

KostasThesis2025 at SemEval-2025 Task 10: A Continual Learning Approach to Propaganda Analysis in Online News

Konstantinos Eleftheriou
eleftheriou.konst@gmail.com

Panos Louridas
louridas@aueb.gr

John Pavlopoulos
annis@aueb.gr

Abstract

In response to the growing challenge of propagandistic presence through online media in online news, the increasing need for automated systems that can identify and classify narrative structures in multiple languages is evident. We present our approach to the SemEval-2025 Task 10 Subtask 2, focusing on the challenge of hierarchical multi-label, multi-class classification in multilingual news articles. Working with a two-level taxonomy of narratives and subnarratives, in the Ukraine-Russia War and Climate Change domain, we present methods to handle long articles based on how they are naturally structured in the dataset, propose a hierarchical classification neural network model with respect to the narrative taxonomy structure, and establish a continual learning training approach that leverages cross-lingual knowledge transfer. Our system was evaluated in five languages, achieving competitive results while demonstrating low variance compared to similar systems in our leaderboard position.

1 Introduction

From early days, propaganda has been a tool in shaping people’s beliefs, actions, and behaviours. The most effective propaganda techniques often act undetected, influencing readers without even their knowledge (Muller, 2018). With the rapid growth of the Internet and the Web revolutionizing the way people share and access information, it has also opened doors to propagandistic techniques being disseminated more effectively, reaching vast audiences worldwide (Tardaguila et al., 2018). The large volume of online content expanding with more than 500 million tweets per day on Twitter (or X) (Sayce, 2022) by 2022, makes manual identification of propaganda impractical in the digital era, highlighting the need for automated tools and systems.

At research level, most work on propaganda detection has focused on high-resource languages, such as English, and little effort has been made to detect propaganda for low-resource languages. Previous work examined content at the document

level (Rashkin et al., 2017), where they work focused on analyzing entire articles to differentiate between propaganda, trusted news, and satire rather than analysing specific narrative structures. SemEval-2020 Task 11, which focused on propaganda and news analysis, was introduced to address this, (Da San Martino et al., 2020) featuring the classification of portions of documents across 44 propagandistic techniques.

SemEval-2025 Task 10¹ was introduced as a significant advancement that focuses on the automatic identification of specific narrative structures, their classification, and the roles of entities involved in online articles in a multilingual setting. It offers three subtasks on news articles: Entity Framing, Narrative Classification and Narrative Extraction.

This study focuses on the Narrative Classification subtask of SemEval-2025 Task 10. Unlike previous tasks, previously discussed, it centers around the identification of both the broader narratives of articles and their specific subnarratives. We explore how hierarchical neural networks can model this nested taxonomy structure, investigate methods for handling long article inputs, and examine how different language orders can affect model performance in a continual learning training strategy. Researchers have studied whether language order affects catastrophic forgetting in continual learning, but optimal order could vary across tasks and language sets. Our research builds on their findings, attempting to address the following research question: "Is there an optimal language order in language-specific continual learning for narrative classification? If so, which is the best and which is the worst?"

During our participation in the challenge, our primary approach, consisting of an ensembled version of a continual learning training strategy was evaluated in five languages with strong results in Portuguese (3rd), English (5th), Russian (5th), Hindi (5th) and Bulgarian (5th) out of on average 18

¹<https://propaganda.math.unipd.it/semEval2025task10/index.html>

teams². Our analysis revealed (Subsection 3.2) that the order in which our model is trained matters significantly in model performance, with certain language orders outperforming others. An important aspect of our approach was stability, showing lower variance in predictions compared to similarly ranked systems for certain languages.

Narrative(s)	Subnarrative(s)
Discrediting Ukraine	Ukraine is associated with nazism, Discrediting Ukrainian military.
Praise of Russia	Praise of Russian military might.

Table 1: Example classification of narratives and subnarratives for a Ukraine-Russian War article.

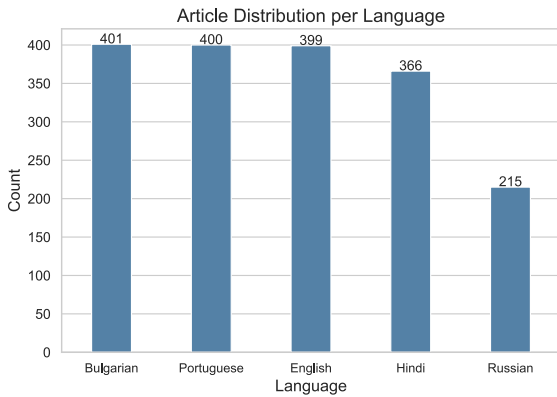


Figure 1: Distribution of articles across the five languages in the training dataset.

1.1 Related Work

1.1.1 Related Tasks

Coan et al. (2021) presents a classification task that focuses solely on contradictory claims of climate change in a similar hierarchical taxonomy. Their work emphasized on structuring claims into multiple levels of specificity following a parent-child relationship. Piskorski et al. (2022) builds upon climate change. It showcases an effective and interesting way of handling limited training data using data augmentation techniques by maintaining the meaning intact, and thus preserve also label consistency, when trying to generate synthetic data. In a different but conceptually relevant domain, Kotseva et al. (2023) developed a hierarchical classification system for COVID-19 misinformation narratives, spanning over 58,000 articles in a similar multilingual setting.

²<https://propaganda.math.unipd.it/semeval2025task10/leaderboardv2.html>

1.1.2 Continual Learning in NLP

Continual learning has gained attention in NLP for its ability to transfer knowledge across different tasks. It is a known method in which input data is continuously used to extend the existing model’s knowledge, i.e. to further train the model. The major challenge towards this goal is catastrophic forgetting, meaning that a continually trained model tends to forget the knowledge it has learned before (Wang et al. (2023); Kirkpatrick et al. (2017)).

In the NLP area, continual learning has been established in different domains. Mi et al. (2020) demonstrates this with a dialogue system that is trained on sequential data. Their approach overcame the catastrophic forgetting challenge by re-playing important past examples. In a next token prediction task, Gogoulou et al. (2024) experimented with training a model in different languages using continual learning. That is, the model was trained first on a single language, then the training would continue with a different language, and so on. They discovered that the language order in which the model is trained, plays a crucial role – a carefully selected language order also showed to reduce catastrophic forgetting.

1.2 Dataset

The dataset is composed of news articles in five languages: Bulgarian, Russian, Portuguese, English, and Hindi. These articles were collected and annotated specifically for the Ukraine-Russia War and Climate Change domains.

The data is divided into training (1.781 articles), development (178), and test (460) sets. The distribution of training data across languages is shown in Figure 1.

The articles vary significantly in length, a characteristic that introduces challenges we discuss in Section 1.3 and address in Subsection 2.1.

1.3 Challenges

The initial challenge resides in the way our labels are structured for classification. Each article can belong to one or more narratives that each map to one or more subnarratives, creating this two-level hierarchy. This presents challenges not only in hierarchical depth, since both levels must be predicted correctly, but also in managing the large number of possible labels. The scale of this can be better understood by looking at Figure 2, which shows a partial taxonomy of narratives and subnarratives for the Ukraine-Russia War domain.

Cross-lingual Variations: Working across multiple languages presents several specific challenges.



Figure 2: Sample taxonomy for Ukraine-Russia War, showing the hierarchical relationship between narratives (inner ring) and their corresponding subnarratives (outer ring).

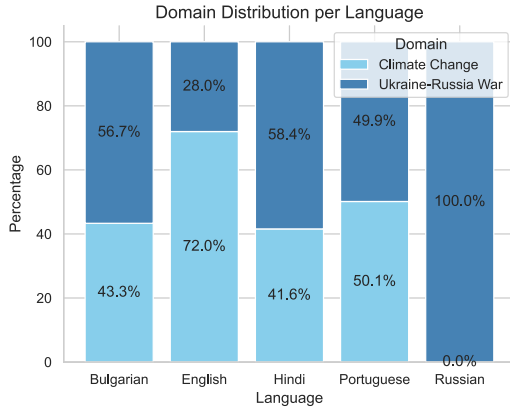


Figure 3: Distribution of domain across the five languages in the training set.

Different languages tend to favour certain narrative patterns over others due to geopolitical factors. Figure 3 demonstrates this, with the Russian language focusing solely on the Ukraine-Russia War Taxonomy, while others exhibit a more balanced distribution between domains. The combination of cross-lingual variations with limited training data poses data imbalance issues, where certain narrative-subnarrative pairs appear much more frequently than others, something we address in Subsection 2.3.

The combination of cross-lingual variations with limited training data poses data imbalance issues, where certain narrative-subnarrative pairs appear much more frequently than others, something we address in Subsection 2.3.

Article Length Variability: Articles vary significantly in length, ranging from short to extensive (mean 403 words, std dev 237 words; between

88 to 924 words across languages). Most (best) text classification models are specifically trained (or fine-tuned) to give good sentence embeddings; however, these models typically have a maximum token limit (usually 512 or 1024 tokens), which becomes problematic when processing large articles into representations that our classification models can then understand. We carefully handle longer news articles in Section 2.1 to overcome a situation where article representation adversely affects the classification task.

1.4 Evaluation Metrics

Our models are evaluated on two F1-score metrics:

- **Fine-grained F1:** Averaged samples F1 computed for complete narrative-subnarrative pair labels, where both the narrative and subnarrative parts must be correct for the predicted label to be considered correct.
- **Coarse-grained F1:** Averaged samples F1 computed for narratives only, by ignoring the subnarrative parts of the narrative:subnarrative predicted and gold labels.

Both metrics are averaged across test articles.

2 System Overview

2.1 Article Representation

We propose a chunking approach that follows the natural structure of news articles in the dataset.

Chunking Approach: Rather than employing an arbitrary paragraph-based splitting, we leveraged the inherent structure of articles, that is, their header/body/footer structure (Table 2).

Section	Characteristics
Header	Short to medium length, introduces main topic
Primary Body	Extended content, core narrative
Secondary Body	Optional extended content
Tertiary Body	Optional medium to long content
Footer	Optional closing content, short to medium

Table 2: Structural composition of news articles in the dataset.

Combining the separated sections into a single embedding that describes the whole article is also something we ought to take care of.

We explored various strategies for doing so:

- Average pooling between sections: Average of all section embeddings, preserving each section equally.

- Weighted average based on section length: Similar to averaging, but sections contribute proportionally to their length.
- Sum of section embeddings: Element-wise addition of all section embeddings, essentially preserving all information.

We analyze the impact of these strategies in Paragraph 3.1 and Subsection 3.2.

2.2 Model Architecture

The problem itself is structured in such a way that it differs from a two-head classification model, where we have a head for classifying narratives and a separate for subnarratives. Each narrative has its own set of subnarratives creating this natural hierarchy.

We developed a base multi-head, multi-task model approach where we have a single head for predicting narratives, then multiple heads for predicting the subnarratives for the given narrative hierarchy. We then explored several variants of this model as for experiments.

We then explored several variants of this model as for experiments.

2.2.1 Multi-Head Base Architecture

Our base architecture consists of three main components:

- A shared base layer that learns features and provides its output to the lower layers.
- A narrative head for predicting the top-level narratives.
- Multiple heads, one per narrative hierarchy, each predicting the corresponding subnarratives for that hierarchy.

2.2.2 Hierarchical Variants

Concatenation Model: Our base model treated narrative and subnarrative predictions independently. That is, subnarrative predictions were computed as:

$$P(subnarr_j|x) = \sigma(h(x)) \quad (1)$$

where $h(x)$ the output of the shared layer (the gray box) for article embedding x .

We enhanced this by concatenating the narrative probabilities with the shared layer output:

$$P(subnarr_j|x) = \sigma([h(x); P(narr_i|x)]) \quad (2)$$

where $narr_i$ is the parent narrative of $subnarr_j$.

This is intuitive, because:

Multi-Head (Base) Model Architecture

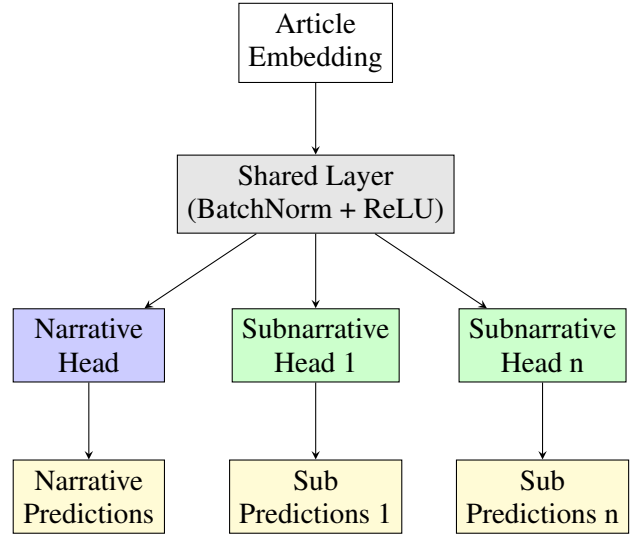


Figure 4: Architecture of the base multi-head model showing the flow from article embedding through shared layer to narrative and subnarrative heads.

- If the probability of the narrative is high, the subnarrative head will be more likely to predict the relevant subnarratives.
- If the probability is low, the model will learn to ignore the corresponding subnarratives.
- At the same time, the shared output of the shared layer will help determine which subnarrative is most appropriate for the given article.

Multiplication Model: As an alternative to concatenation, we implemented element-wise multiplication between the output of the shared layer and the narrative probabilities.

$$P(subnarr_j|x) = \sigma(h(x) \odot P(narr_i|x)) \quad (3)$$

where $h(x)$ is the shared layer output for article embedding x .

This conceptually creates a stronger hierarchical dependency, acting as a natural "gate" in the hierarchy:

- If the narrative probability is close to 0, the corresponding subnarrative head's input will be scaled down, effectively disabling that subnarrative head.
- If the narrative probability is close to 1, the shared layer output passes through somewhat unaffected.

Concatenation Model Architecture

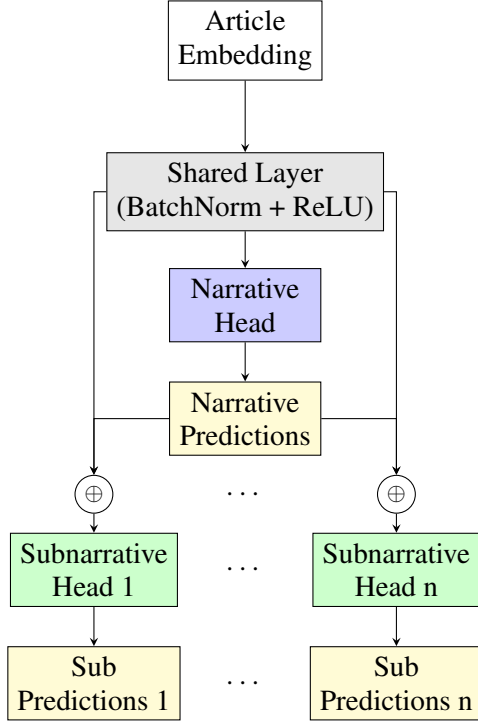


Figure 5: Architecture overview of the architecture for the concatenation model, showing how narrative predictions are combined with shared layer output to feed into subnarrative heads.

Multiplication Model Architecture

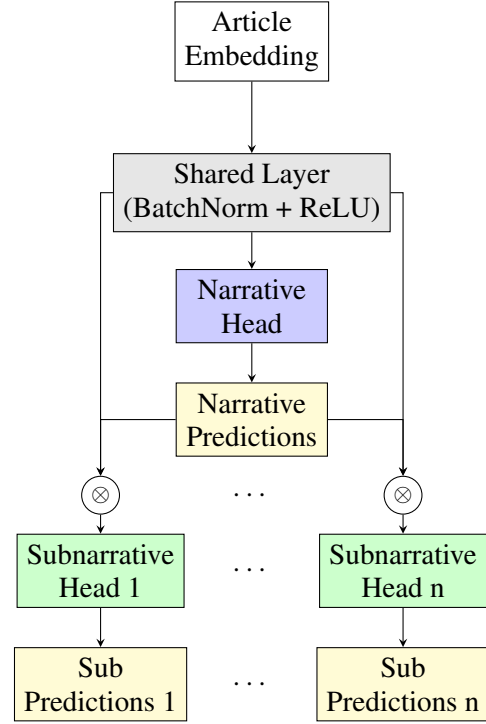


Figure 6: A simplified architecture overview of the multiplication model, showing how narrative predictions act as gates by multiplying with shared layer output before feeding into subnarrative heads.

2.2.3 Heatmap Model

We also developed a simplified approach that views the problem as a 2D classification task. Instead of managing multiple prediction heads, this model directly outputs a probability matrix $H \in \mathbb{R}^{n \times c}$ where:

- Rows (n) represent narratives.
- Columns (c) represent subnarratives.
- Each cell $H_{i,j}$ represents the probability of the narrative-subnarrative pair.

2.3 Loss Function

Our loss function is designed to handle both imbalanced labels and the need to remain consistent in our hierarchical predictions.

Weighted BCE: We use a weighted version of BCE (Binary Cross Entropy) to account for the class imbalance. Each label is assigned a weight that is proportional to its frequency in the dataset. This way, rare labels contribute proportionally more to the loss.

Hierarchy and Miss-classifications: We penalize inconsistencies in the hierarchy and label miss-classifications. More specifically, the loss consists of:

$$\mathcal{L}_{\text{total}} = (1 - W_{\text{sub}}) \cdot \mathcal{L}_{\text{narr}} + W_{\text{sub}} \cdot \mathcal{L}_{\text{sub}} + W_{\text{cond}} \cdot \mathcal{L}_{\text{cond}} \quad (4)$$

$\mathcal{L}_{\text{narr}}$ represents the weighted BCE loss for narrative predictions, while \mathcal{L}_{sub} captures the weighted BCE loss for subnarrative predictions. The term $\mathcal{L}_{\text{cond}}$ serves as a conditioning term that enforces hierarchical relationships.

The conditioning term enforces the hierarchical structure through:

$$\mathcal{L}_{\text{cond}} = \text{mean}(|p_{\text{sub}} \cdot (1 - p_{\text{narr}})| + p_{\text{narr}} \cdot |p_{\text{sub}} - y_{\text{sub}}|) \quad (5)$$

The first part ($|p_{\text{sub}} \cdot (1 - p_{\text{narr}})|$) penalizes the model for predicting subnarratives when their parent narrative is inactive. The remaining part ensures that the subnarrative predictions match ground truth when their parent narrative is active.

2.4 Training Strategies

Our initial experiments with the base architectures revealed significant performance instability across training runs (Section 3). This motivated us to explore the idea of ensembling, which has been

shown to reduce variance and improve generalization in classification tasks (Dietterich, 2000).

In addition, we establish several training strategies, each one tackling the problem from a different angle.

2.4.1 N-fold Cross-Ensembling

We employed an n-fold cross ensembling approach that leverages our data more effectively. This training strategy splits the combined train and validation data into n folds and then trains n distinct models. Each model learns from a different subset of the data, thus capturing slightly different patterns. During inference, we average the predictions from all models to produce the final output.

2.4.2 Continual Learning

The instability problem motivated us to try an alternative approach, one that changes the way the model learns from the training data. For the training phase of our models, we tried a more sequential approach, matching closely with how we, humans, learn different concepts.

Just as learning Ukrainian becomes easier when you know Russian (by having similar grammar and vocabulary), we hypothesized that this sequential order can help our model find meaningful patterns per language.

In particular, for our problem:

- Russian language can provide a good base for the URW taxonomy.
- Bulgarian builds on top of Russian as both are Slavic languages.
- Every single language that follows keeps enriching the model’s understanding with its unique characteristics.

Upon reaching our target language during the training phase, we give the model more time to adapt by increasing its training patience and lowering the learning rate.

We created an ensemble combining multiple models trained upon different order. During inference, better performing language orders get more weight in the final prediction.

3 Evaluation

Below we present the comparison results across model variants, embedding models, and aggregation strategies. We report both Coarse-F1 (for narratives) and Fine-F1 (for subnarratives), along with their standard deviations. However, the primary focus of the task is on the Fine-F1 score.

All comparisons are performed specifically for the English validation dataset, as it demonstrated the most balanced distribution of narratives in the dataset across the two domains and is widely recognised as the most prominent language in NLP research. Each model was run five times, and the results were aggregated to ensure a fair comparison.

We evaluated our experiments with two embedding models: KaLM³ and Stella⁴. We specifically chose these embedding models because they are both multilingual, instruction-based that achieved high performance on the MTEB (Massive Text Embedding Benchmark) leaderboard⁵. During the stage of transforming our article sections into meaningful numbers that our classification models can understand, we instructed the embedding models to:

“Produce an embedding useful for detecting relevant war- or climate-related narratives from a taxonomy.”

3.1 Architecture Variants Comparisons

Metric	Simple	Concat	Mult
Coarse-F1	0.489 ± 0.03	0.497 ± 0.02	0.477 ± 0.02
Coarse std	0.385 ± 0.01	0.386 ± 0.01	0.384 ± 0.01
Fine-F1	0.329 ± 0.03	0.333 ± 0.02	0.311 ± 0.02
Fine std	0.320 ± 0.02	0.327 ± 0.02	0.321 ± 0.01

Table 3: Mean performance comparison between the base hierarchical model and its variants (averaged over 5 runs).

Table 3 shows the mean performance across model base variants. The high standard deviation (± 0.02 - 0.03) indicates run-to-run instability.

Concat variant shows a sign of effectiveness in comparison to the Simple model by slightly outperforming it. Multiplication variant lags behind for both approaches, indicating that the hard-gating mechanism might be too restrictive. If our narrative predictions are not confident or even, and most importantly, not correct, the subnarrative head will receive very weak input because of the hard gating.

Embedding Model Comparisons: Table 4 shows performance between embedding models.

KaLM embeddings consistently appear to outperform Stella in all metrics.

³<https://huggingface.co/HIT-TMG/KaLM-embedding-multilingual-mini-instruct-v1.5>

⁴https://huggingface.co/NovaSearch/stella_en_1.5B_v5

⁵<https://huggingface.co/spaces/mteb/leaderboard>

Metric	KaLM	Stella
Coarse-F1	0.497 \pm 0.02	0.450 \pm 0.02
Fine-F1	0.333 \pm 0.02	0.298 \pm 0.02

Table 4: Performance comparison across embedding models.

Aggregation Strategy Comparisons: Tables 5 and 6 present Fine-F1 scores (our primary goal is to improve subnarrative classification, we limit this analysis to solely Fine-F1 scores for simplicity) across model variants and aggregation strategies per embedding model.

Model	Sum	Mean	Weighted
Simple	0.329 \pm 0.03	0.285 \pm 0.01	0.325 \pm 0.02
Concat	0.333 \pm 0.02	0.305 \pm 0.01	0.300 \pm 0.02
Mult	0.311 \pm 0.02	0.287 \pm 0.02	0.283 \pm 0.01
CK-Ens	0.338	0.335	0.312

Table 5: Fine-F1 scores for KaLM embeddings across model variants and aggregation strategies.

Model	Sum	Mean	Weighted
Simple	0.309 \pm 0.01	0.259 \pm 0.01	0.343 \pm 0.01
Concat	0.298 \pm 0.02	0.256 \pm 0.02	0.338 \pm 0.02
Mult	0.260 \pm 0.01	0.260 \pm 0.01	0.327 \pm 0.01
C-Ens	0.300	0.308	0.322

Table 6: Fine-F1 scores for Stella embeddings across model variants and aggregation strategies.

Sum aggregation strategy appears to perform best across all other strategies for the KaLM Embeddings. This shows that KaLM benefits from preserving all information.

On the other hand, the weighted strategy seems to suit well with Stella, consistently outperforming all other strategies.

Threshold Optimizations: Our previous experiments tried to find the most optimal thresholds separately for narratives and subnarratives, exploring values up to 0.6. These thresholds determine the minimum probability for a narrative or subnarrative to be considered active in the predictions. We discovered that the weighted aggregation strategy benefits from increasing the threshold range up to 0.9, with the most noticeable improvement for Stella Embeddings. Table 7 presents these results.

Other strategies with also different embedding models seem to receive relatively higher variance.

3.2 Continual Learning

Table 8 shows the results between several language sequences and embedding combination strategies using the Concat variant.

Model	C-F1	F-F1	F-std
Simple	0.538 \pm 0.021	0.426 \pm 0.010	0.375 \pm 0.008
Concat	0.554 \pm 0.025	0.442 \pm 0.019	0.375 \pm 0.016
Mult	0.556 \pm 0.014	0.426 \pm 0.017	0.362 \pm 0.011
CK-Ens	0.566	0.410	0.343

Table 7: Performance metrics for Stella embeddings with weighted aggregation and a threshold range up to 0.9.

C-F1: Coarse-F1, F-F1: Fine-F1, F-std: Fine-std

Order	Sum	Avg	W. Avg
RU \rightarrow BG \rightarrow PT \rightarrow HI \rightarrow EN	0.378	0.351	0.316
RU \rightarrow BG \rightarrow HI \rightarrow PT \rightarrow EN	0.356	0.323	0.341
BG \rightarrow RU \rightarrow PT \rightarrow HI \rightarrow EN	0.314	0.343	0.316
HI \rightarrow PT \rightarrow RU \rightarrow BG \rightarrow EN	0.302	0.312	0.330
PT \rightarrow HI \rightarrow RU \rightarrow BG \rightarrow EN	0.300	0.289	0.352
Ensemble of All Orders	0.348	0.349	0.357

Table 8: Impact of language ordering on Fine-F1 scores across different embedding combination strategies using Stella embeddings and a threshold range up to 0.6.

Impact of Aggregation Strategy: At first glance, we see that the combination strategy is sensitive to the language order:

- Sum strategy shows drastic response to the language ordering, with Fine-F1 scores ranging from 0.300 to 0.378.
- Mean strategy shows similar-to-moderate sensitivity, with Fine-F1 scores ranging from 0.289 to 0.351.
- Weighted average demonstrates the most balanced performance across orders, with Fine-F1 scores ranging from 0.316 to 0.357.

Specifically the weighted average strategy performs consistently better across different orders. In contrast to other strategies, it focuses on certain sections which might help the classification task, making thus the order less significant. However, when evaluating the effectiveness of a language order, we should primarily focus on the Sum and Avg strategies (which do not introduce any weighting). Both of these strategies agree that the first order (RU \rightarrow BG \rightarrow PT \rightarrow HI \rightarrow EN) produces the best results.

Impact of Language Order: When evaluating for English data, the sequence that starts with Russian followed by Bulgarian outperforms every other sequence. Even swapping between these languages shows a performance drop. This suggests that when training the model with sequential data, starting with certain languages helps it build strong foundation patterns, strongly influencing final performance.

Impact of Embedding Choice: Interestingly, while KaLM embeddings outperformed Stella in our stand-alone experiments (Section 4), we observed different behavior in continual learning, with KaLM model under performing. This might suggest that Stella embeddings might be more appropriate in a knowledge transfer setup.

Threshold Optimization for Continual Learning: Following our discovery that weighted aggregation strategy benefits from higher thresholds we applied the same approach to our continual learning training method.

Table 9 presents these results.

Language Order (Thresh)	C-F1	F-F1	F-std
RU→BG→PT→HI→EN (0.75/0.50)	0.614	0.449	0.349
RU→BG→HI→PT→EN (0.75/0.55)	0.608	0.437	0.352
RU→HI→PT→BG→EN (0.80/0.60)	0.600	0.444	0.359
BG→RU→PT→HI→EN (0.70/0.55)	0.575	0.404	0.364
PT→HI→RU→BG→EN (0.75/0.60)	0.586	0.424	0.359
HI→PT→RU→BG→EN (0.70/0.50)	0.561	0.376	0.371
Ensemble (0.75/0.60)	0.570	0.424	0.362

Table 9: Performance of continual learning models, using Stella embeddings with weighted aggregation strategy and a threshold range up to 0.9

C-F1: Coarse-F1, F-F1: Fine-F1, F-std: Fine-std

Higher thresholds (0.75/0.55 or 0.80/0.60) lead to better performance, with the RU→BG→PT→HI→EN sequence at 0.75/0.50 yielding the highest Fine-F1 score of 0.449.

In Section 5 we do an in-depth analysis to show order significance.

3.3 N-fold Cross-Ensembling

Table 10 shows the coarse and fine F1 scores across the Simple and Concat variants for each embedding model using the sum aggregation strategy approach. Table 11 presents the same results but with the weighted aggregation strategy approach.

Model	Threshold	C-F1	F-F1	F-std
Stella-Simple	0.60/0.55	0.506	0.385	0.378
Stella-Concat	0.60/0.55	0.499	0.390	0.371
KaLM-Simple	0.60/0.55	0.457	0.341	0.344
KaLM-Concat	0.50/0.45	0.503	0.373	0.346

Table 10: Cross-validation ensemble performance with sum aggregation strategy with a threshold range up to 0.9.

C-F1: Coarse-F1, F-F1: Fine-F1, F-std: Fine-std

The weighted aggregation strategy shows strong performance compared to the sum approach across all model variants and embedding models. The results should be taken into account, knowing that the

Model	Threshold	C-F1	F-F1	F-std
Stella-Simple	0.60/0.55	0.555	0.406	0.353
Stella-Concat	0.65/0.60	0.560	0.419	0.362
KaLM-Simple	0.55/0.30	0.518	0.354	0.338
KaLM-Concat	0.60/0.35	0.529	0.382	0.355

Table 11: Cross-validation ensemble performance with weighted aggregation strategy with a threshold range up to 0.9

C-F1: Coarse-F1, F-F1: Fine-F1, F-std: Fine-std

approach utilizes all available data. This approach also reduces overfitting to language-specific patterns by exposing each fold or model to different language distributions.

3.4 Heatmap Classifier

Table 12 presents the performance of our Heatmap classifier approach which is not multi-head based but follows a rather simpler architecture.

Configuration	Coarse-F1	Fine-F1	Fine std
KaLM-Sum	0.474 ± 0.016	0.340 ± 0.014	0.319
Stella-Sum	0.482 ± 0.011	0.329 ± 0.009	0.323
KaLM-Weighted	0.442 ± 0.017	0.299 ± 0.009	0.332
Stella-Weighted	0.450 ± 0.013	0.293 ± 0.018	0.335

Table 12: Heatmap classifier performance with different embedding models and aggregation strategies with 0.6 thresholds.

Like our base model and its variants, the Heatmap classifier shows similar instability issues. Increasing threshold values up to 0.9 further increased variance substantially.

4 Discussion

4.1 Test Set Performance

For our final submission, we used an ensembled version of the continual learning model using the concat variant of the base model. We positioned each target language, as the final stage of the learning sequence, which we give more patience and a lower learning rate.

The training configuration used Stella embeddings with a searching threshold of up to 0.6 and a sum aggregation strategy for section embeddings.

The results for the test set are shown in Table 13.

Lang	Rank	C-F1	std-C	F-F1	std-F
EN	16/30	0.409	0.314	0.239	0.243
PT	4/14	0.478	0.201	0.309	0.153
RU	6/15	0.596	0.257	0.333	0.234
BG	7/13	0.510	0.322	0.333	0.300
HI	6/14	0.384	0.418	0.282	0.402

Table 13: Test set performance across languages.

C-F1: Coarse-F1, F-F1: Fine-F1, std-C/F: Standard deviation for coarse/fine metrics.

An important aspect of our results is stability. The proportion of F1 score and std is lower in comparison to teams near our entry. This shows a sign that our model is able to generalize and learn robust features. In comparison however to top teams, it’s architecture is not enough to capture more complex ones.

Previous experiments showcased that the weighted avg strategy with increased thresholds yielded better performance, at least in the validation set; however, this discovery occurred after our test submission deadline. Post-competition analysis revealed that it would have yielded slightly better results (Tables 14 and 15).

Language	F1 samples	F1 std samples
EN	0.287	0.296
PT	0.329	0.171
HI	0.340	0.434
BG	0.355	0.311
RU	0.398	0.292

Table 14: Post submission comparison of test set performance using threshold 0.6 with weighted strategy and Stella Embeddings.

Language	F1 samples	F1 std samples
EN	0.362	0.370
PT	0.326	0.208
HI	0.341	0.450
BG	0.357	0.349
RU	0.400	0.283

Table 15: Post submission comparison of test set performance using threshold 0.9 with weighted strategy and Stella Embeddings.

The results show improvements in all languages when using the weighted strategy. The increased range of threshold values up to 0.9 proved significant for the English dataset. However, for the rest of the languages, having an increased threshold did not seem to contribute to better performance, with some languages even experiencing higher variance.

Our updated positions for the test set are shown in Table 16.

Lang	Rank	C-F1	std-C	F-F1	std-F
EN	5/27	0.556	0.396	0.362	0.370
PT	3/14	0.539	0.214	0.329	0.171
RU	5/15	0.571	0.344	0.400	0.283
BG	5/13	0.523	0.371	0.357	0.349
HI	5/14	0.453	0.441	0.341	0.456

Table 16: Updated test set performance across languages. C-F1: Coarse-F1, F-F1: Fine-F1, std-C/F: Standard deviation for coarse/fine metrics.

4.2 Future Work

Embeddings: Our section-based approach opens opportunities for instruction-tuned embedding models. Different sections also contain different contexts, with headers as a more generic content for the article, and body adding supportive details - a future study can experiment with a more targeted processing per section with instruction-based embedding models where different instructions can be applied to different sections.

Architecture: While our approach used powerful pre-trained embeddings, a clear limitation is that we did not perform any fine-tuning on pre-trained models, something that was time and resource consuming for this research. A promising direction would be to see how fine-tuning could enhance performance, and how continual learning would respond to this new architecture.

5 Qualitative Analysis

5.1 Language Order Analysis

For testing the significance of language order, we performed 25 independent experiments (5 random data batches per language \times 5 random seeds per order) to ensure stability and performed statistical significance for the theoretically best order, against the other variants.

Order	Fine	Coarse	p-value
RU→BG→PT→HI→EN	.350 \pm .017	.513 \pm .013	6.89×10^{-5}
RU→BG→HI→PT→EN	.323 \pm .022	.485 \pm .020	.601
HI→PT→RU→BG→EN	.312 \pm .005	.479 \pm .007	.025
RU→HI→PT→BG→EN	.210 \pm .016	.369 \pm .027	1.45×10^{-23}
PT→HI→RU→BG→EN	.289 \pm .011	.476 \pm .011	1.17×10^{-7}

Table 17: Impact of language order on model performance across different article batches and random seeds for sum aggregation strategy.

Language Order Analysis with Sum Strategy:

The in-theory best sequence

(RU→BG→PT→HI→EN) achieved the highest score for the Fine F1 score. The variant that starts with Bulgarian and follows Russian, led to a slight decrease in performance.

Our hypothesized worst language order (RU→HI→PT→BG→EN) gave poor performance, with a very small p-value (1.17e-07), meaning it’s very unlikely this poor performance occurred by chance.

Overall, the results show that when trying to create a model for English data, having certain languages early on in the sequence tends to help the model perform better.

Language Order Analysis with Weighted Strat-

egy: While we are at it, we also did a thorough

analysis for the weighted strategy, which outperformed the sum strategy.

Order	Fine	Coarse	p-value
RU→BG→PT→HI→EN	.423 ± .006	.583 ± .020	.068
BG→RU→PT→HI→EN	.355 ± 0.034	.501 ± .015	1.10×10^{-9}
HI→PT→RU→BG→EN	.398 ± 0.014	.571 ± .021	9.17×10^{-6}
RU→HI→PT→BG→EN	.440 ± .013	.611 ± .018	3.09×10^{-6}
PT→HI→RU→BG→EN	.405 ± 0.014	.576 ± .015	.0029

Table 18: Impact of language order using weighted average strategy across different article batches and random seeds.

Weighted strategy revealed different patterns compared to sum.

Both RU→BG→PT→HI→EN and RU→BG→HI→PT→EN orders maintain strong performance, their difference is not statistically significant ($p = 0.068$). Language order RU→HI→PT→BG→EN performs surprisingly well, better than our best order for sum strategy and contrasting with its poor performance under the same approach.

However, the weighted strategy appears to be more robust to order variations, showing generally higher performance across all orderings compared to sum strategy. This shows that embedding aggregation affects the importance of language order. Sum aggregation preserves all article information equally, making language order clear and much more significant. Weighted average weights sections by their length, it shows more balanced performance across different orders, making language order less significant to performance.

6 Acknowledgments

This research was conducted as an undergraduate semester thesis project.

I would like to thank Panos Louridas and John Pavlopoulos for their guidance and AUEB for supporting this work.

References

Travis G Coan, Constantine Boussalis, John Cook, and Mirjam O Nanko. 2021. Computer-assisted classification of contrarian claims about climate change.

Giovanni Da San Martino, Alberto Barrón-Cedeño, Henning Wachsmuth, Rostislav Petrov, and Preslav Nakov. 2020. *Semeval-2020 task 11: Detection of*

propaganda techniques in news articles. In *Proceedings of the Fourteenth Workshop on Semantic Evaluation*, pages 1377–1414. International Committee for Computational Linguistics.

Thomas G. Dietterich. 2000. Ensemble methods in machine learning.

Evangelia Gogoulou, Timothée Lesort, Magnus Boman, and Joakim Nivre. 2024. Continual learning under language shift. *arXiv preprint arXiv:2402.18449*. Correspondence: evangelia.gogoulou@ri.se.

James Kirkpatrick, Pascanu, and 1 others. 2017. Overcoming catastrophic forgetting in neural networks.

Bonka Kotseva, Irene Vianini, and Nikolaos Nikolaidis. 2023. Analyzing covid-19 misinformation narratives across multilingual sources.

Fei Mi, Liangwei Chen, Mengjie Zhao, Minlie Huang, and Boi Faltings. 2020. Continual learning for natural language generation in task-oriented dialog systems. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Online. Association for Computational Linguistics.

Robert Muller. 2018. Indictment of internet research agency. pages 1–37.

Jakub Piskorski, Nikolaos Nikolaidis, Nicolas Stefanovitch, Bonka Kotseva, Irene Vianini, Sopho Kharazi, and Jens P Linge. 2022. Exploring data augmentation for classification of climate change denial: Preliminary study.

Hannah Rashkin, Eunsol Choi, Jin Yea Jang, Svitlana Volkova, and Yejin Choi. 2017. Truth of varying shades: Analyzing language in fake news and political fact-checking. In *Proceedings of the 2017 Conference on Empirical Methods in Natural Language Processing, EMNLP '17*, pages 2931–2937, Copenhagen, Denmark.

David Sayce. 2022. [The number of tweets per day in 2022](#).

Cristina Tardaguila, Fabrício Benevenuto, and Pablo Ortellado. 2018. [Fake news is poisoning brazilian politics. WhatsApp can stop it](#). *The New York Times*.

Liyuan Wang, Xingxing Zhang, Hang Su, and Jun Zhu. 2023. [A comprehensive survey of continual learning: Theory, method and application](#). *arXiv preprint arXiv:2302.00487*.