мено або псевдонім.
- Назва файлів выгенерованих в час аналізы.
- Полный час награвня (в секундах) каждой звуковой пробкы.

#### б. Розпознаваня мови бесідующего зо схоснуваньом моделю OpenAI Whisper

До аналізы і категоризації даних опрацувано скрипт в мову *Python*, котрый уможливят автоматычний выбір і переробку звукових файлів, мовну аналізу і реєстрацію результатів.

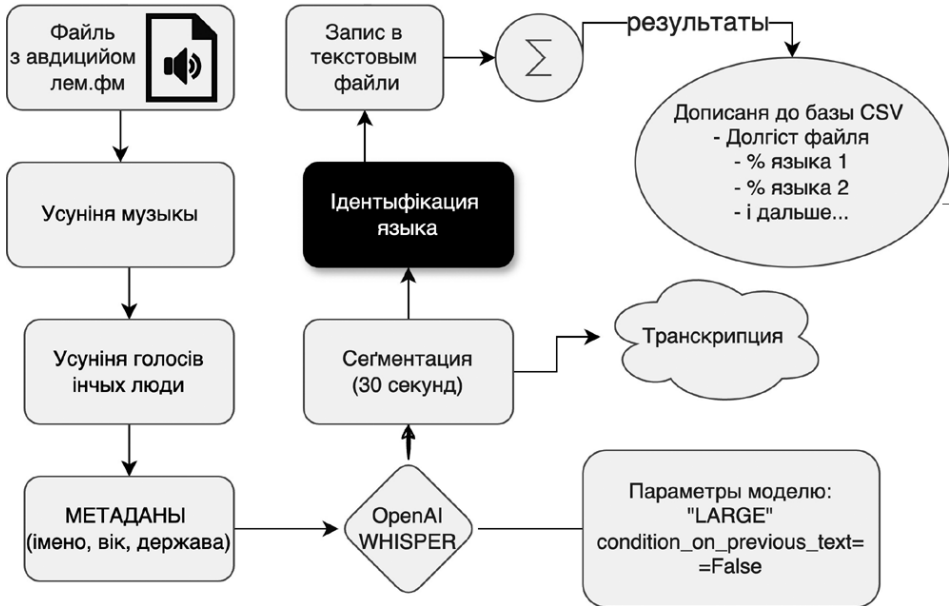


Рис. 2. Схема процесу ідентифікації языка радийовой авдиції

Діяння скрипту схемово представлено на Рис. 2, а поедны етапы його діяння сут слідуочы:

- выбір каталогу: выбір каталогу з авдіо файлами при схоснувано графічного інтерфейсу (GUI),
- сегментация звуку: каждый авдіо файл ест ділений на часткы долгости 30 секунд,
- діяння моделю: до транскрипції бесіды і вызначыня языка схоснувано модель Whisper в варіянті «large». Архітектура моделю запроектувана ест з 32 верств сіти шырокости 1280 неврону і 20 голу (анг. heads), што разом дае 1,55 мільярда параметрів. Ест то найвекшы модель в родні Whisper, што перекладат ся на высокую прецизию транскрипції і розпознаваня языка завдякы його барз зложеній структурі. Модель робит транскрипцію бесіды для каждого сегменту авдіо, а на выході моделю сут додатковы інформації, в тым вызначеный язык. Параметры моделю сконфігуруваны сут на вартости: *no\_speech\_threshold=0.8* і *condition\_on\_previous\_text=False*, завдякы чому каждая наступна 30-секундова пробка аналізувана ест незалежні од попередньої. Шырокоств сіти дефініюе чысло неврону в каждой верстві моделю, значыт розмір вектора репрезентации, што є перетворювана через модель на даній верстві,



а число голів односит ся до числа «голів увагы» (анг. attention heads) в кожій верстві механізму увагы моделю. Тот механізм дозволять зосередити ся модельови на ріжних частях секвенції на вході в час переробкы. Кожда головка аналізує даны під інчим оглядом, завдякы чому модель може ліпше выіматн зложены залежности в языковых секвенциях,

- реєструваня результатів: результати транскрипції і вызначений язык для каждого сегменту сут записуваны запоряд до текстовых файлів,
- экспорт до CSV: результати сут записуваны в файли CSV, в котрым сут первістны метадамы і інформації о розпознаных языках і їх процентова участ в полном часі награня.

Скрипт додатково хоснує бібліотеку *librosa* до ладуваня звуку і *pydub* до сегментації даных.

#### 4. Результаты

Даны переаналізувано в тот спосіб, же в кожій авдиції окрислено шор розпознаных языків, зачынаючы од найчастійше розпознаваного (#1), ведно з ідентифікатором державы бесідуючого. Пак порахувано процентову участ каждого з языків в роли языка #1, #2 ітд. Результати графічні передставлено на выкресах (Рис. 3а–г), вказуючы найчастійше розпознаваны языки в бесіді Русинів, што жыют в Польщы, Словації, Україні і Сербії.

В аналізуваных авдициях з Польщы найчастійше розпознаваный язык (#1) становил 84 проц. полной долгости авдиції, а з того аж 92 проц. становил польскій язык, а 7 проц. українскій язык. Другій найчастійше розпознаваный язык (#2) становил 12 проц., з чого 63 проц. приписано українському языкови, а 21 проц. польському (Рис. 3а).

В авдициях, де бесідуют особы з України (Рис 3б), язык розпознаваный як першый (#1) становил 91 проц. цілой базы, з чого аж 99 проц. тых сегментів приписано українському языкови.

Для Русинів на Словації найчастійше розпознаваный язык становил 66 проц., з чого 77 проц. приписано словацкому языкови. Другій найчастійше розпознаваный язык (#2) становил 19 проц., де польскій і словацкій становили 27 проц., а українскій язык 16 проц. (Рис. 3в).

В примірі авдицій з Сербії (Рис. 3г), перший найчастіше розпознаваний язык (#1) становил 71 проц. награнь, з чого 87 проц. тых авдицій приписано хорватскому языкови. Другій найчастіше розпознаваний язык (#2) становил 22 проц. часу, з чого 78 проц. становил словенський язык.



а)



б)

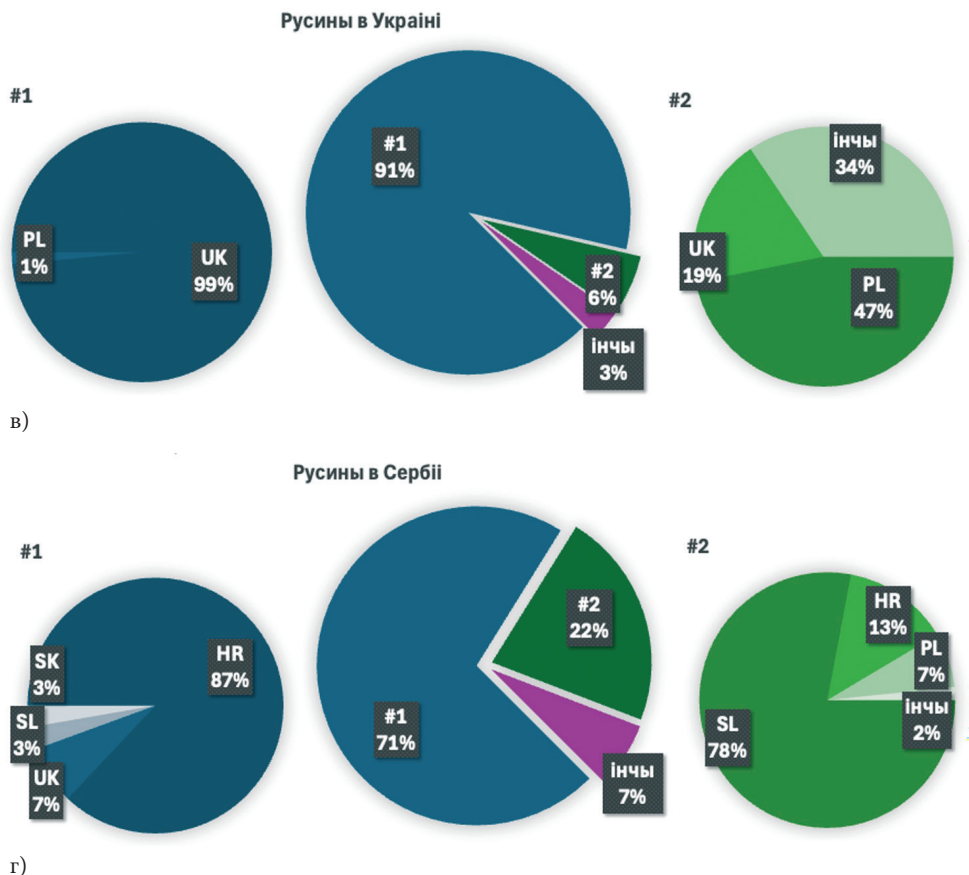


Рис. 3 а–г. Процентове сопоставління розпознаних мов Русинів з різних держав

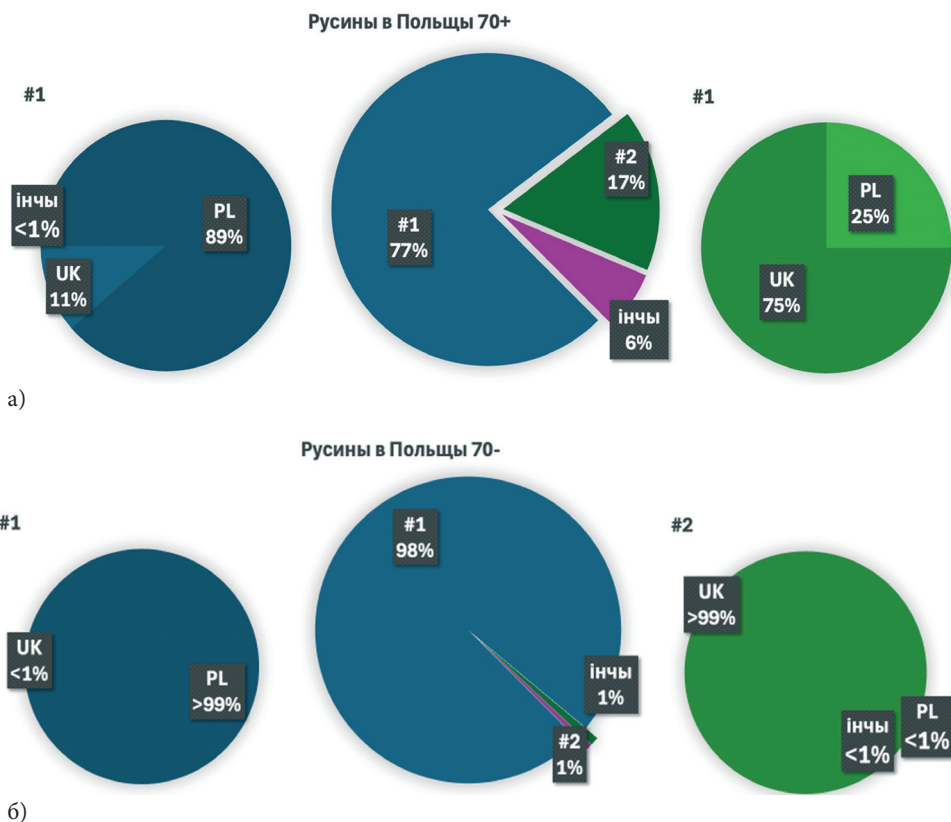
Додатковий аналіз зроблено для Русинів з Польщі. Вивіски (Рис. 4 а, б) представляють аналіз розпознаних аудіцій Русинів з Польщі, котрих поділено на групи за роками: особи старші як 70 років «70+» і особи молодші як 70 років «70-». Треба додати, що база награнь лекторів «70+» становит понад 60 проц. лекторів з Польщі.

Ограничній дослідженні з поділом на групи за роками виключені для лекторів з Польщі, виникат з двох причин. По перше, кількість інших груп була за мала, жебы отримати відповідно великі проби. По друге, проблемом было рішити роки лекторів з інших держав, бо част редакторів, што робили награня, не робит уж в радію.

В групі «70+» (Рис. 4а) найчастіше розпознаваним мовою (#1) был польській, котрий становил 89 проц. аудіцій, обнимаючи 77 проц.

цілковитого аналізованого часу. Як другий найчастіше розпознаваний язык (#2) вкрито українській (75 проц. в тій категорії), што становит 17 проц. вшиткых аналізуваных авдиций. Сут гев тіж присутны інчы языки, яки появляють ся в шестьох процентах розпознань.

В молодшій групі «70-» найчастіше розпознаваним языком (#1) был польскій, обнимаючи веце як 99 проц. часу в тій категорії, з 98 проц. вшиткых аналізуваных авдиций (Рыс. 4б). Інчы языки незначні фігурують в тій групі.



Рыс. 4 а–б. Процентове сопоставліня розпознаваных языків Русинів з Польщы з поділом на групы подля років (веце або менше як 70 років)

## 5. Дискусія

На основі досягнутих результатів явні видно вплива домінуючого языка (урядового) даной державы на класифікацію русинського языка через языковий модель. З вынятком Русинів з Войводіны, бесіда Русинів з України, Польщы і Словаціи была класифікувана через алгоритм як найбарже зближена до домінуючого в даній державі языка. Хоц флексийны, лексикальны і синтактычны ріжниці медже варіантами русинського языка сут з перспектывы його хоснователи або языкознавців рішучо меншы як ріжниці медже русинськым а домінуючым языком даной державы, алгоритм, што комплексово аналізує языкову подібніст, выказує значучо векше зближыня до домінуючых языків. Особливо українській язык розпознаваний єст релятивні часто серед Русинів з Польщы і Словаціи, але в Сербії є інакше. В примірі Русинів з Войводіны найвекшу подібніст выказуют хорватській і словенській, при незначній участі сербского языка. Треба єднак зазначыти, же база пробок з Сербії, на якій оперто аналізу, была найменше чысленна з увагы на ограниченный доступ до материялів.

Не сут знаны докладны механізмы і критеріи, на основі яких модель рішат о класифікації языка. Діяня сіти тыпу *Transformer* єст барз зложеным процесом статистычної аналізу, опертым на ідентифікації взірців в даных языках, што дозволят вызначати подібніст без фактычного розумліня змісту.

Традиційны моделі ASR, основаны на вкрытых моделях Маркова (НММ, англ. Hidden Markov Models), характеризували ся деяком міром детермінізму, што злегшало аналізу вплива акустычных примет на результат транскрипції/класифікації языка. Сучасны моделі ASR, основаны на глубоких невронных сітях, в тым рекурентных невронных сітях (RRN, англ. Recurrent Neural Networks), конволюційных невронных сітях (CNN, англ. Convolutional Neural Networks), і архітектурі тыпу *Transformer*, сут неявны і недетерміністычны. Означат то, же тоты самы входовы даны (сигнал бесіды) можуть в ріжных примірах вести до кус інчых транскрипцій. Тота недетерміністычність выникат з великой кількосты параметрів моделю і зложеных, неленіовых реляцій медже нима. Утруднят то безпосередню аналізу і выокремліня особливых примет акустычного сигналу, які рішают о конкретным результаті

транскрипції. Інакше бесідуючы, тяжко є докладні окрислити, котры фрагменты сигналу бесіды і яким способом мали влияния на клясифікації поедных фонемів ци слів. Істніють єднак методы, што дозволяють на досліджяня і інтерпретатицію моделів ASR, хоц сут зложены і вымагають заавансуваного інформатычного знаня. Належат до них м.ін. (Rahate і ін. 2022):

- аналіза верств увагы (attention mechanism) дозволят візуалізувати звязок медже фрагментами сигналу бесіды а елементами транскрипції, вказуючы, на які части звуку модель «зверат увагу» в час розпознаваня,
- методы основаны на градієнтовым пошырюваню, дозволяють ідентифікувати фрагменты сигналу на вході, які мают найвекше влияния на активацию окремых нейронів і в тот спосіб на кінцьовый результат транскрипції,
- методы пертурбації входу основаны сут на тым, же вводит ся невеликы зміны в сигналі бесіды і обсервує, як влияют на результат транскрипції, што дозволят ідентифікувати основны акустычны приметы.

Результаты вказуют праві нульове вызначаня російского языка, што дакус засакаує, беручы до увагы його сутьове влияния і хоснуваня в Україні, особливо на сході державы, одкале походила част бесідуючых осіб.

Аналіза лемківского языка выказала тіж, што результаты ріжнят ся в залежности од років бесідуючых осіб. В звязку з тым проведено другу аналізу. Поділено пробкы на дві групы подля років: бесідуючы особы, што мают менше і веце як 70 років. Подібніст лемківского языка до польского серед старшых осіб была значучо менша як серед молодшых осіб, што сугерує поступуючу языкову асиміляцію, особливо в обшыри фонетыкы. В лексикальным аспекті старшы бесідуючы особы хоснували подібне, аж і векше чысло польонізмів або запожычынь, што може вказувати на зложены механізмы языковой асиміляції медже поколінями.

## 6. Підсумуваня

В статі проаналізувано ци модель OpenAI Whisper годен розпознавати і клясифікувати русинській язык, хоц модель не был непосредньо вчений на русинських материялах. Проведено досліджыня на основі

радийовых авдиций Русинів з Польщы, Украіны, Словаці і Сербії, при увзгляднію ріжниц медже віковыма групами. Результаты вказуют, же модель OpenAI Whisper класифікує русинській язык головні як урядовый язык даной державы, што насуват высновак о сутьовым влияню домінуючых языків на автоматичну класифікацію меншынового языка. Достережено тіж значучу языкову асиміляцію в молодшій групі хоснуватели русинського языка, што свідчыт о поступуючым влияню польского языка в Польщы і інчых урядовых языків в інчых державах.

В будучых досліджынях плянує ся змодифікувати модель OpenAI Whisper так, жебы выеліминувати спомагання для домінуючого языка в механізмі розпознаваня, што дозволит ліпше зрозуміти подібніст русинського языка до інчых славянських языків без сильного влияня офіційного языка даной державы. Додатково плянує ся проаналізувати русинскы пісні і співанкы а тіж інчы славянскы пісні. В співі выступуют інчы механізмы емісії звуку як в бесіді, што може вказати шыршый контекст. Сесы досліджыня можут причынити ся до барже прецизийной і актуальной класифікації меншыновых языків і ліпшого зрозумліня механізмів асиміляції і языковой глобалізації.

Вшыткы первістны даны і дрібницьовы результаты вызначены проведеным експериментом сут доступны для заінтересуваных через майльовый контакт з автором.

## Бібліографія

- Bouamor, Houda, Hassan, Sabit, Habash, Nizar. 2019. «The MADAR Shared Task on Arabic Fine-Grained Dialect Identification». В: *Proceedings of the Fourth Arabic Natural Language Processing Workshop*. Ред. Wassim El-Hajj, Lamia Hadrich Belguith, Fethi Bougares, Walid Magdy, Imed Zitouni, Nadi Tomeh, Mahmoud El-Haj, Wajdi Zaghoulani, 199–207. Florence: Association for Computational Linguistics. <https://doi.org/10.18653/v1/W19-4622>.
- Kushko, Nadiya. 2007. «Literary Standards of the Rusyn Language: The Historical Context and Contemporary Situation». *The Slavic and East European Journal* 51, ч. 1: 111–132.
- Moser, Michael. 2016. «Rusyn: A New-Old Language In-between Nations and States». В: *The Palgrave Handbook of Slavic Languages, Identities and Borders*. Ред. Tomasz Kamusella, Motoki Nomachi, Catherine Gibson, 124–139. London: Palgrave Macmillan. [https://doi.org/10.1007/978-1-137-34839-5\\_7](https://doi.org/10.1007/978-1-137-34839-5_7).
- Nikitin, Alexey G., Kochkin, Igor T., June, Cynthia M., Willis, Catherine M., Mcbain, Ian, Videiko, Mykhailo Y. 2009. «Mitochondrial DNA Sequence Variation in the Boyko,

- Hutsul, and Lemko Populations of the Carpathian Highlands». *Human Biology* 81, ч. 1: 43–58. <https://doi.org/10.3378/027.081.0104>.
- Plišková, Anna. 2008. «Practical Spheres of the Rusyn Language in Slovakia». *Studia Slavica Academiae Scientiarum Hungaricae* 53, ч. 1: 95–115. <https://doi.org/10.1556/SSLav.53.2008.1.6>.
- Rabus, Achim, Scherrer, Yves. 2017. «Lexicon Induction for Spoken Rusyn – Challenges and Results». В: *Proceedings of the 6th Workshop on Balto-Slavic Natural Language Processing*. Ред. Tomaž Erjavec, Jakub Piskorski, Lidia Pivovarova, Jan Šnajder, Josef Steinberger, Roman Yangarber, 27–32. Valencia: Association for Computational Linguistics. <https://doi.org/10.18653/v1/W17-1405>.
- Radford, Alec, Kim, Jong Wook, Xu, Tao, Brockman, Greg, McLeavey, Christine, Sutskever, Ilya. 2023. «Robust Speech Recognition via Large-Scale Weak Supervision». В: *Proceedings of the 40th International Conference on Machine Learning (ICML'23)*. Ред. Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, Jonathan, 1–28 (28492–28518). Honolulu: JMLR.org.
- Rahate, Anil, Walambe, Rahee, Ramanna, Sheela, Kotecha, Ketan. 2022. «Multimodal Co-learning: Challenges, Applications with Datasets, Recent Advances and Future Directions». *Information Fusion* 81: 203–239. <https://doi.org/10.1016/j.inffus.2021.12.003>.
- Scherrer, Yves, Rabus, Achim. 2019. «Neural Morphosyntactic Tagging for Rusyn». *Natural Language Engineering* 25, ч. 5: 633–650. <https://doi.org/10.1017/S1351324919000287>.
- Zampieri, Marcos, Nakov, Preslav, Scherrer, Yves. 2020. «Natural Language Processing for Similar Languages, Varieties, and Dialects: A Survey». *Natural Language Engineering* 26, ч. 6: 595–612. <https://doi.org/10.1017/S1351324920000492>.