

Designing a Speaking Assessment Module for the SELF Language Placement Test

WORK IN PROGRESS

Pinxun HUANG, Eli STAFFORD, Sylvain COULANGE

LIDILEM (*Laboratory of Linguistics and Didactics of Foreign and Mother Tongues*)

LIG (*CNRS, Institute of Engineering, Grenoble Computer Science Laboratory*)

Université Grenoble Alpes

SELF - Formative Language assEssment System

Online placement test developed at Université Grenoble Alpes



SELF
**INNOVA
LANGUES**

A screenshot of the SELF assessment interface. On the left, there is a speech bubble icon and the text "Reste 2 écoutes". On the right, a green panel contains a question icon (a drop with a question mark and smile) and three radio button options labeled 1, 2, and 3. Option 2 is selected. Above the green panel, there are two progress indicators: "1/2 OK" and "2/2 OK". Below the green panel is a "Valider" button. At the bottom of the interface, a progress bar shows 40% completion.

SELF - Formative Language assEssment System

Online placement test developed at Université Grenoble Alpes



7 languages

(English, Italian, Mandarin, Japanese, Spanish, French, German)



~1h



A1-C1



~40 000 tests per year



~40 universities

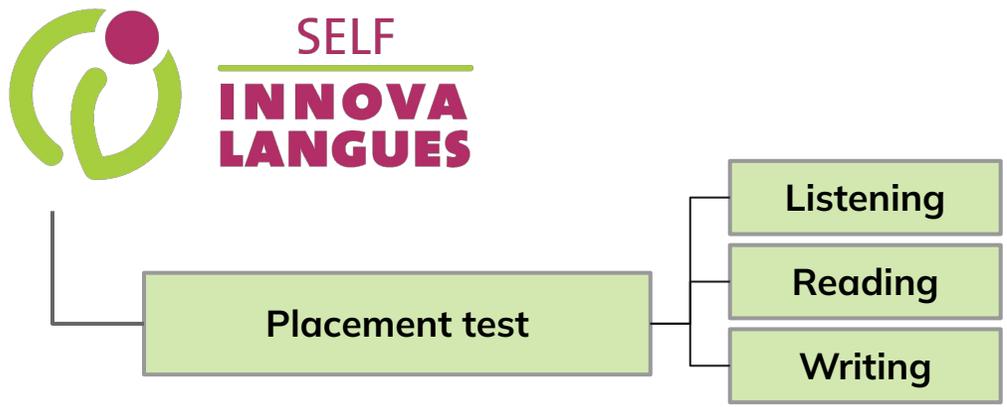


adaptive algorithm



More info here: <https://self.univ-grenoble-alpes.fr/>

SELF - Formative Language assEssment System



Reste 1 écoute

1/1

1
2
3

ブログ

Léa Girard

1/1

レアさんは、来週から何をしますか。

- 日本に旅行します
- フランス語を説明します
- 外国語で案内します

Reste 2 écoutes

1/2 2/2

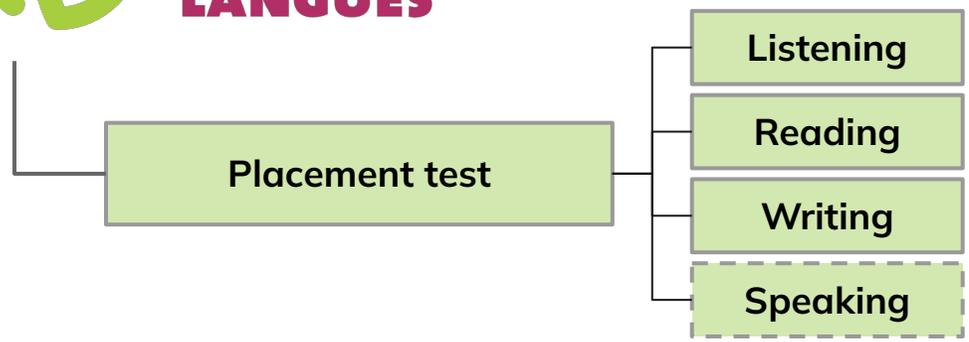
Recette de cuisine sur votre blog

Soupe à l'oignon

- c les oignons en rondelles
- faire cuire les oignons dans une casserole avec du beurre et de l'huile puis a de la farine
- verser le vin blanc dans la casserole
- couvrir la casserole pendant 1h30 et mélanger de temps en temps en tournant avec une cuillère

Bon appétit ! 🍴

SELF - Formative Language assEssment System



The interface displays various assessment tasks:

- Reste 1 écoute:** A listening exercise with a speech bubble icon and a progress indicator of 1/1.
- Blog:** A reading comprehension task featuring a blog post by Léa Girard about her experience in Japan. It includes a progress indicator of 1/1 and a question: "レアさんは、来週から何をしますか。" (What will Léa do starting next week?).
- Reste 2 écoutes:** A listening exercise with a video player showing Catherine Taylor and a progress indicator of 2/2.
- Recette de cuisine sur votre blog:** A reading task with a recipe for "Soupe à l'oignon" (Onion Soup). It includes a progress indicator of 2/2 and a question: "c _____ les oignons en rondelles" (cut the onions into rounds).

SELF - Speaking modules (English & French)

Developing a Speaking assessment module for SELF placement tests
May 2025 - Nov. 2026



Alex Carr
Instructional
Design
SELF English



Alireza Pournouri
Instructional
Design
SELF English



Pinxun Huang
Speech
Processing



Eli Stafford
Speech
Processing



Anne-Cécile Perret
Coordinator
SELF Français



Sylvain Coulange
Technical and Scientific
Coordinator

SELF - Speaking modules (English & French)

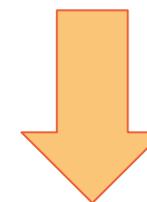
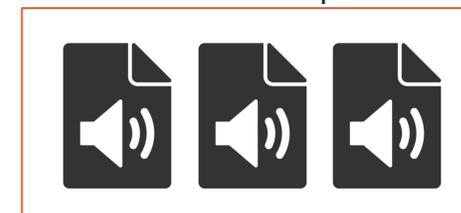
Developing a Speaking assessment module for SELF placement tests
May 2025 - Nov. 2026

1/1

Record

Save

student = Marine Dupont



Speaking score: B1

CO + CE + EEC + PO = "On the way to B2.1"

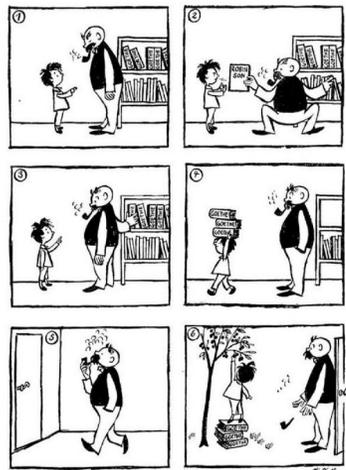
Task Examples



You've just met Jonathan, describe him to your friend.



Send a voice message to your friend in response. Describe the activities you did. You can use the pictures below if you want.



Take a look at these pictures and tell the story.



"Excuse me... Hi! We're doing a report on what people think of what the Prime Minister said last night about climate change. He said quote "The government should not be held responsible for climate change. It's the individual citizens that must be held accountable." Do you agree with this statement?"

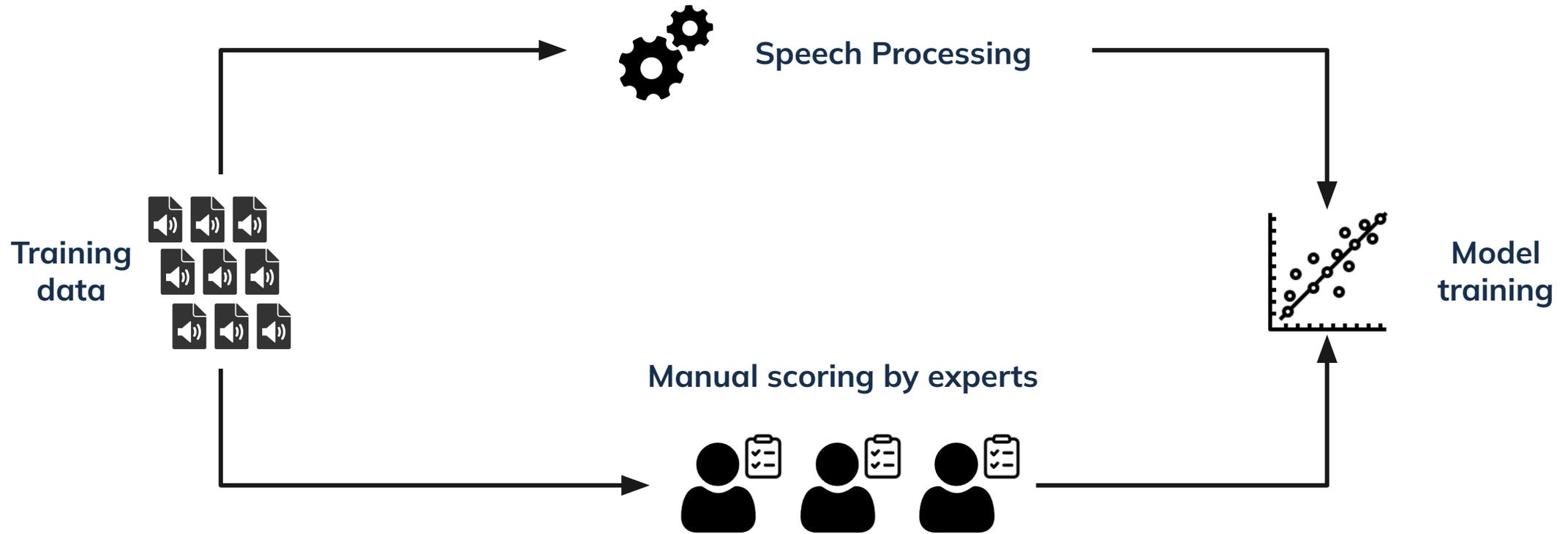
Answer the reporter's question and explain your reasons.



"Hey. So, I talked to my parents yesterday and they think it's not a good choice for me to go to art school next year. They said I should choose a profession that would be stable and that I should go to medical school to become a doctor like them. But I don't know, painting is my biggest passion, you know? So I don't really know what to do right now. What do you think?"

Help Alex making a choice, giving her advice about the pros and cons of both choices.

Training a Score prediction model



Using the model to predict scores on new data



Manual Scoring by Experts

- 12 raters recruited
- 150 recordings / rater
- 600 recordings
- ~3 ratings / recording
- = 1800 ratings

SELF PO Rating Interface IGOR LOGOUT

Progress: 85 / 150 recordings

Task name: Describe_Jonathan Task code: 10520

Context 0:00 / 0:02

Jonathan, 25

Question 0:00 / 0:04

Recording ID: 1762545566 Student ID: 133508

Recording to Rate 0:00 / 0:34 100%

Rating Form

Non-scorable
 Off-topic

Task Completed NO PARTIALLY YES

Pragmatics A1 A2 B1 B2 C1 C2

Vocabulary and Grammar A1 A2 B1 B2 C1 C2

Pronunciation / Prosody / Fluency A1 A2 B1 B2 C1 C2

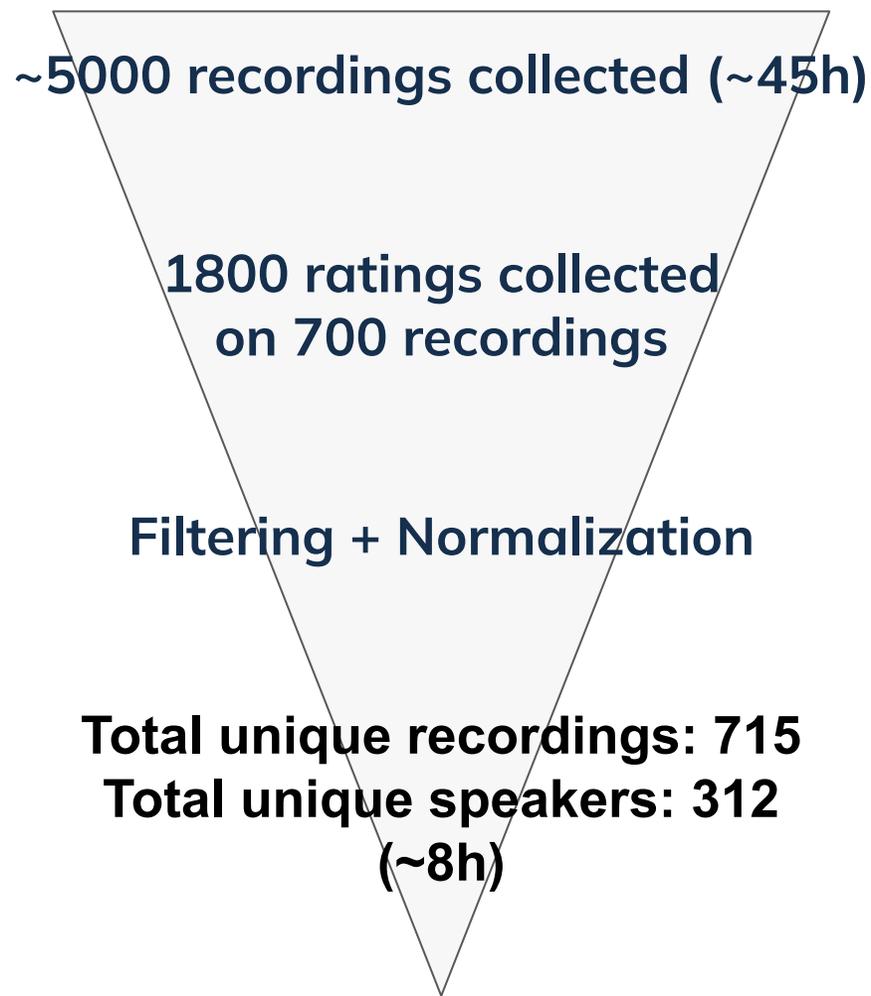
Recording's Global Level A1 A2 B1 B2 C1 C2
According to you, what is the proficiency level of this recording?
Value: 4

Rating Confidence Score LOW MEDIUM HIGH

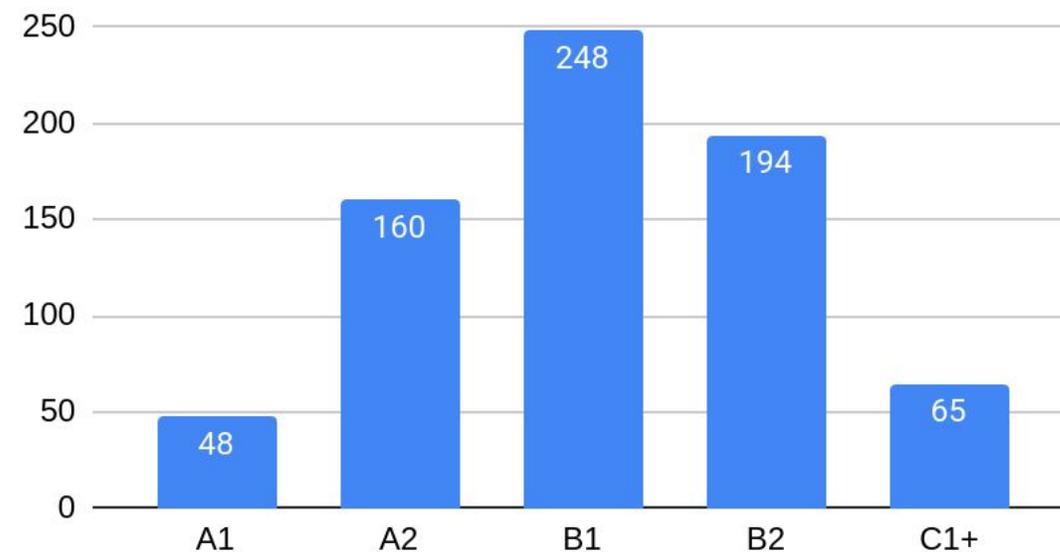
Comments (optional)
 Add any additional comments...

SUBMIT RATING AND LOAD NEXT

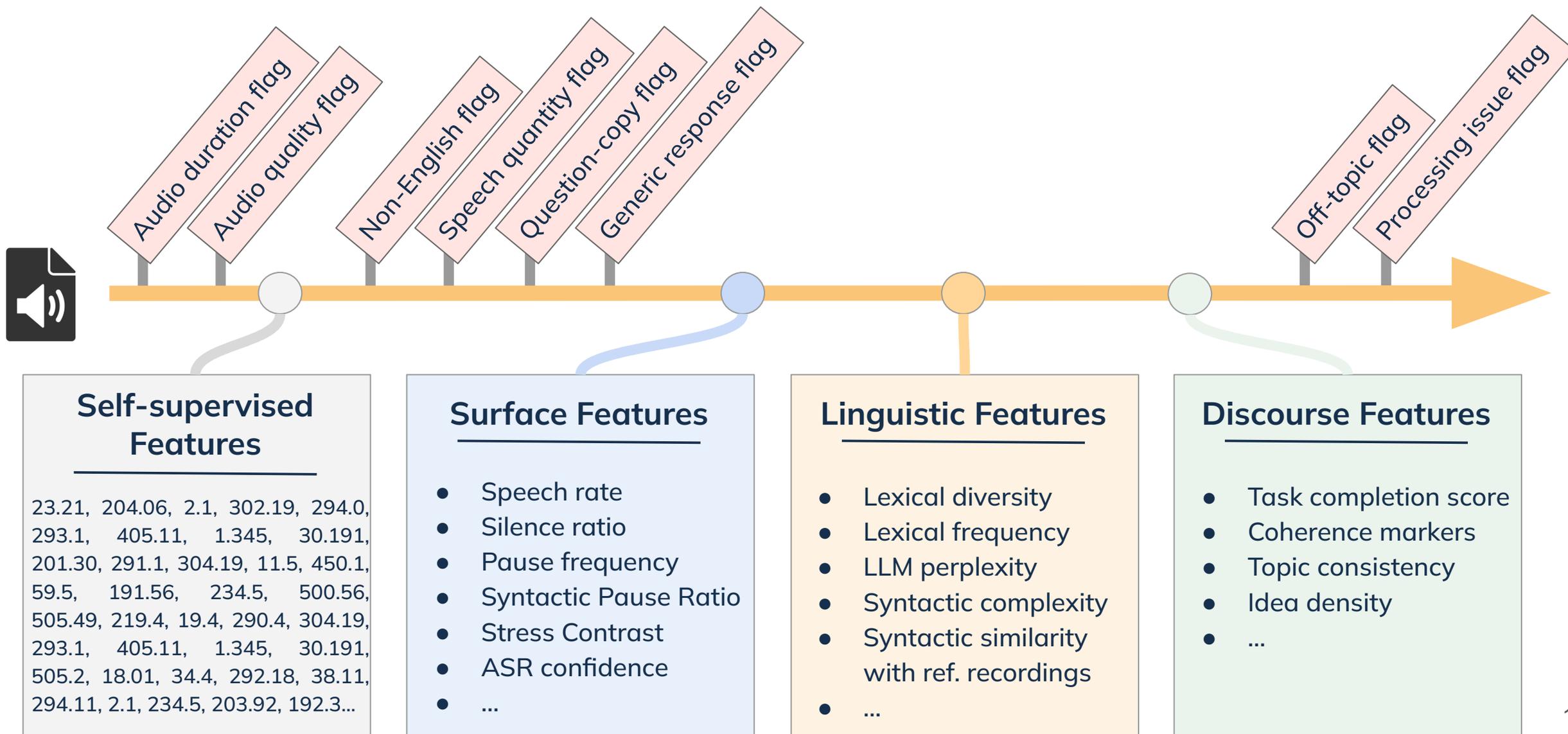
Final Dataset (Recordings with Human Scores)



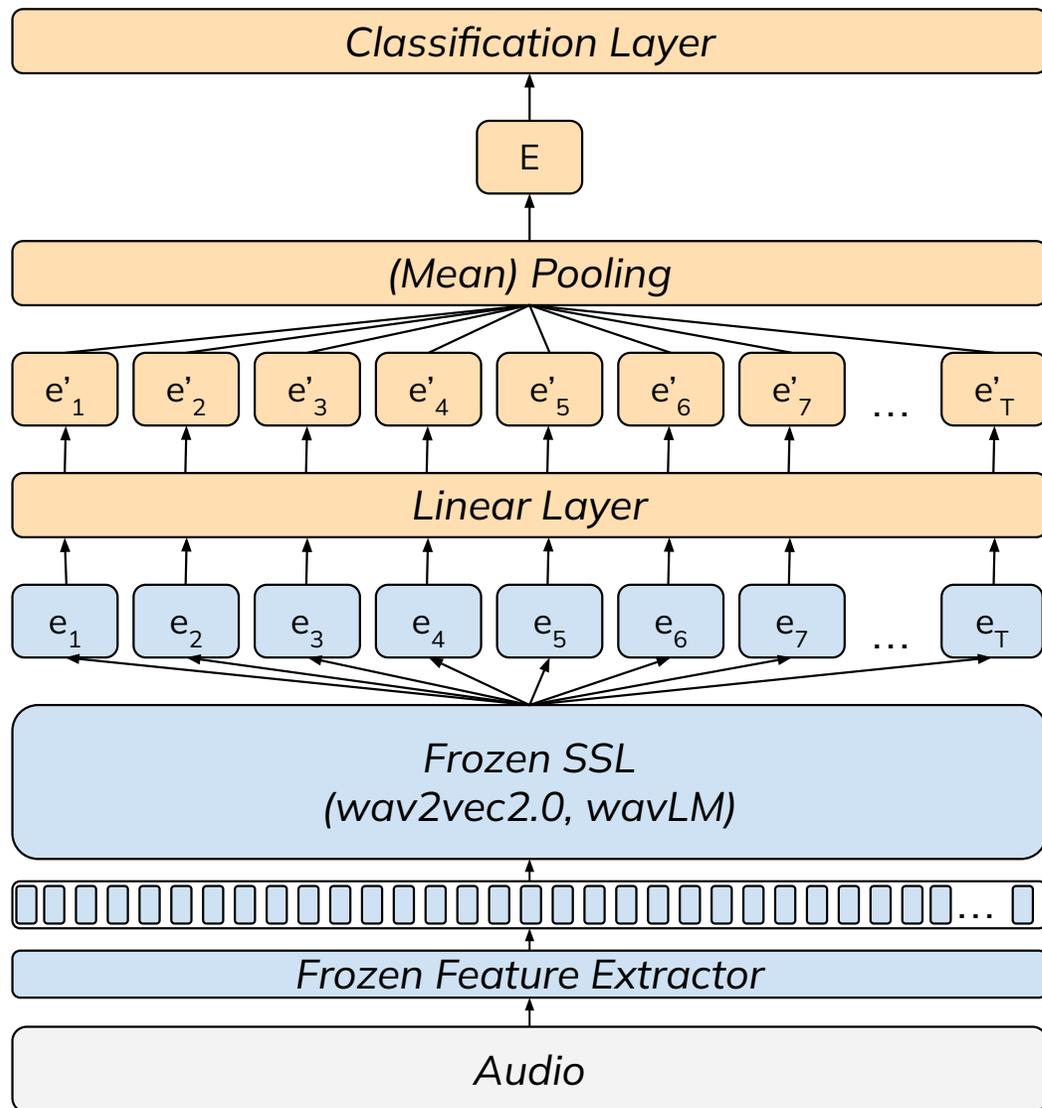
Number of recordings with gold score



Evaluation Criteria (for English)



Using Pre-Trained Speech Embeddings



Pre-Trained Embeddings

- Models like wav2vec2 and wavLM, which produce sequences of vectors representing slices of audio, trained on unlabeled specs
- Low explicability

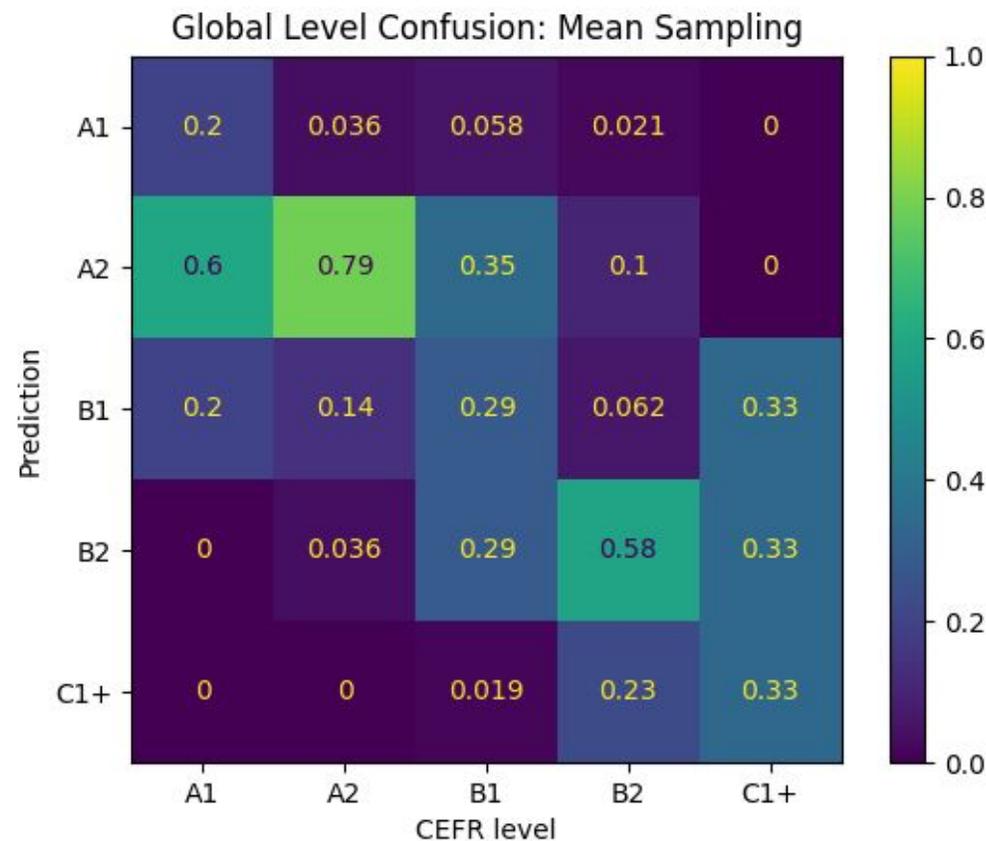
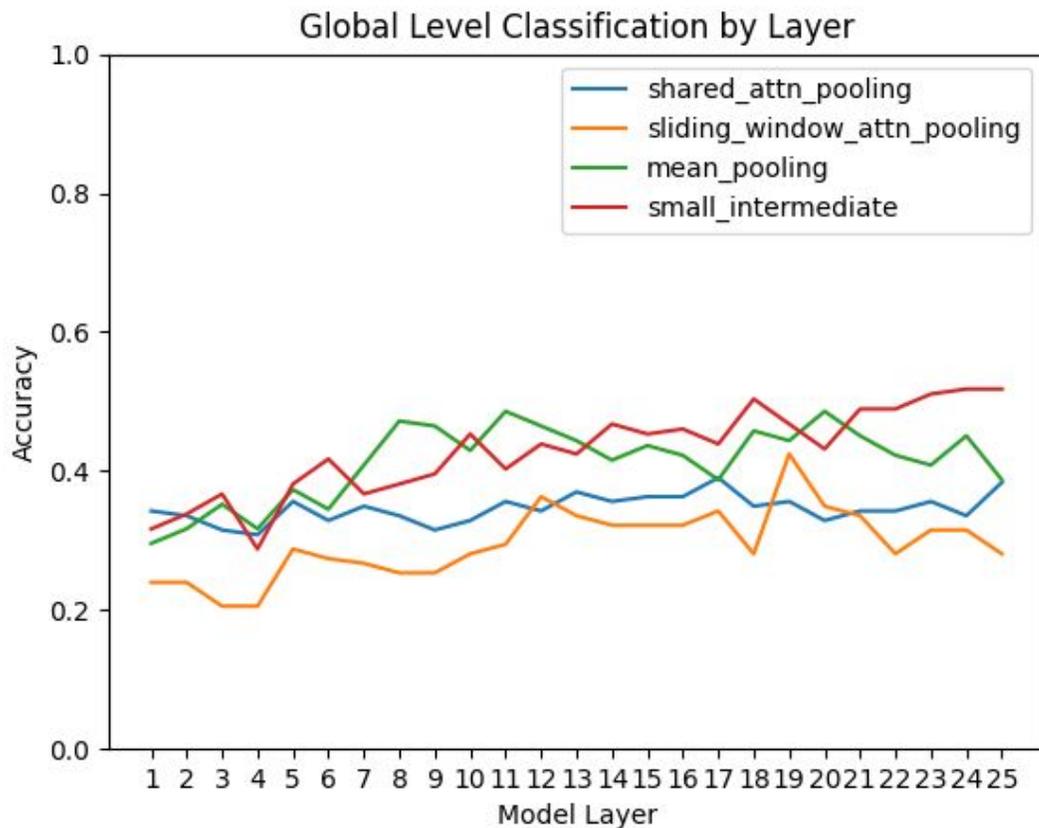
End-to-End Task Architecture

- Frozen SSL + (pooling +) classifier head
 - Leverage SSL embeddings with lightweight classifier layers
- Utterance level classification (~Emotion Recognition)

Implications of pooling

- Lack of generic sentence-vector-like encoders for speech, usually task specific
- Unlikely to retain word level linguistic information (pronunciation, stress, semantics)
- Proven performance on non-linguistic tasks (Emotion recognition, speaker recognition)

Global Level Results - wavlm-large



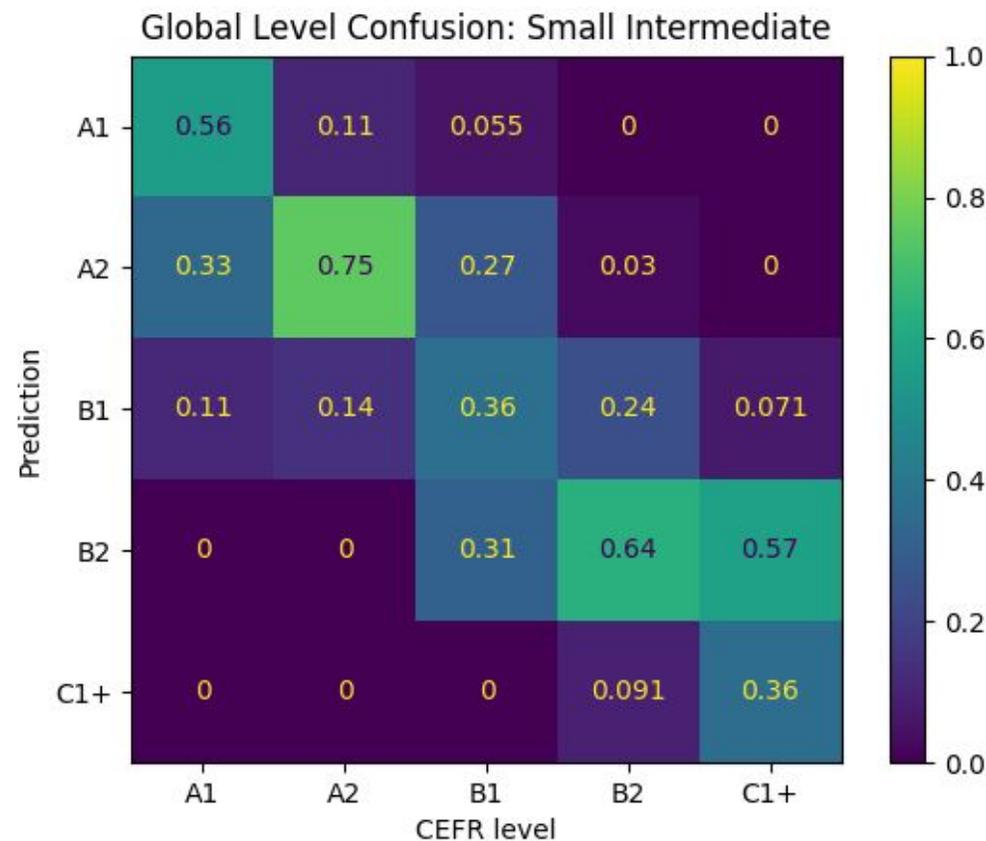
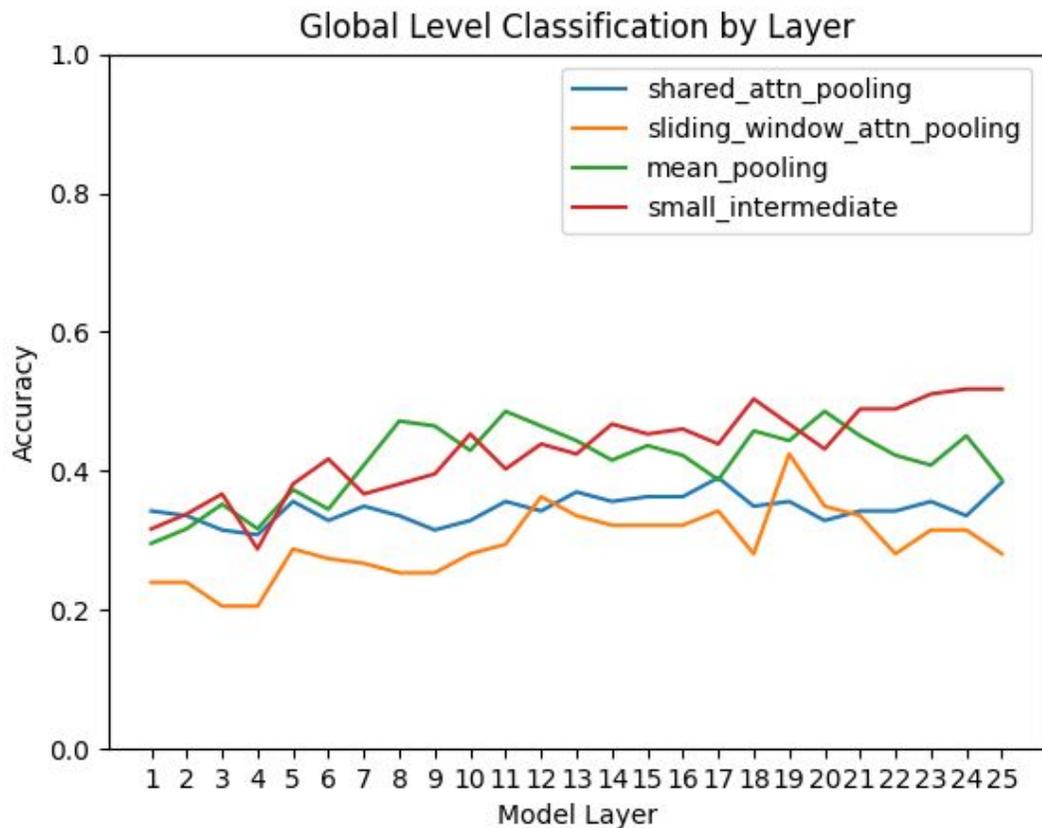
Best Results: (naive baseline: 39.5%)

Layer 11 with weighted sampling

Accuracy: 48.5% Macro (unweighted) F1: .397

W/ ± 1-level: 82.0% Train Acc : ~85%

Global Level Results - wavlm-large

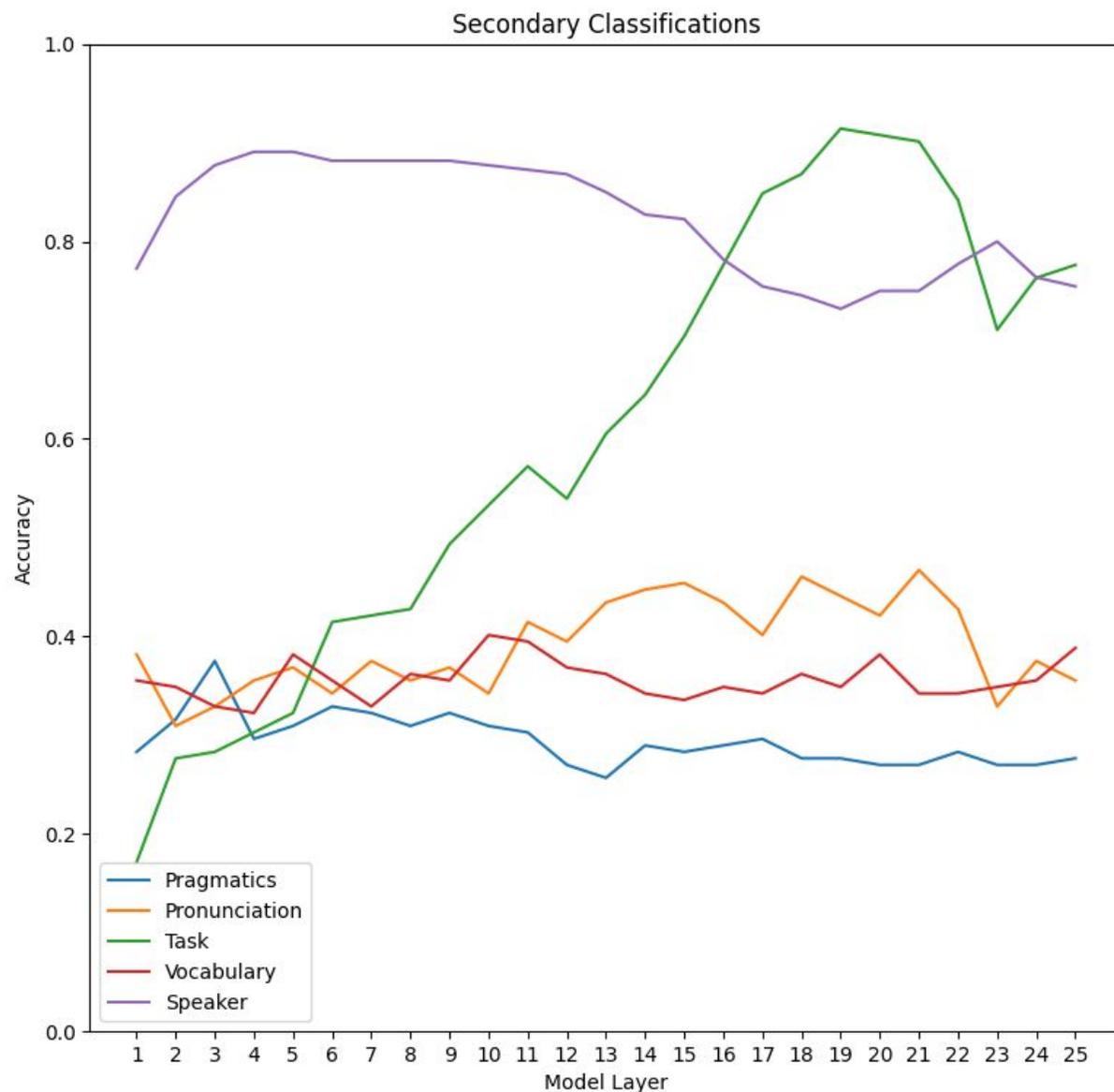


Best Results:

Layer 25 with small intermediate

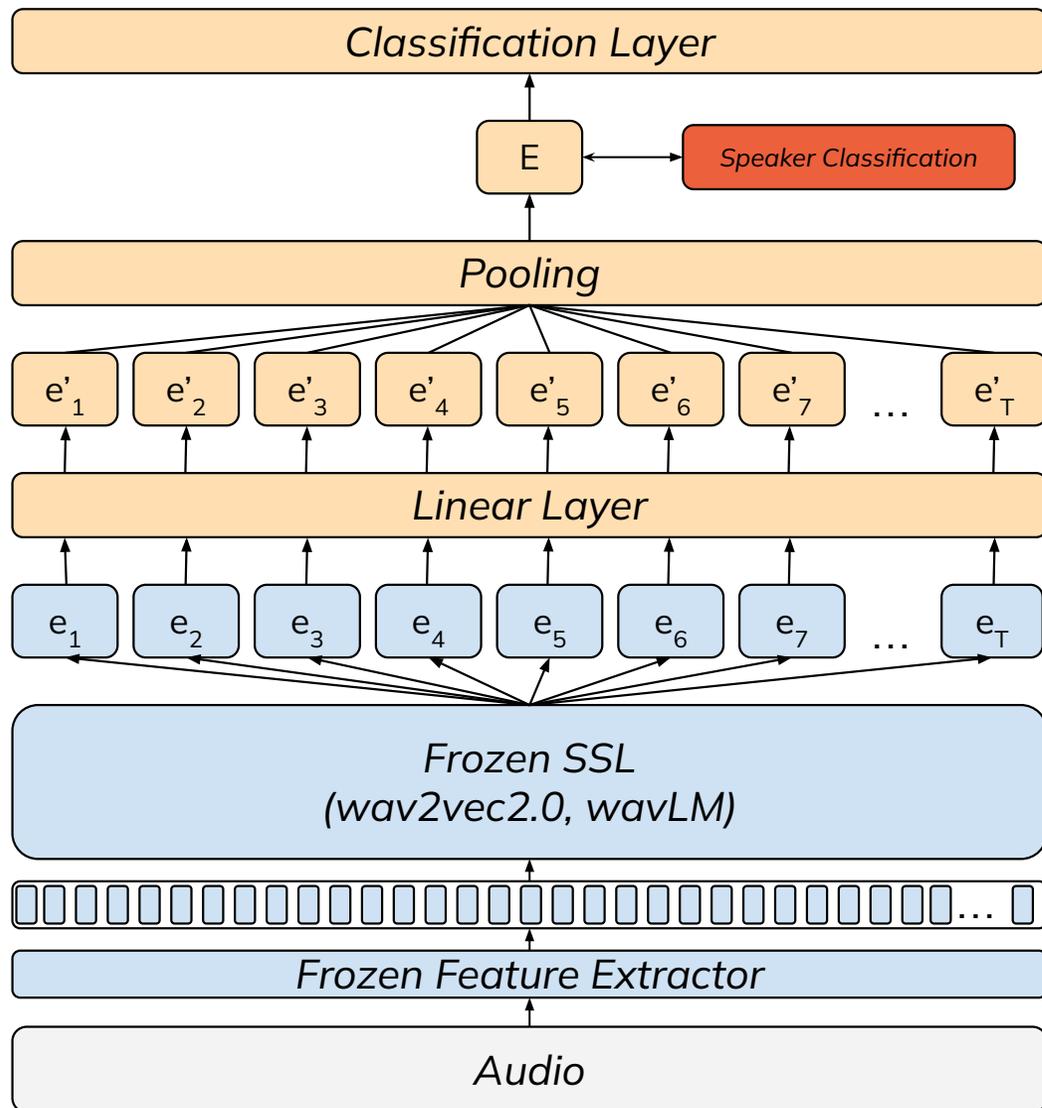
Accuracy: 51.8% Macro (unweighted) F1: .511
 W/ ± 1-level: 89.4% Train Acc: ~75%

Secondary Characteristics



- Stronger predictions on pronunciation than pragmatics or vocabulary
- No significant strengths compared to the score based classifier (with this architecture)
- Much better at speaker identification
 - Can rely on this for level identification
 - With many single example speakers, this ability can hurt generalization
- Task identification
 - o Keyword identification

Adversarial Training



Penalize intermediate representations which allow for speaker classification

Algorithm

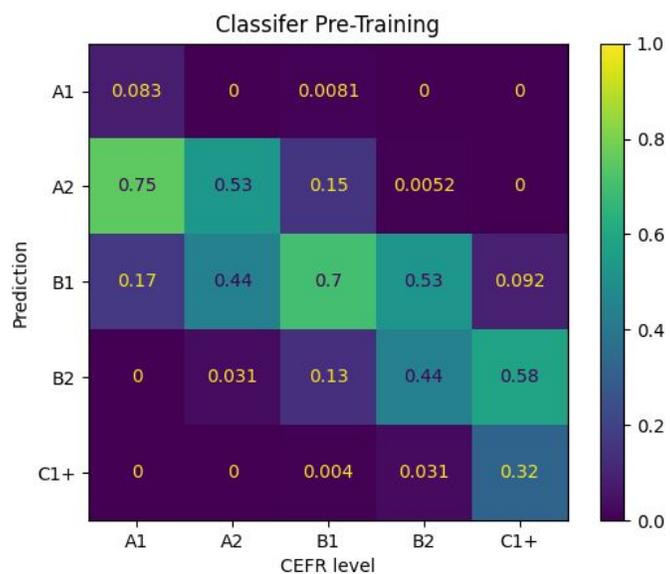
- Pre-train base classification
- Pre-train a speaker classification head on the intermediate embedding
- Perform adversarial training (Hardt et al, 2016)
- Fine-tune classification layer on new intermediate embeddings

Ensures both Equality of Opportunity (as described in Hardt et al) and better generalization across speakers

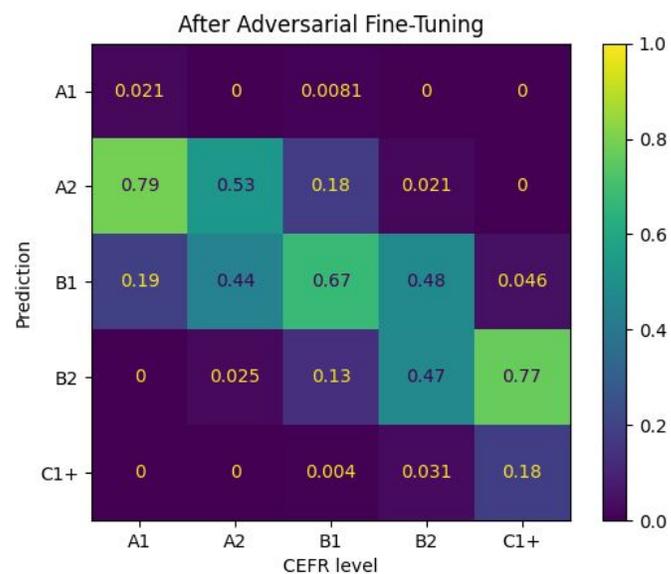
Preliminary Results

10-fold cross validation
 wavLM layer 25
 Intermediate embedding of size 64
 Heavily parameter dependant

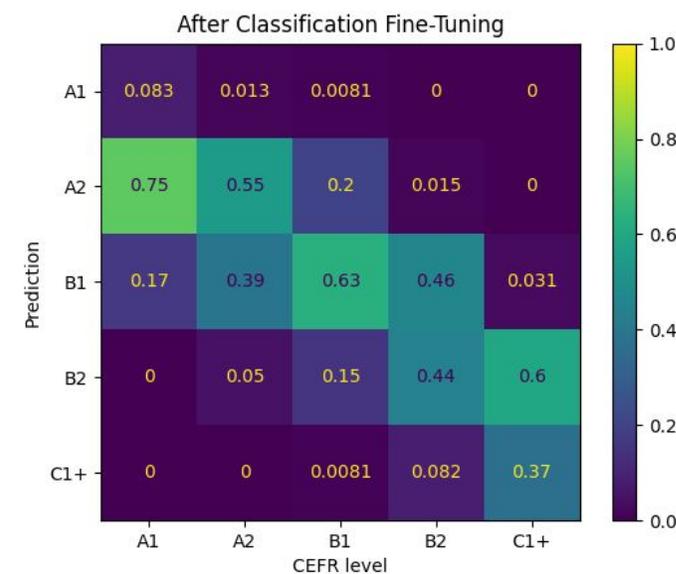
No change in generalization power
 Strong theoretical guarantees



Accuracy: 51.6%
 Macro F1: .436

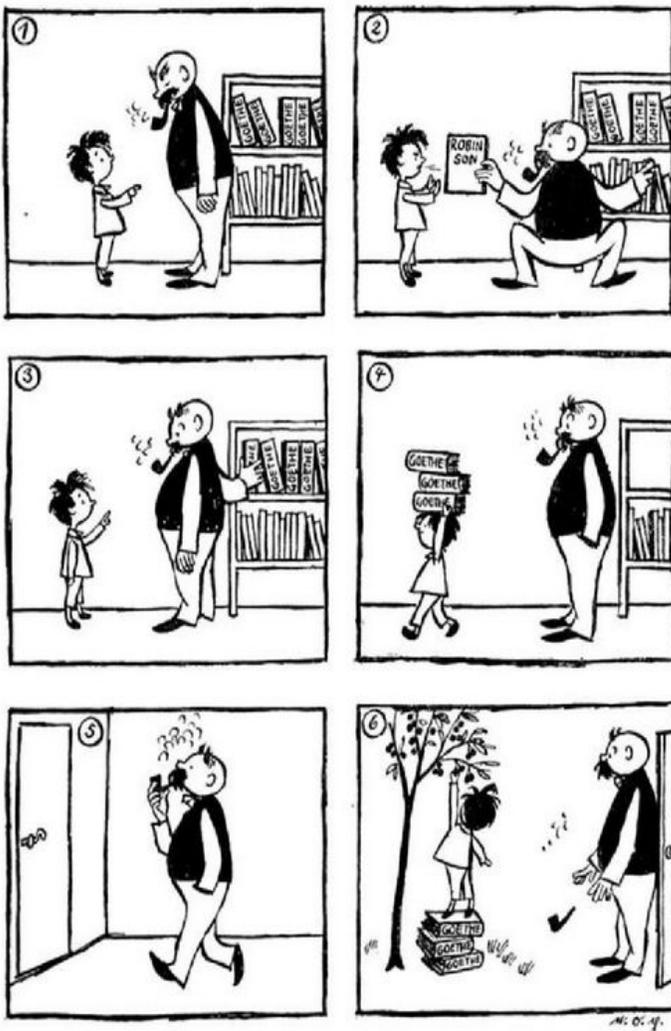


Accuracy: 49.8%
 Macro F1: .378



Accuracy: 50.2%
 Macro F1: .428
 Train Acc: ~63%

On-Topic or Off-Topic ?



Task:

Telling a story. Take a look at these pictures and tell the story.

A) Student answer:

A child asks some books to a man to put them in the garden and reach a branch of a tree.

B) Student answer:

A kid walks up to a man asking for a book the man hands him what looks like a comic book the kid refuses saying he wants bigger books the man surprised gives him three books when he walks out he sees that the kid was only using them to reach a higher place.

C) Student answer:

It's a children and his parents and the parents give some knowledge to his children and the children use this knowledge to arrive at his objective of thinking.

Approaches to Off-Topic Detection

1) GPT2 Likelihood

- Convert Task Question and Image into Text Prompt
- Measure how likely the Student's Answer is given the Prompt for a Large Language Model-GPT2

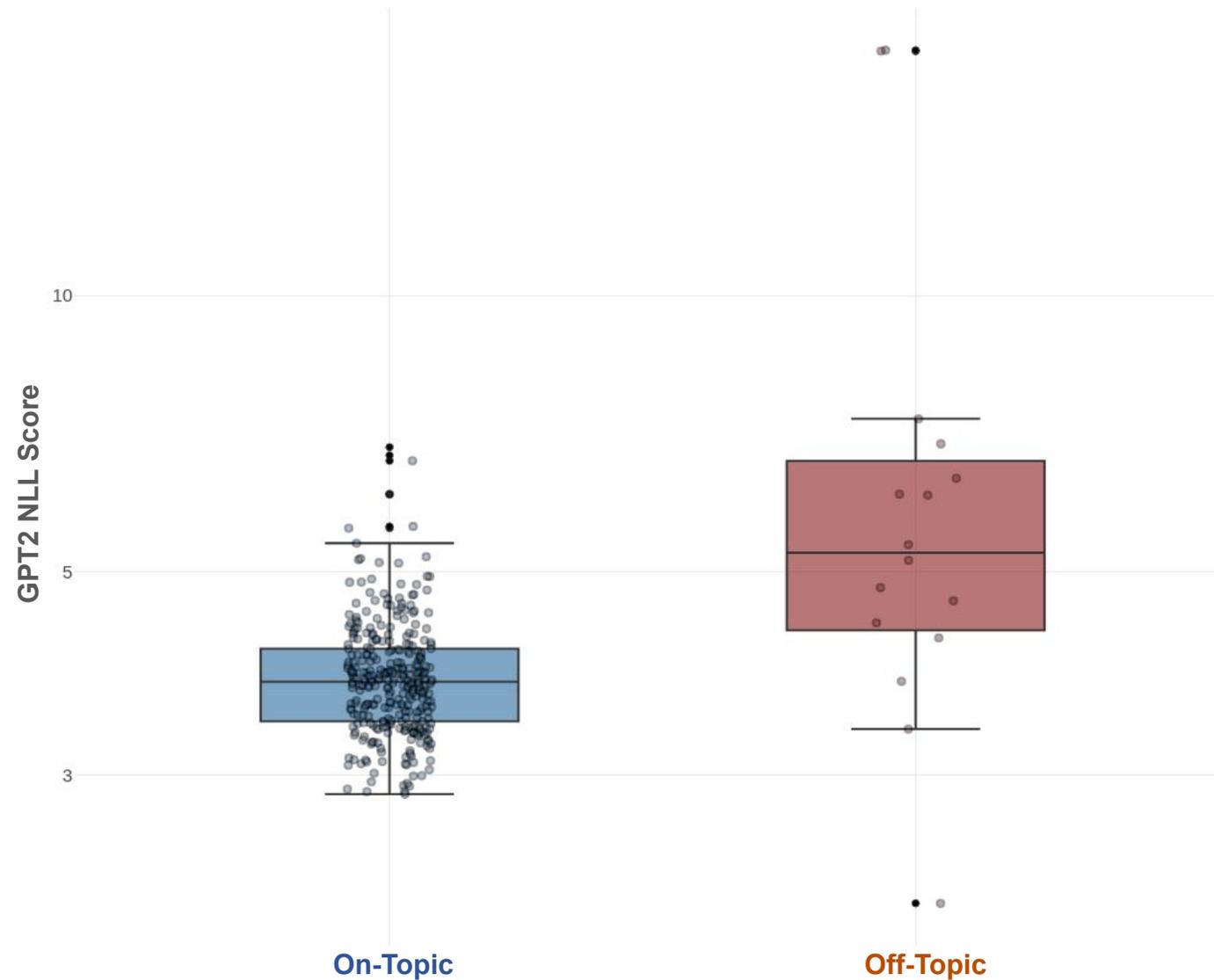
2) Sentence-BERT Semantic Similarity

- Measure the Semantic Similarity between the Prompt and the Student's Response

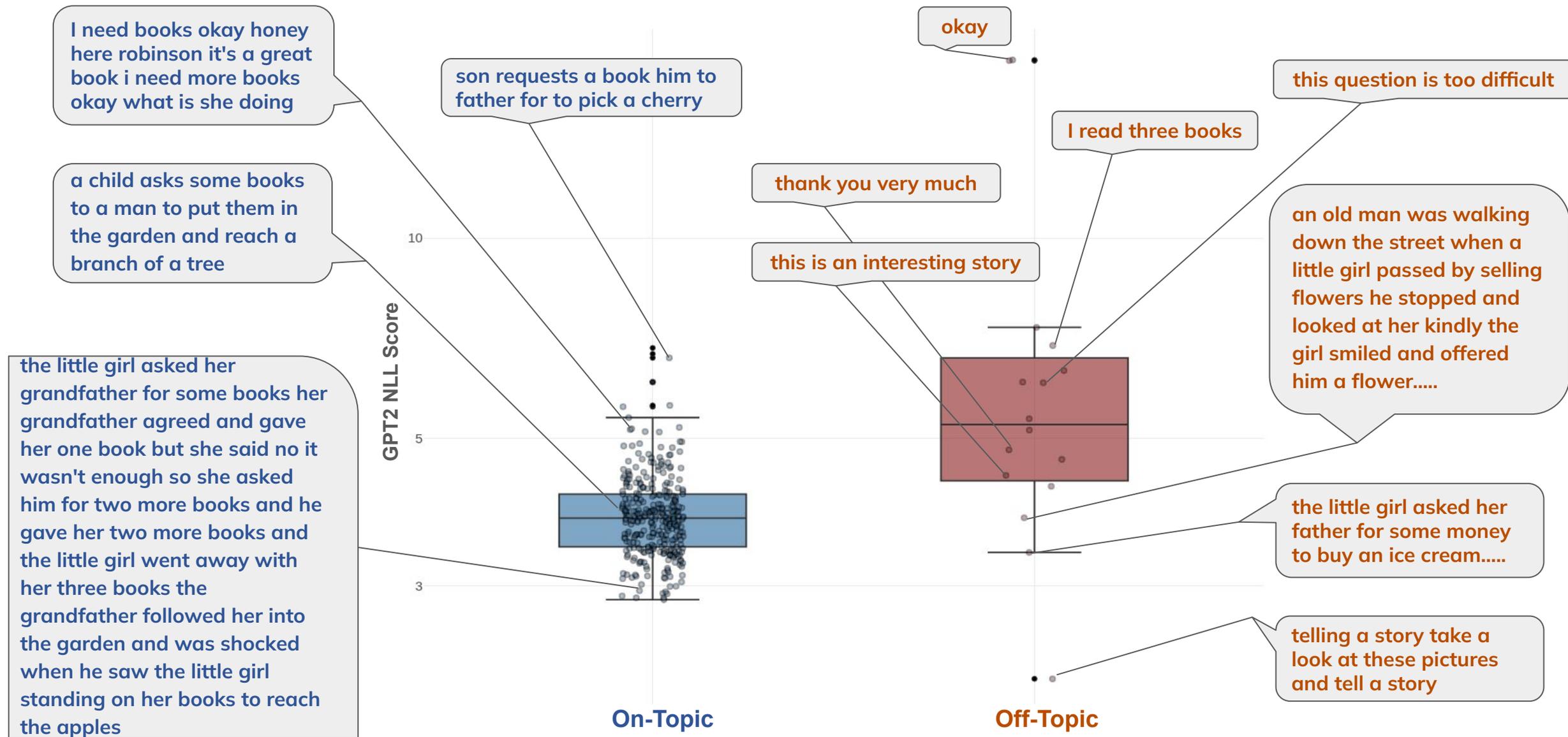
3) LLM as a Judge

- Give Task Question and Student's Response to LLM-Mistral
- Use LLM to judge Response Relevance

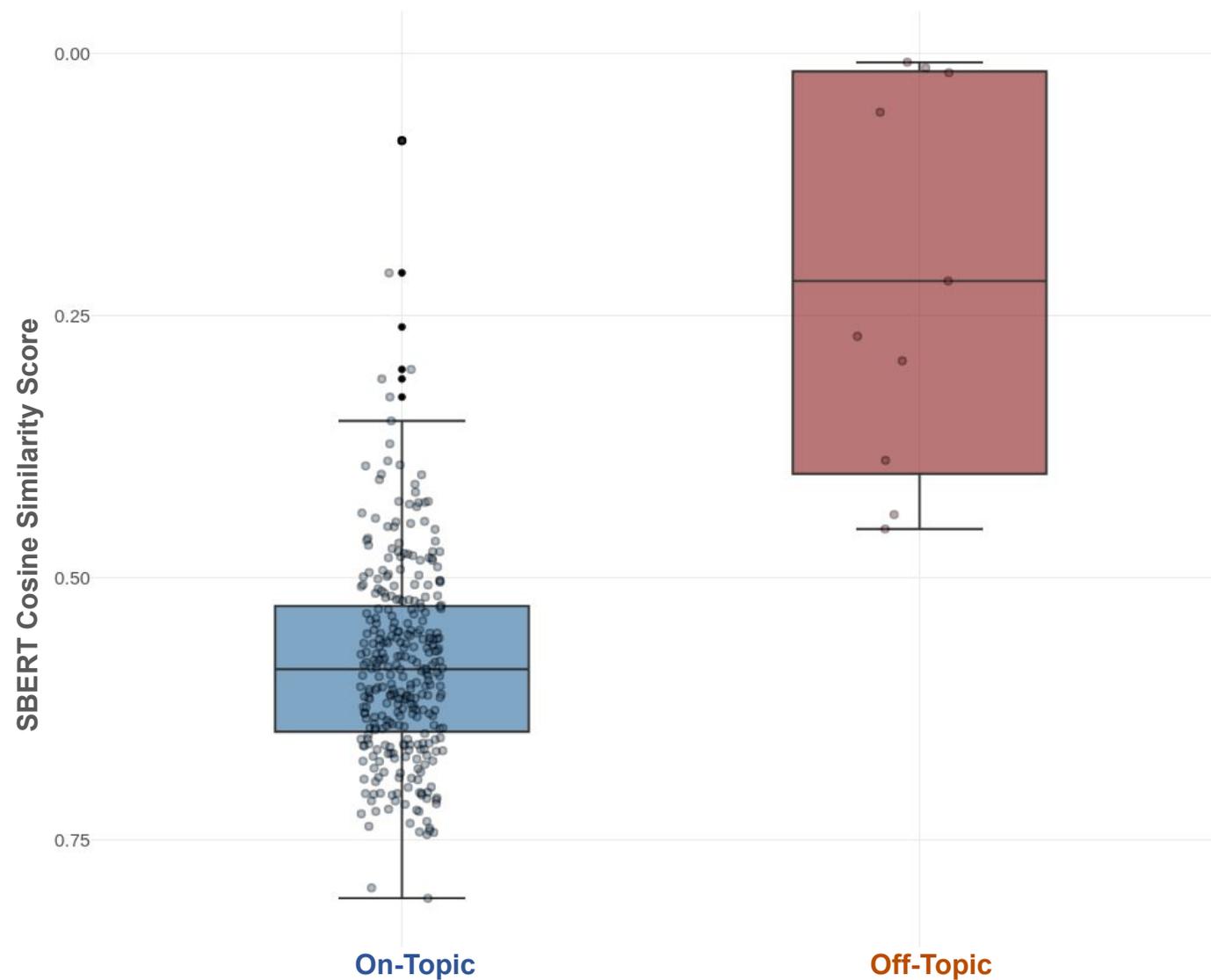
1) GPT2 Likelihood (how likely the response given the prompt)



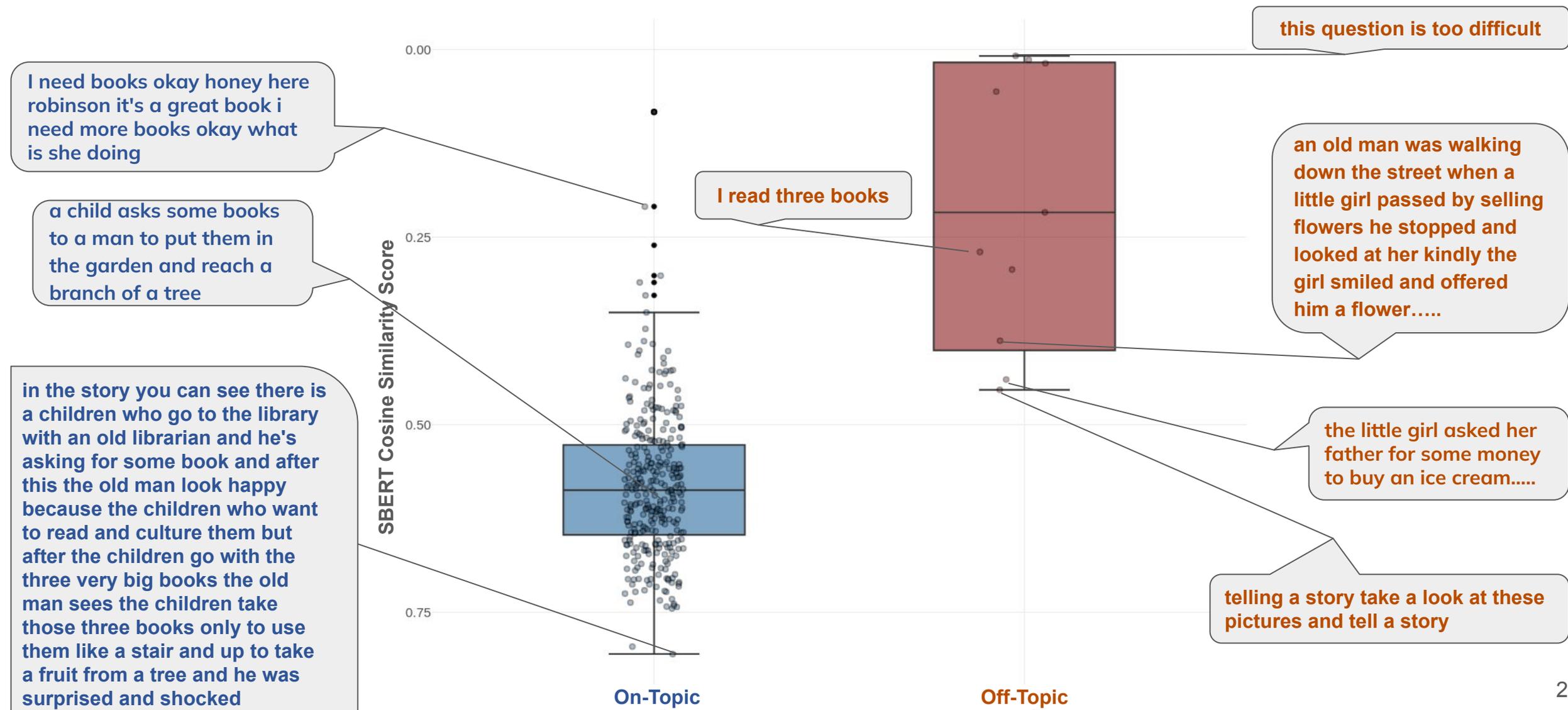
1) GPT2 Likelihood (how likely the response given the prompt)



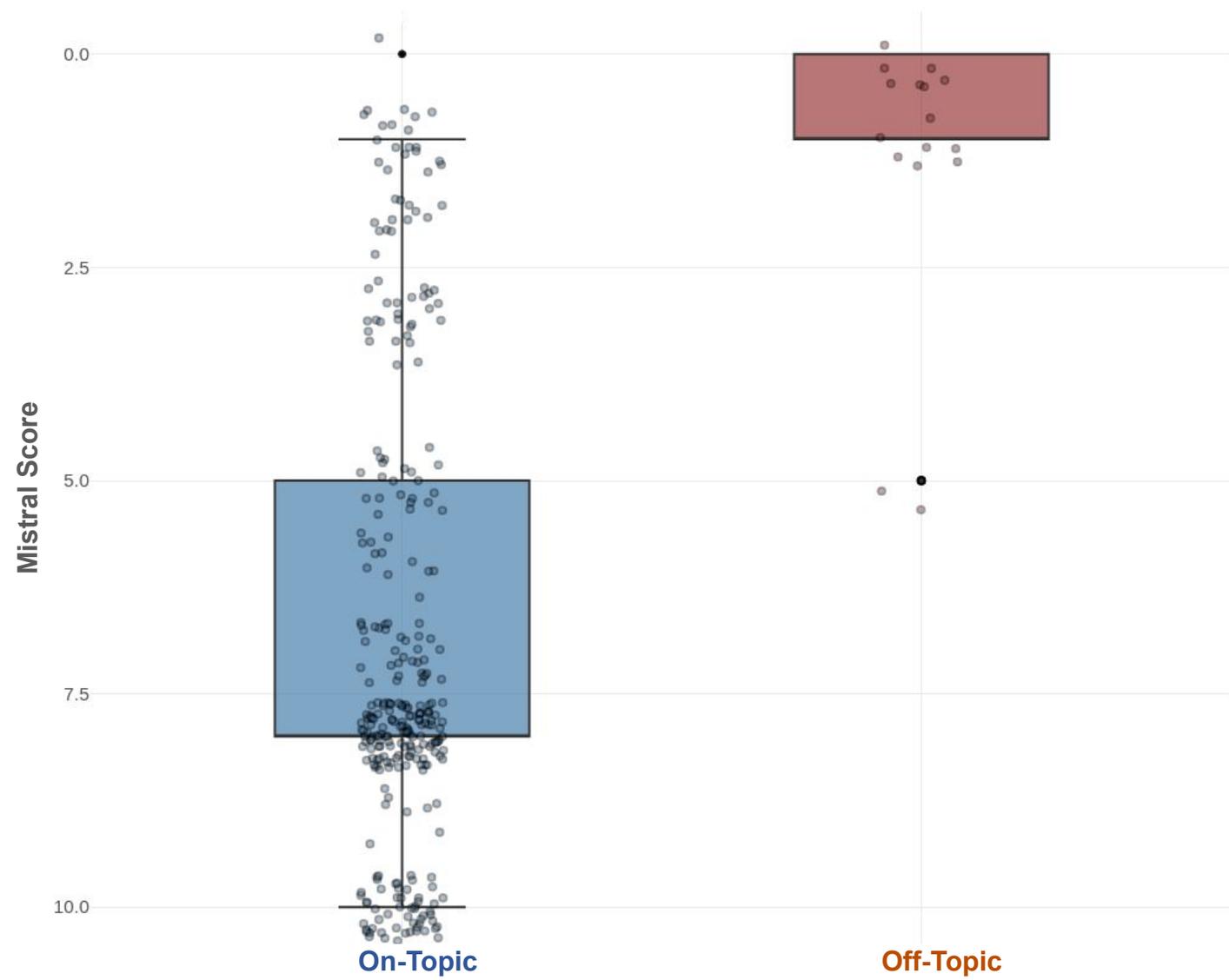
2) Sentence-BERT Semantic Similarity between Prompt and Answer



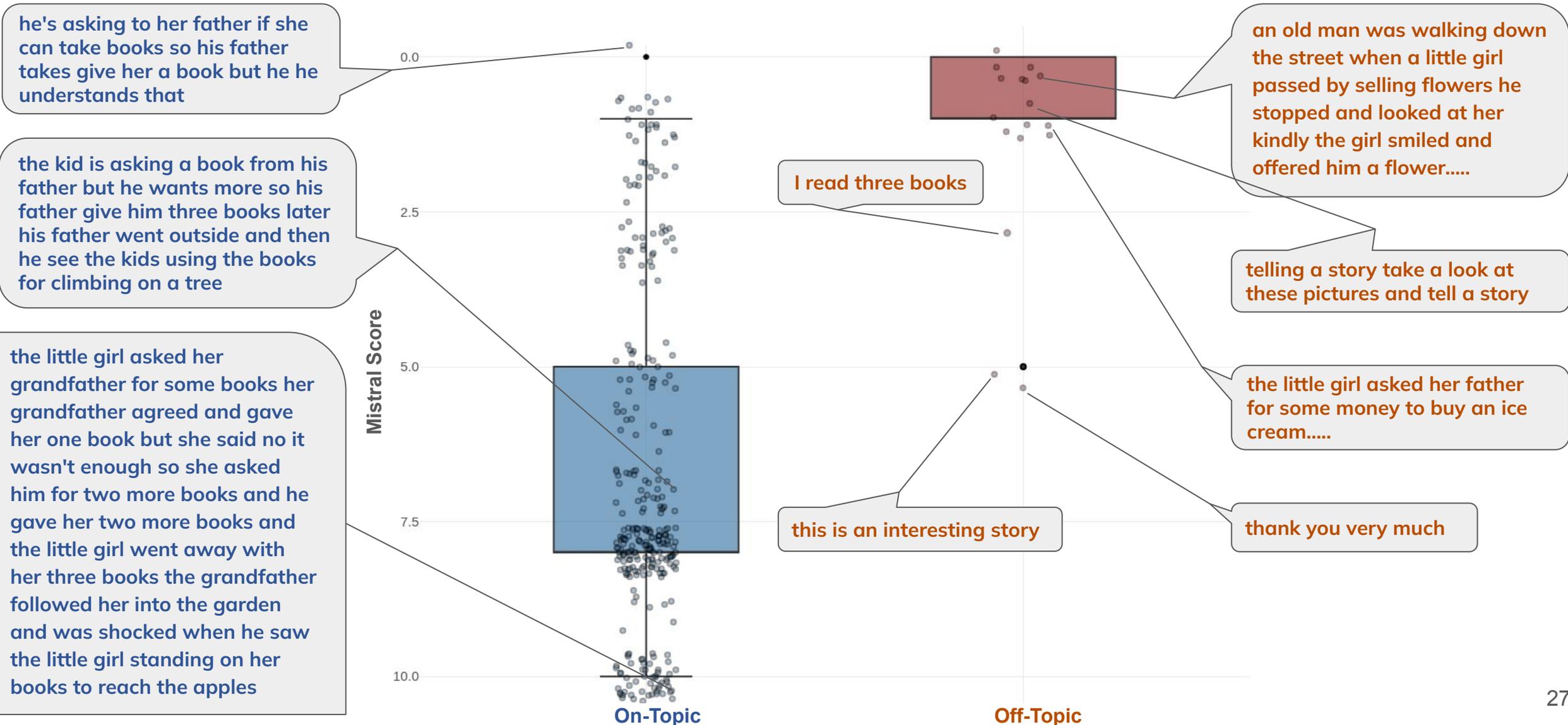
2) Sentence-BERT Semantic Similarity between Prompt and Answer



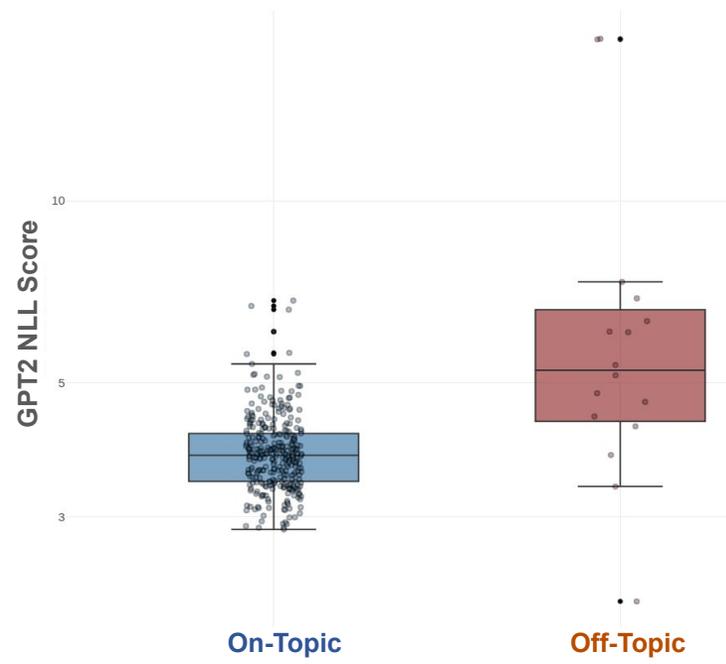
3) LLM as a Judge



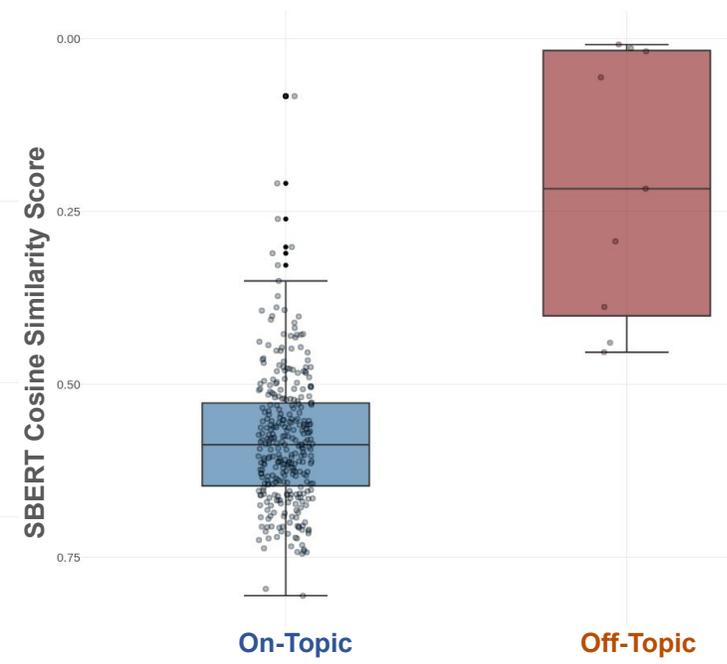
3) LLM as a Judge



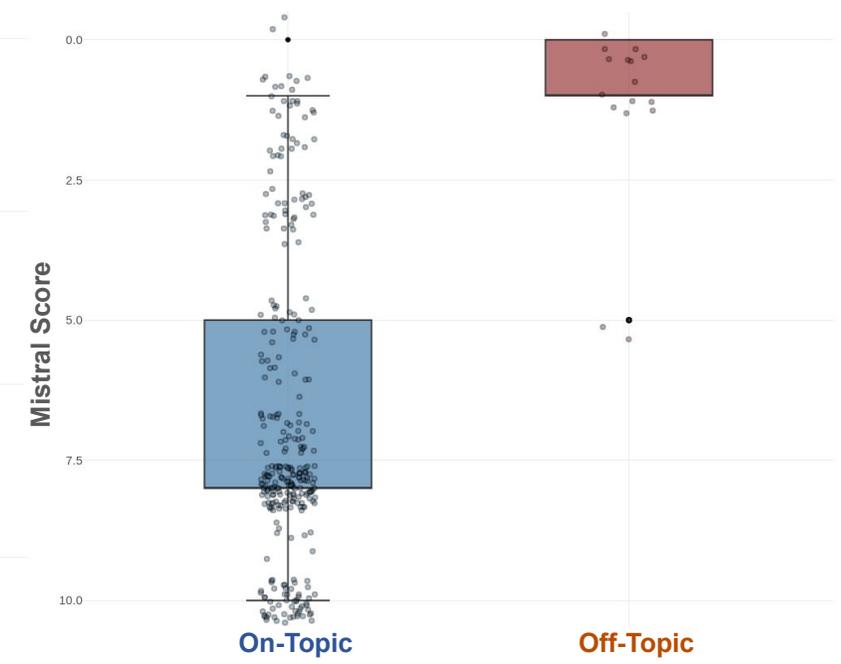
Comparison across GPT2 likelihood, Sbert similarity and LLM judgement



F1 score: 0.722

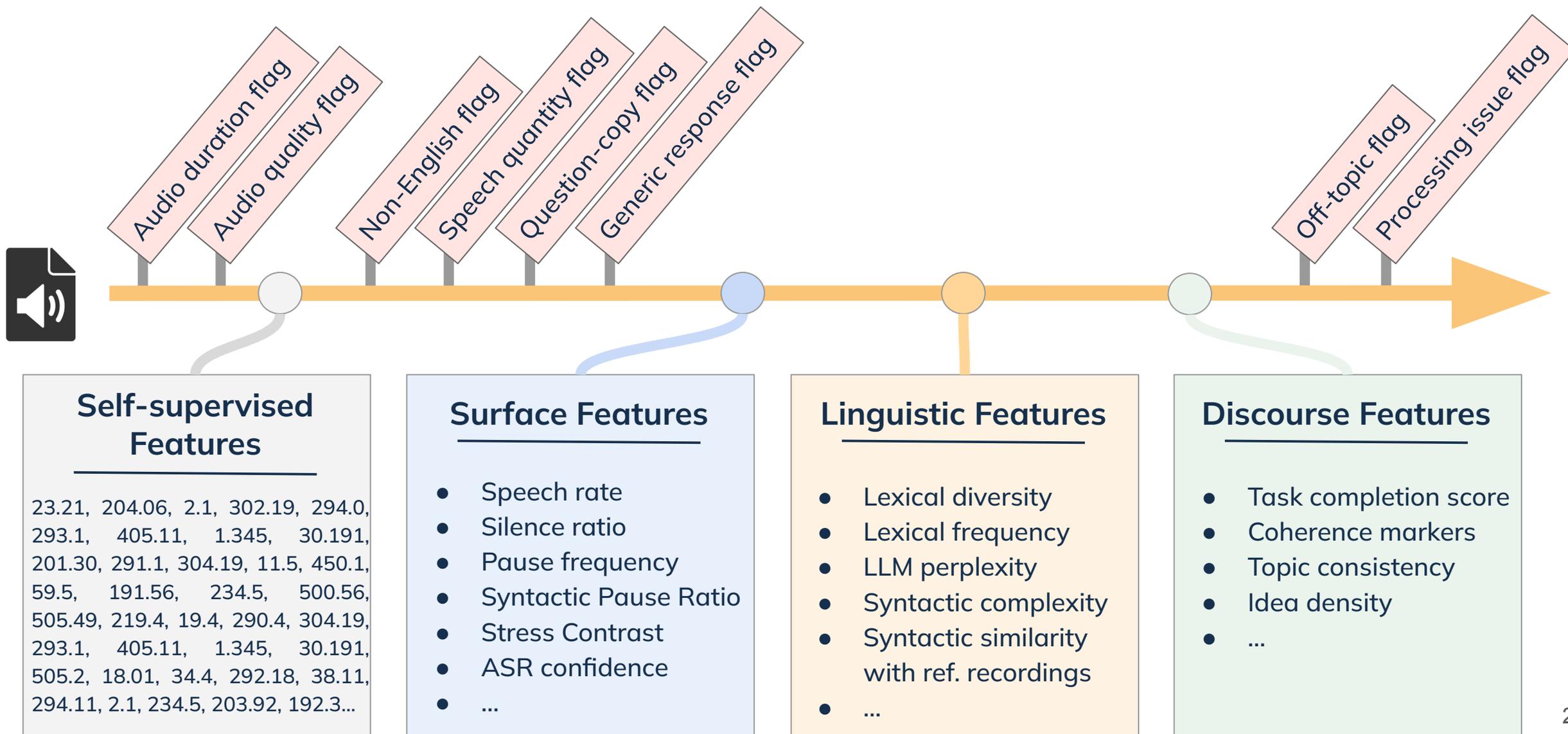


F1 score: 0.909



F1 score: 0.75

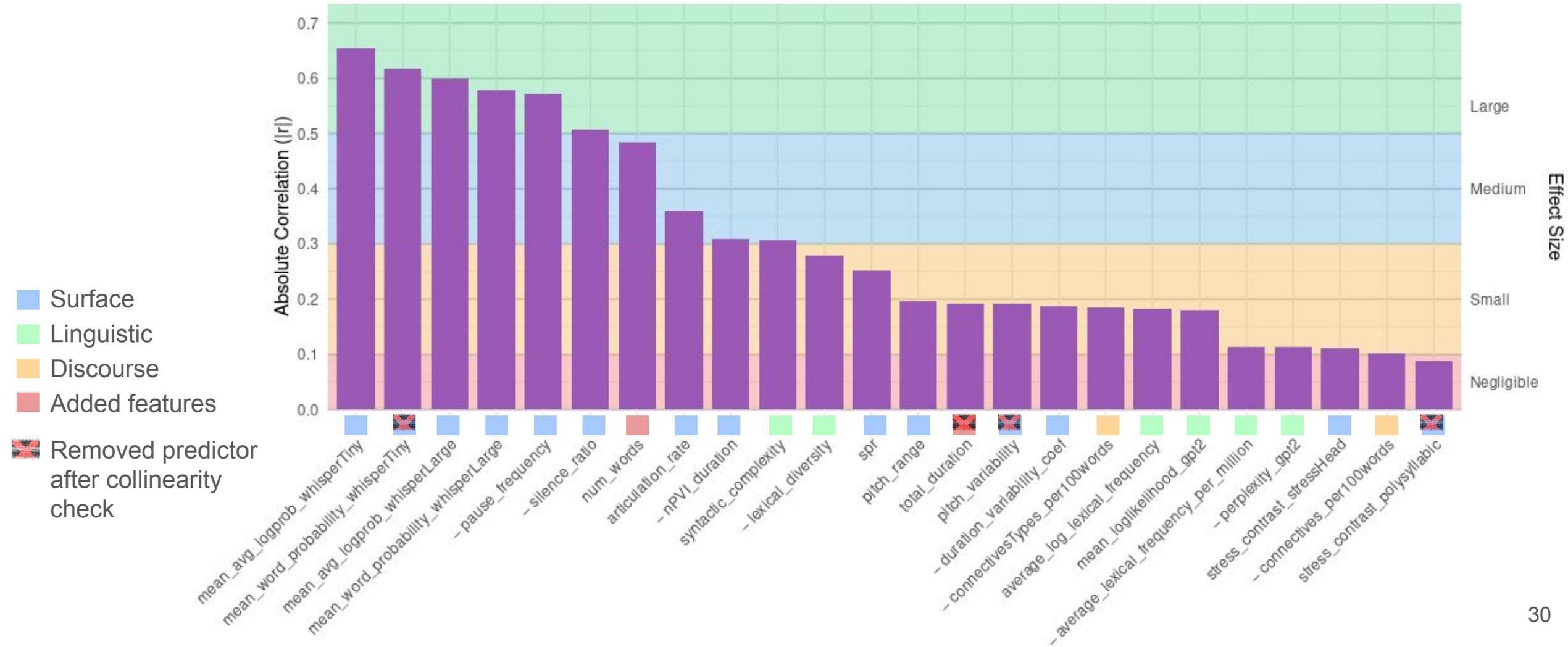
Evaluation Criteria (for English)



Correlations between Features and Human Scores

Significant Pearson Correlations with Global Level

Showing absolute values for significant correlations ($p < 0.05$) | Measures with '-' prefix have negative correlations



Correlations between Features and Human Scores

Category	Sub-category	Measure	r	Effect Size	Sig	Description
surface	Pronunciation	mean_avg_logprob_whisperTiny	0.654	Large	***	Mean segment-level log-probability score from Whisper tiny.en
surface	Pronunciation	mean_word_probability_whisperTiny	0.618	Large	***	Mean word-level confidence score from Whisper tiny.en
surface	Pronunciation	mean_avg_logprob_whisperLarge	0.599	Large	***	Mean segment-level log-probability score from Whisper large
surface	Pronunciation	mean_word_probability_whisperLarge	0.577	Large	***	Mean word-level confidence score from Whisper large
surface	Fluency	pause_frequency	-0.57	Large	***	Number of pauses / number of words
surface	Fluency	silence_ratio	-0.508	Large	***	Total silence duration / total audio duration
surface	Meta	num_words	0.485	Medium	***	Absolute number of words
surface	Fluency	articulation_rate	0.359	Medium	***	Number of words / (duration - pauses)
surface	Prosody	nPVI_duration	-0.309	Medium	***	Rhythm (Normalized Pairwise Variability Index for vowel duration)
surface	Fluency	spr	0.251	Small	***	Syntactic Pause Ratio (Coulange & de Jong, 2025)
surface	Prosody	pitch_range	0.196	Small	***	5th and 95th percentile of semitone pitch values after bandpass filtering
surface	Meta	total_duration	0.193	Small	***	Duration in seconds
surface	Prosody	pitch_variability	0.193	Small	***	Standard deviation of semitone pitch values after bandpass filtering
surface	Prosody	duration_variability_coef	-0.187	Small	***	Coefficient of variation for vowel durations (SD/Mean, Dellwo, 2006)
surface	Prosody	stress_contrast_stressHead	0.112	Small	**	S score for words with initial and medial stress only
surface	Prosody	stress_contrast_polysyllabic	0.088	Negligible	*	S score for polysyllabic words only (Coulange, 2025)
surface	Prosody	content_function_stress_contrast	-0.055	Negligible		Stress contrast between content and function words: mean of expected primary stressed syllable of content words / mean of syllables of function words
surface	Prosody	stress_contrast	-0.047	Negligible		S score for all words (Coulange, 2025)
surface	Prosody	nPVI_f0	0.019	Negligible		Normalized Pairwise Variability Index for vowel F0 (pitch pattern metric)
linguistic	Syntax	syntactic_complexity	0.308	Medium	***	Mean number of closing and opening constituents between adjacent words
linguistic	Vocabulary	lexical_diversity	-0.28	Small	***	Type-token ratio
linguistic	Vocabulary	average_log_lexical_frequency	0.182	Small	***	Lexical types average log word frequency based on the COCA corpus
linguistic	Syntax	mean_loglikelihood_gpt2	0.18	Small	***	Overall English "Likeliness" based on GPT2
linguistic	Vocabulary	average_lexical_frequency_per_million	-0.114	Small	**	Lexical types average per-million word frequency based on the COCA corpus
linguistic	Syntax	perplexity_gpt2	-0.114	Small	**	Perplexity score from GPT2
linguistic	Syntax	syntactic_similarity	(pending)			Similarity of syntactic structures between response and model answers
linguistic	Vocabulary	target_identification	(pending)			Percentage of target words identified in the response based on task-specific keywords
discourse	Coherence	connectivesTypes_per100words	-0.184	Small	***	Number of conjunction types (e.g., "however", "therefore") per 100 words
discourse	Coherence	connectives_per100words	-0.103	Small	**	Number of conjunction tokens (e.g., "however", "therefore") per 100 words
discourse	Topic	task_completion	(pending)			Percentage of task requirements met in the response
discourse	Topic	topic_consistency	(pending)			Topic consistency measure (definition pending)
discourse	Topic	idea_density	(pending)			Number of content words / total number of words

Training Simple Logistic Classifiers based on Features

A) Linear Mixed Model (lmm) Trained on continuous scores (1-7)

B) Cumulative Link Mixed Model (clmm) Trained on CEFR levels (A1-C2)

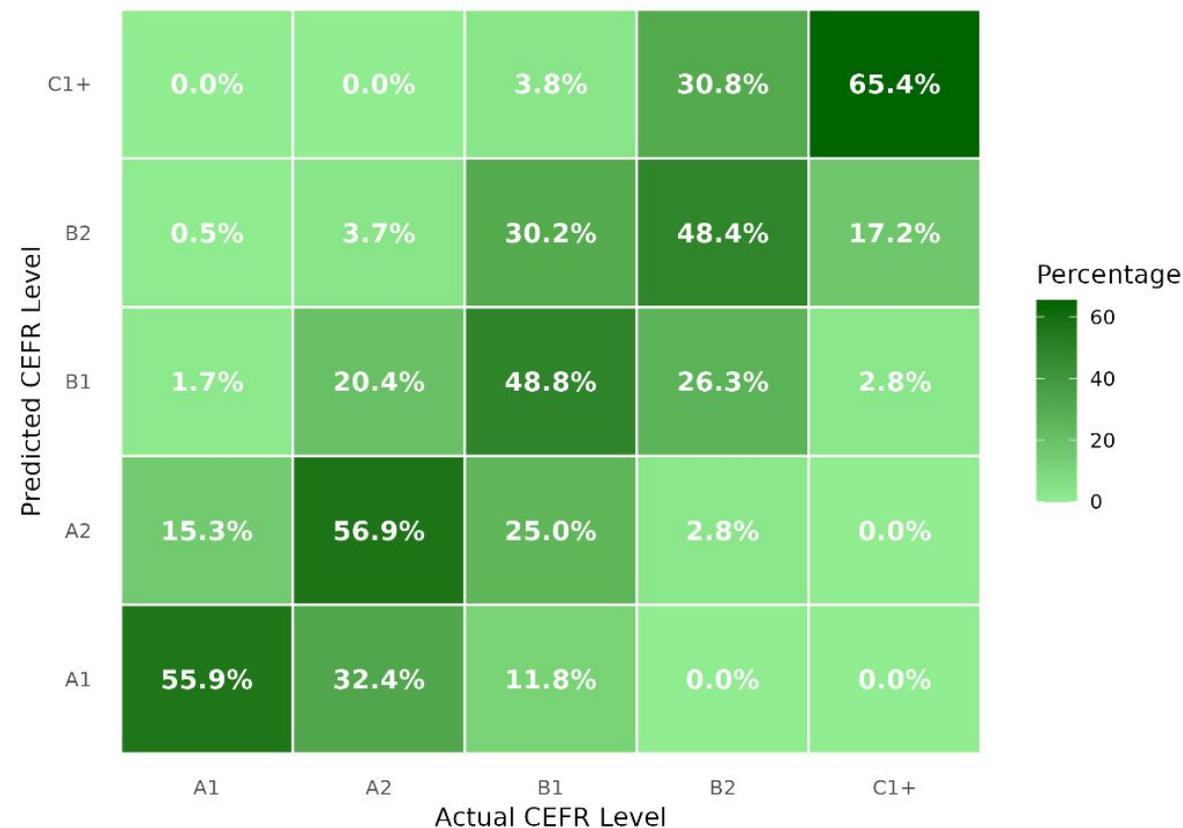
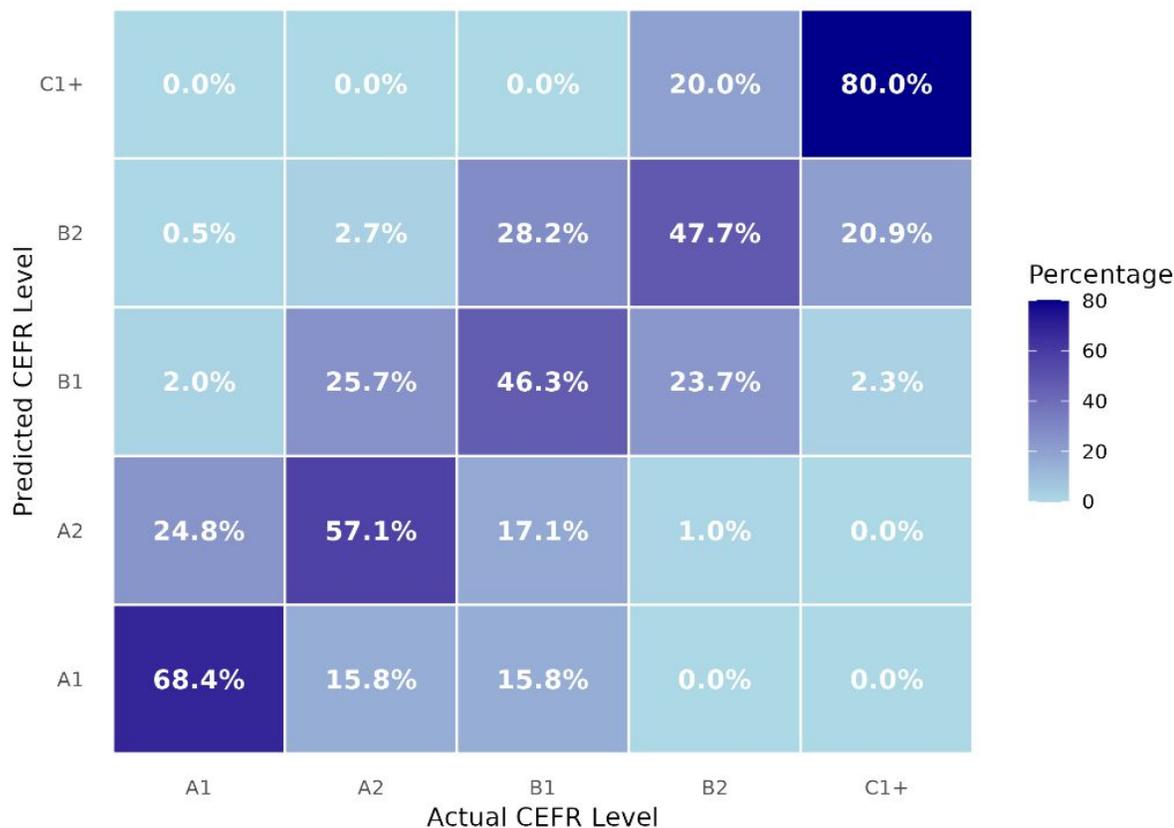
Cross-Validation Confusion Matrices (10 folds)

Linear Mixed Model

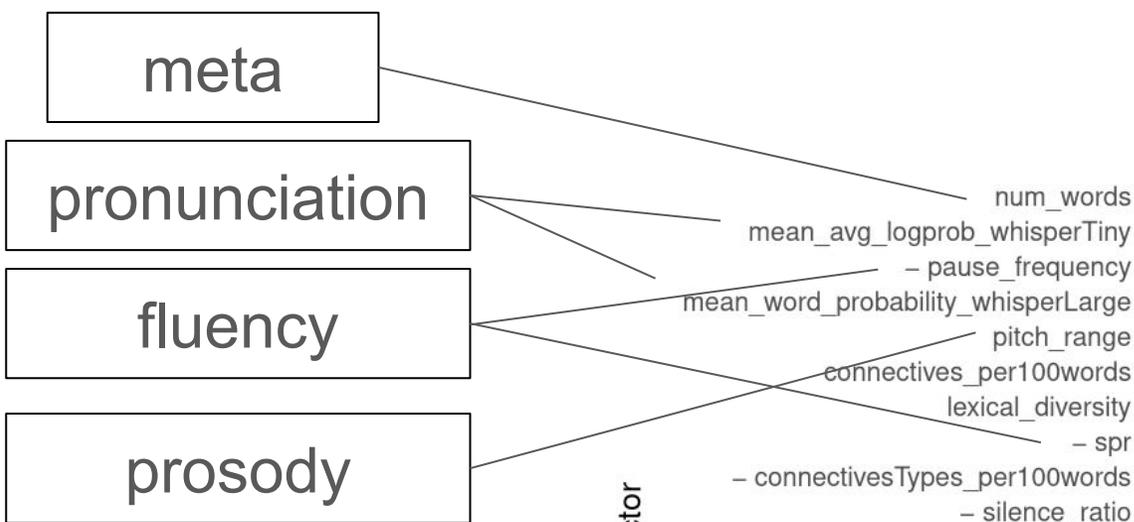
10-fold CV | Weighted F1: 0.515 (0.436-0.586)
 Exact accuracy: 49.4% (41.8-57.1) | Within ±1 level: 96.4% (93.3-100.0)

Ordinal Model

10-fold CV | Weighted F1: 0.519 (0.461-0.585)
 Exact accuracy: 51.2% (46.3-58.6) | Within ±1 level: 95.6% (90.3-100.0)



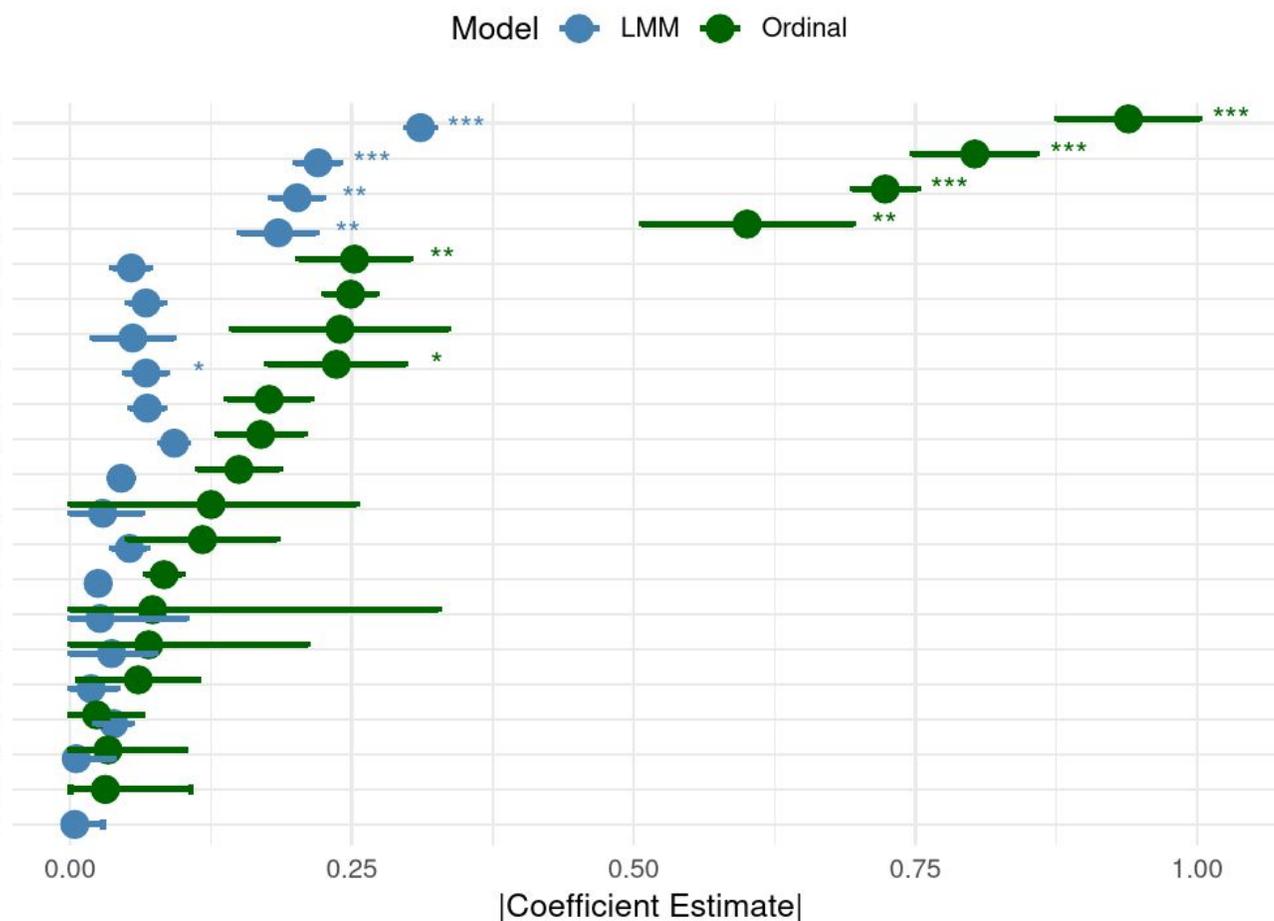
Features used by the Classifier



- num_words
- mean_avg_logprob_whisperTiny
- pause_frequency
- mean_word_probability_whisperLarge
- pitch_range
- connectives_per100words
- lexical_diversity
- spr
- connectivesTypes_per100words
- silence_ratio
- average_lexical_frequency_per_million
- mean_avg_logprob_whisperLarge
- nPVI_duration
- stress_contrast_stressHead
- perplexity_gpt2
- mean_loglikelihood_gpt2
- duration_variability_coef
- syntactic_complexity
- articulation_rate
- average_log_lexical_frequency
- average_log_lexical_frequency

Model Comparison: Absolute Coefficient Estimates

Mean \pm SD across 5 folds | '-' prefix = negative effect
 Significance: *** = all folds $p < 0.001$, ** = 80%+ folds $p < 0.01$, * = 80%+ folds $p < 0.05$



Merci !

Pinxun HUANG, Eli STAFFORD, Sylvain COULANGE

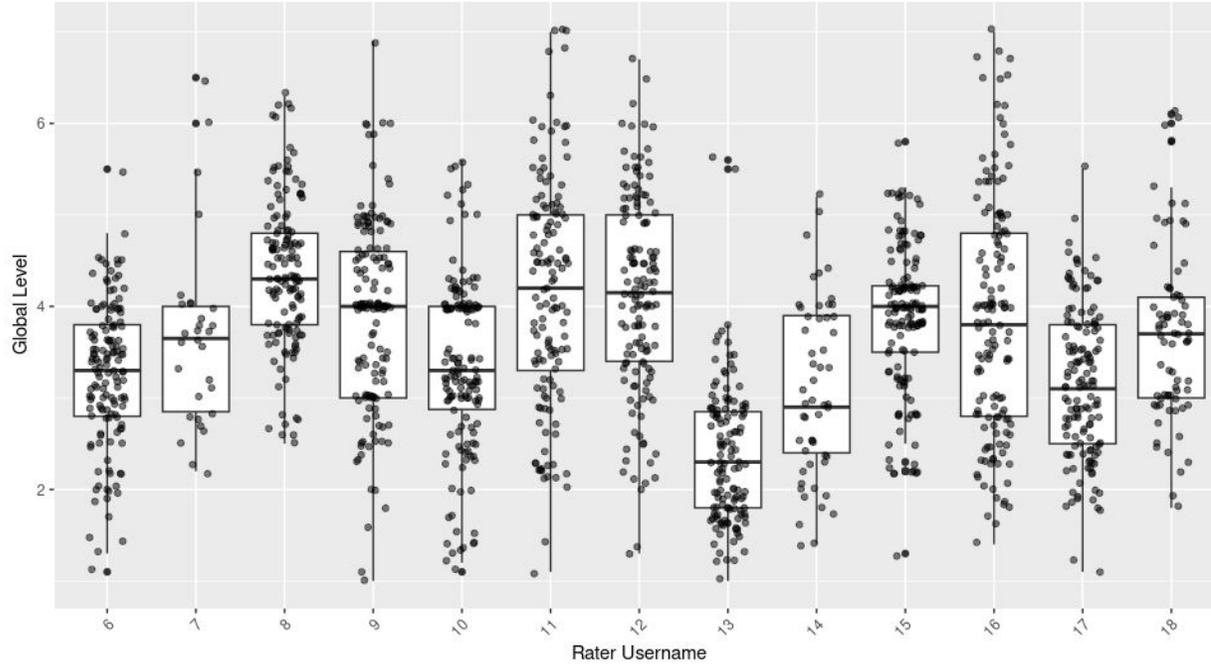
LIDILEM (*Laboratory of Linguistics and Didactics of Foreign and Mother Tongues*)

LIG (*CNRS, Institute of Engineering, Grenoble Computer Science Laboratory*)

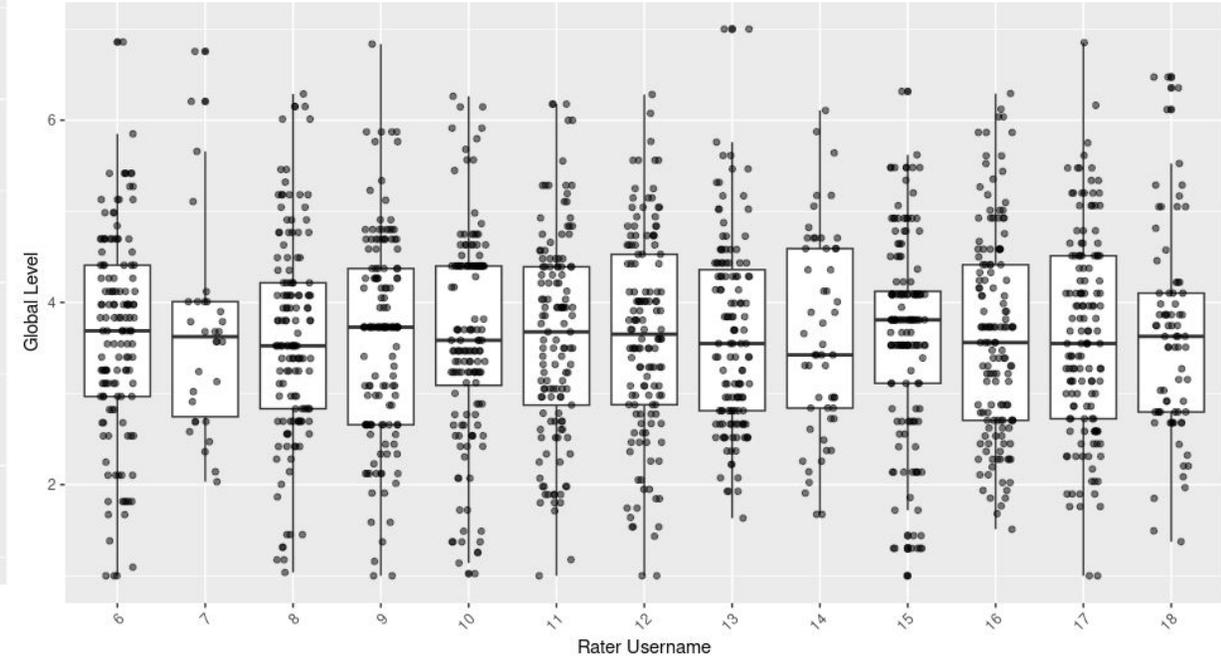
Université Grenoble Alpes

Simple CEFR prediction model

Distribution of Global Level per Rater ID (raw)

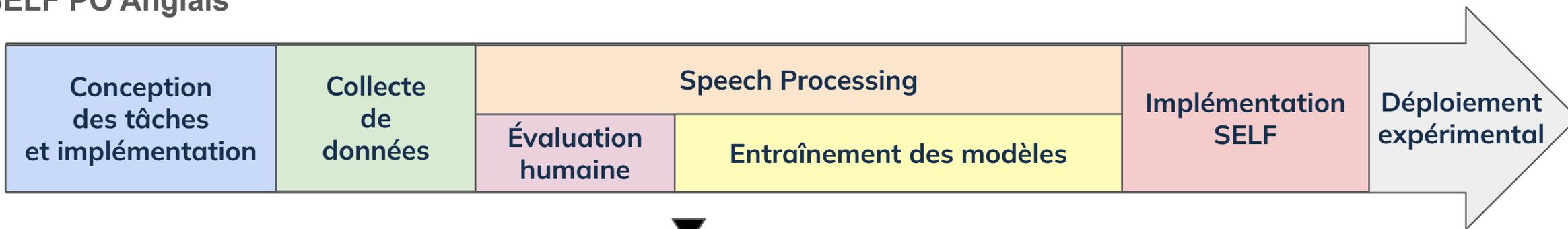


Distribution of Global Level per Rater ID (aligned to pooled mean/SD, no Rater excluded)



Rater normalization method: Linear equating
(linear alignment to pooled mean/SD)
+ clipping [1.0, 7.0]

SELF PO Anglais



mai juin juill août sept oct nov déc jan fév mars avr mai juin juill août sept oct nov déc jan fév mars
 2025 2026 2027

SELF PO Français

