

Temporal fine-tuning for early risk detection

Horacio Thompson^{1,2}, Esaú Villatoro-Tello³, Manuel Montes-y-Gómez⁴, and Marcelo Errecalde¹

¹ Universidad Nacional de San Luis, San Luis, Argentina

² Consejo Nacional de Investigaciones Científicas y Técnicas (CONICET)

³ Idiap Research Institute, Martigny, Switzerland

⁴ Instituto Nacional de Astrofísica, Óptica y Electrónica (INAOE), Puebla, Mexico

hthompson@unsl.edu.ar, esau.villatoro@idiap.ch, mmontesg@inaoep.mx,
merreca@unsl.edu.ar

Abstract. Early Risk Detection (ERD) on the Web aims to identify promptly users facing social and health issues. Users are analyzed post-by-post, and it is necessary to guarantee correct and quick answers, which is particularly challenging in critical scenarios. ERD involves optimizing classification precision and minimizing detection delay. Standard classification metrics may not suffice, resorting to specific metrics such as $ERDE\theta$ that explicitly consider precision and delay. The current research focuses on applying a multi-objective approach, prioritizing classification performance and establishing a separate criterion for decision time. In this work, we propose a completely different strategy, *temporal fine-tuning*, which allows tuning transformer-based models by explicitly incorporating time within the learning process. Our method allows us to analyze complete user post histories, tune models considering different contexts, and evaluate training performance using temporal metrics. We evaluated our proposal in the depression and eating disorders tasks for the Spanish language, achieving competitive results compared to the best models of MentalRiskES 2023. We found that temporal fine-tuning optimized decisions considering context and time progress. In this way, by properly taking advantage of the power of transformers, it is possible to address ERD by combining precision and speed as a single objective.

Keywords: Intelligent Systems · Machine Learning · Transformers · Early Risk Detection · Mental Health.

1 Introduction

One of the problems that has become relevant in recent years is Early Risk Detection (ERD) on the Web, which consists of correctly identifying risk users as soon as possible. It incorporates a significant complexity to standard classification problems since the users are analyzed post-by-post rather than processing the complete history. Challenges such as CLEF eRisk [8–12, 20, 17, 21] and MentalRiskES [15] have emerged to solve ERD problems in scenarios where users suffer from different mental disorders. In particular, MentalRiskES 2023 was

the first edition to propose tasks exclusively for the Spanish language. These challenges defined a testing environment for ERD problems within a post-round scheme. In every round, a post from each user is received, and it is necessary to respond to each to continue to the next round. Whether the user is detected as positive, a risk alarm is issued, and the analysis is ended in that post-round. Otherwise, if the user's responses were successively negative and the post history ends, it is concluded that the user is negative. The solutions are evaluated considering standard classification metrics, as well as temporal metrics such as $ERDE\theta$ [8] and F-latency [24].

Although precision is essential for ERD problems, as time progresses and decisions are delayed, speed becomes increasingly important, taking priority over precision. Because both are fundamental objectives to solve these problems, two possible approaches arise: multi-objective and combined single-objective. In the first, each objective is solved independently according to the priorities of the problem, typically addressing the precision and then optimizing the decision time. With the emergence of *transformers* [30], numerous studies have focused on this approach by prioritizing the correct classification of users. These models are trained and validated differently than in the testing environment, and this disparity complicates the optimal model selection. Moreover, the limited number of tokens in each architecture frequently leads to a partial observability scenario regarding the complete user's post history.

On the contrary, in a combined single-objective approach, precision and speed are concurrently considered during the learning process. That aspect has significant advantages over the standard (separated) multi-objective approach. First, considering input samples and time progress would allow for defining a single learning component where models can receive and integrate all the information needed for ERD. Furthermore, it would allow a scenario similar to the one used in the testing stage to be elegantly reproduced during training. As a consequence, optimal models can be obtained using the same temporal evaluation metrics utilized for assessing the ERD systems. However, despite these advantages, few works have applied this approach.

In this work, we propose *temporal fine-tuning*, a novel method that allows tuning transformer-based models and simultaneously optimizing precision and speed for ERD problems. The main contributions are: 1) transforming the input from the users to keep track of the delays, 2) defining a loss function to be optimized according to a temporal evaluation metric, and 3) implementing a training and validation procedure that allows selecting the optimal hyperparameters for a particular ERD problem. This work presents an innovative strategy that addresses a gap in the research domain, showing promising results that support the viability of our proposal when compared to state-of-the-art methods. In this way, our study offers new perspectives for future research.

2 Related work

Several studies have addressed ERD through transformer-based methods that follow the multi-objective approach. Most strategies solve ERD as a user classification problem by applying the fine-tuning process on models such as BERT and RoBERTa [1, 18, 3, 22], as well as pre-trained models for Spanish such as BETO and RoBERTuito [5, 19, 6]. Other proposals include embeddings extracted from transformers followed by different types of classifiers [16, 23, 27, 25, 31]. Although these strategies focused on optimizing classification performance, they also achieved temporal efficiency in an ERD environment, suggesting that classification robustness could impact decision speed. A more representative method within this approach was proposed in [14, 28, 29], addressing precision and speed separately. A BERT-based classifier was applied, followed by an independent component for decision-making, which was not considered within the learning process of the predictor. Sun et al. [26] argue that an early classifier should be accurate, but when decisions are unclear and delayed, adding a lookahead component to the classifier is necessary. At this point, the combined single-objective approach gains importance by considering precision and speed concurrently during learning. However, there is a limited number of works that use this approach. A study proposed a memorization network for affective states that is updated considering time [7], and the EARLIEST architecture for time series was adapted to balance precision and speed in decisions using reinforcement learning [13, 14]. In this work, we propose an original strategy that contributes to the research domain by applying the combined single-objective approach to solve ERD problems.

3 Temporal fine-tuning

In this section, we will address the most important aspects of our proposal to incorporate time during the learning process and then evaluate the models in an ERD environment. Fig. 1 shows the pipeline of our proposal.

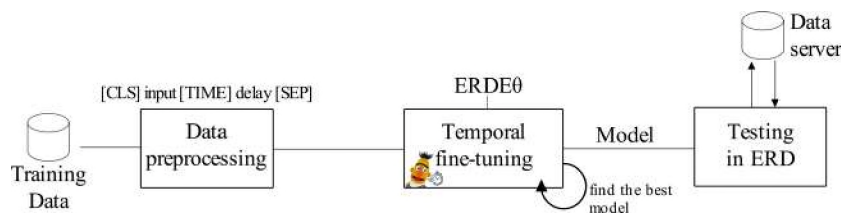


Fig. 1: Pipeline of the *temporal fine-tuning* process. The samples are modified including time, *temporal fine-tuning* is applied according to an $ERDE\theta$, and the best model is chosen to be evaluated in an ERD environment.

3.1 Data preprocessing

Time is explicitly included in the input samples, considering the number of posts that have been read until the moment (*delay*). For an arbitrary sample (input, label), the new input is defined as $\text{input}_{\text{delay}} = [\text{CLS}] \text{input} [\text{TIME}] \text{delay} [\text{SEP}]$. The [TIME] token is added to the model architecture, separating the post content of the moment it was read. For example, *hoy es un día triste* (today is a sad day) evaluated at $\text{delay} = 10$, the model would receive [CLS] *hoy es un día triste* [TIME] 10 [SEP].

3.2 Temporal fine-tuning process

Unlike standard classification tasks, where systems are evaluated with metrics such as F1 score, in ERD problems, the performance of the models is assessed according to specialized metrics that explicitly consider classification accuracy and delay in detecting positive cases. One of the most popular is ERDE_{θ} , defined as:

$$\text{ERDE}_{\theta}(d, k) = \begin{cases} c_{fp} \\ c_{fn} \\ lc_{\theta}(k) \cdot c_{tp} \\ 0 \end{cases} \quad (1)$$

c_{fp} : for false positives (FPs)
 c_{fn} : for false negatives (FNs)
 $lc_{\theta}(k) \cdot c_{tp}$: for true positives (TPs)
 0 : for true negatives (TNs)

where the latency cost (lc_{θ}) is:

$$lc_{\theta}(k) = 1 - \frac{1}{1 + e^{k-\theta}} \quad (2)$$

It means that, for a decision d that is a TP at time $k > \theta$, the penalty will depend on the c_{tp} value, where typically $c_{tp} = c_{fn} = 1$, i.e., the maximum penalty. In this way, the ERDE_{θ} metric could be helpful during the learning process since it evaluates the correctness and delay of the final decisions.

The *temporal fine-tuning* process is carried out in epochs with the typical training and validation stages. We propose to incorporate in both stages an evaluation scheme similar to the testing environment for ERD problems, where time is measured according to the number of posts requested to issue a response (*delay*). In each *delay*, post windows of length M are evaluated, which consist of concatenating the current post and the previous $M-1$. The *delays* are configured according to the window size. Fig. 2 shows that, for $M=10$, $\text{delay}=10$ evaluates from post 0 to post 9, $\text{delay}=20$ from post 10 to post 19, and continues until all users u_i are analyzed.

The *loss* function is calculated at the end of each *delay*, evaluating according to the post window over those users who have not yet completed their analysis,

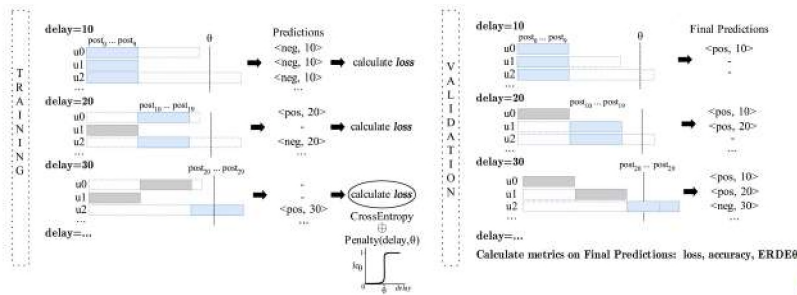


Fig. 2: *Temporal fine-tuning* scheme for an epoch. The training and validation stages are subdivided into *delays*, where users are evaluated according to post windows. In training, *loss* is calculated at the end of each *delay* considering CrossEntropy and lc_θ . In validation, the model performance is evaluated considering *loss*, *accuracy*, and $ERDE\theta$.

and the gradient is propagated to the entire transformer. We used the $ERDE\theta$ metric to design the *loss* function, but it was not included as indicated in (1) since it would not be differentiable. Instead, a linear and differentiable function was implemented by considering the classification performance (CrossEntropy) and harshly penalizing TPs that are delayed (using θ as a limit, according to (2)). In this way, as the *loss* is minimized, $ERDE\theta$ is also reduced, establishing it as the training objective. Listing 1 shows the procedural details for the *loss* function. The validation stage is evaluated with the same *delay* scheme, concluding the epoch by calculating the *loss*, *accuracy*, and $ERDE\theta$ metrics, which can be weighted to select the optimal model for an ERD problem.

3.3 Model testing in ERD

The best model obtained can be evaluated using a mock-server tool [14], which simulates ERD environments through rounds of posts and answers submissions, and it calculates the final results using different metrics. A client application was defined to interact with the server: when it receives a round of posts, the system preprocesses them by adding time, invokes the predictive model, and returns a response. We used a sliding post window, configured as in the learning stage, and a simple decision policy: if the probability exceeds a *threshold*, a user at-risk alarm is issued; otherwise, the analysis should continue.

4 Experimental results

The experiments were carried out on the *depression* and *eating disorders* tasks of MentalRiskES 2023, using the following datasets (supplied by the MentalRiskES Organizers and cannot be disclosed or shared publicly; more details in [15]): *train* and *trial* to train the models and *test* to evaluate them with the mock-server and

```

def temporal_loss(preds_labels, preds_times,
                 real_labels, real_times,  $\theta$ ):
    """ Calculates the temporal loss function.
    :param preds_labels: List of 1's and 0's (predicted labels).
    :param preds_times: List of integers (#posts read).
    :param real_labels: List of 1's and 0's (real labels).
    :param real_times: List of integers. (#total posts).
    :param  $\theta$ : Integer. Time limit for optimization (ERDE $\theta$ ).
    :return: Mean of the total loss. """
    # Calculate classification loss
    cls_loss = CrossEntropy(preds_labels, real_labels)
    # Calculate penalty for delay
    delay_loss = []
    for pred, real, pred_time, real_time
        in zip(preds_labels, real_labels, preds_times, real_times):
        if (pred==1 and pred==real) # TPs
            and (real_time<pred_time or  $\theta$ <preds_time): # Delayed
            delay_loss.append(1)
        else
            delay_loss.append(0)
    # Calculate total loss using classification and delay loss
    total_loss = []
    for cls, delay in zip(cls_loss, delay_loss):
        if delay==1:
            total_loss.append(delay)
        else
            total_loss.append(cls)

    return total_loss.mean()
    
```

Listing 1: Implementation of the *loss* function (pseudo-code)

compare results with other teams (Table 1). It can be noted that they present a relatively balanced distribution between classes. In addition, a model with an acceptable performance should complete the user evaluation before the mean number of posts (34.7 for depression and 27.9 for eating disorders in *test*). The maximum number of posts in *test* indicates the total number of rounds for the evaluation stage.

We applied *temporal fine-tuning* to the BETO model, a BERT variant trained on a large Spanish corpus [2], and ERDE30 as the training objective since the participants in MentalRiskES were ranked with this metric. *Delays* and post windows were configured using 10 posts, and the rest of the hyperparameters were the following: *learning_rate*= $3e-5$ (*depression*) and $5e-5$ (*eating disorders*), and *epoch*=10, *batch_size*=8, and *optimizer*=AdamW for both tasks. The best models were chosen by weighting the *accuracy* and ERDE30 metrics. We used *threshold* = 0.7 to detect positive users, and negative users were identified when they no longer had any posts. Besides, the *minDelay* parameter was added to

Table 1: Details of the *depression* and *eating disorders* corpora. The number of users (total, positives, and negatives) and the number of posts per user (mean, minimum, and maximum) are reported.

| | Corpus | #Users | | | #Posts per user | | |
|-------------------------|--------|--------|-----|-----|-----------------|-----|-----|
| | | Total | Pos | Neg | Mean | Min | Max |
| <i>Depression</i> | Train | 175 | 94 | 81 | 35.7 | 11 | 100 |
| | Trial | 10 | 6 | 4 | 62.4 | 11 | 100 |
| | Test | 149 | 68 | 81 | 34.7 | 11 | 100 |
| <i>Eating disorders</i> | Train | 175 | 74 | 101 | 33.9 | 11 | 50 |
| | Trial | 10 | 5 | 5 | 38.9 | 18 | 50 |
| | Test | 150 | 64 | 86 | 27.9 | 11 | 50 |

establish a minimum wait to start issuing alarms, testing $minDelay=10$ and $minDelay=5$.

We included in this study the baseline *sliding_window* model that was trained and validated considering the *delay* scheme without including time in the input samples, along with the same hyperparameters and settings mentioned above. Table 2 shows the results obtained by our models. The three best results are included based on the ranking of the MentalRiskES Organizers according to ERDE30 among more than 25 proposals for both tasks[15]. We also show the UNSL#0 results (obtained by our laboratory in the challenge), which applied classic fine-tuning to the BETO model and a decision-making component based on the history of previous predictions [29].

4.1 Depression results

The models obtained with *temporal fine-tuning* achieved the second-best place by notably outperforming UNSL#1, considering the overall performance among all metrics. In particular, *temporal_ft-minDelay10* obtained the second-best ERDE30 and the best classification performances, and *temporal_ft-minDelay5* achieved the best F-latency. It is highlighted that the *temporal_ft* models outperformed the respective *sliding_window* models in all metrics, with ERDE30 values very similar to each other. Additionally, our proposals remarkably outperformed UNSL#0 and the mean performance among all teams across most metrics.

4.2 Eating disorders results

achieved remarkable results The *temporal fine-tuning* models achieved remarkable results among the top positions, especially *temporal_ft-minDelay5*, which reached the third-best place, notably outperforming CIMAT-NLP-GTO#1 in all metrics. The *temporal_ft-minDelay10* model was competitive with CIMAT-NLP-GTO#1 on ERDE30, outperforming it in classification metrics. Besides, the *temporal_ft* models outperformed the respective *sliding_window* models in

Table 2: Results obtained considering the classification (Precision, Recall, and F1) and early classification (ERDE5, ERDE30, and F-latency) metrics for both tasks. The three best results are shown based on the ranking of MentalRiskES Organizers according to ERDE30, as well as the mean values among all results. Values in bold and underlined depict 1st and 2nd performance for each metric, respectively. Values closer to 0 show better performances for the ERDE metric; for the rest of the metrics, values closer to 1 are preferred.

| (a) Depression | | | | | | |
|------------------------------|-------------|-------------|-------------|--------------|--------------|-------------|
| Rank-Team#Model | P | R | F1 | ERDE5↓ | ERDE30↓ | F-latency |
| 1-SINAI-SELA#0 [5] | 0.78 | 0.74 | 0.72 | <u>0.395</u> | 0.140 | 0.72 |
| 2-UNSL#1 [29] | 0.79 | 0.76 | 0.73 | 0.567 | 0.148 | 0.61 |
| 3-BaseLine-Deberta#0 [15] | 0.79 | 0.69 | 0.64 | 0.303 | 0.153 | 0.72 |
| 14-UNSL#0 | 0.75 | 0.74 | 0.73 | 0.551 | 0.188 | 0.59 |
| <i>MentalRiskES2023-mean</i> | 0.73 | 0.66 | 0.62 | 0.383 | 0.232 | 0.60 |
| sliding_window_minDelay10 | 0.79 | 0.78 | 0.78 | 0.526 | 0.150 | 0.64 |
| sliding_window_minDelay5 | 0.80 | 0.76 | 0.75 | 0.473 | 0.140 | <u>0.71</u> |
| temporal_ft-minDelay10 | 0.83 | 0.83 | 0.83 | 0.486 | <u>0.146</u> | 0.66 |
| temporal_ft-minDelay5 | <u>0.81</u> | <u>0.81</u> | <u>0.81</u> | 0.440 | 0.150 | 0.72 |

| (b) Eating disorders | | | | | | |
|------------------------------|-------------|-------------|-------------|--------------|--------------|-------------|
| Rank-Team#Model | P | R | F1 | ERDE5↓ | ERDE30↓ | F-latency |
| 1-CIMAT-NLP-GTO#0 [4] | 0.96 | 0.97 | 0.97 | 0.334 | 0.018 | 0.86 |
| 2-UNSL#1 | <u>0.91</u> | <u>0.92</u> | <u>0.91</u> | 0.433 | <u>0.045</u> | <u>0.78</u> |
| 3-CIMAT-NLP-GTO#1 | 0.87 | 0.87 | 0.85 | 0.379 | 0.065 | 0.76 |
| 11-UNSL#0 | 0.82 | 0.79 | 0.75 | 0.502 | 0.105 | 0.67 |
| <i>MentalRiskES2023-mean</i> | 0.82 | 0.79 | 0.75 | 0.322 | 0.122 | 0.71 |
| sliding_window_minDelay10 | 0.85 | 0.85 | 0.85 | 0.454 | 0.108 | 0.61 |
| sliding_window_minDelay5 | 0.87 | 0.87 | 0.87 | 0.371 | 0.088 | 0.74 |
| temporal_ft-minDelay10 | 0.90 | 0.90 | 0.90 | 0.448 | 0.069 | 0.65 |
| temporal_ft-minDelay5 | <u>0.91</u> | 0.91 | <u>0.91</u> | <u>0.365</u> | 0.062 | 0.77 |

all metrics. Considering the mean among all teams, our results achieved better performance in Precision, Recall, F1, and ERDE30.

4.3 Analysis of *temporal fine-tuning* learning

We analyze the model behavior during learning using the depression task as an example. Fig. 3(a) shows that the model solves the early classification of users as the epochs progress. In the first epochs, the model tends to issue responses before $\theta=30$ due to the penalty. As epochs progress, final decisions are optimized by adjusting correctness and delay. Therefore, TPs tend to be obtained before θ , while TNs when the user posts end, many of them in instances greater than θ . Fig. 3(b) shows examples of FNs, FPs, and delayed TPs that were detected in epoch 0 and corrected in epoch 1. The best model of the validation stage was

obtained at epoch 5 with optimal weighted performance (Fig. 3(c)). In this way, *temporal fine-tuning* could learn when to end users' analysis during training, making it easier to adapt the model to an ERD environment without adding complex independent decision components.

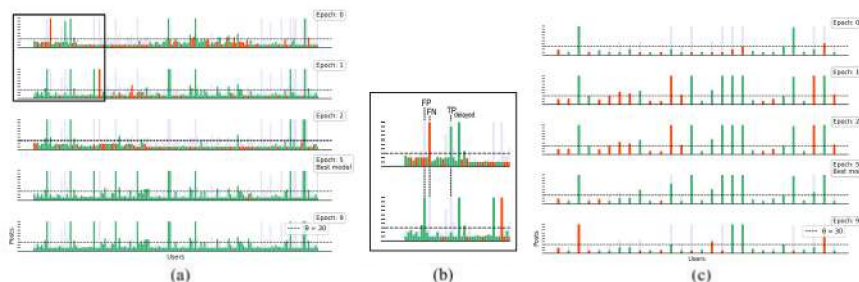


Fig. 3: Model learning during *temporal fine-tuning* for *depression* task. (a) Training stage; (b) Some cases of (a) detected in epoch:0 and corrected in epoch:1; (c) Validation stage. On the x-axis, the vertical bars depict the model's decisions, and their length shows the number of posts read and the instance of the final response. Green bar: correct decision; red bar: wrong decision; gray bar: unread posts during analysis. The y-axis shows *delays* every 10 posts. The dashed horizontal line denotes the limit at $\theta=30$.

4.4 *Temporal fine-tuning* vs. mock-server

Fig. 4 shows the models' behavior when evaluating the *test* data using *temporal fine-tuning* validation and mock-server. In the first, models can only decide every 10 posts, while in the mock-server, they decide post-by-post. Despite this, similar behaviors were obtained, demonstrating that our proposal design would allow us to know the expected performance in an ERD environment and appropriately select the best models.

4.5 Temporal representation

Due to the inclusion of time, the semantics of the words could change depending on the instance where they are analyzed. For example, for the sentence *Hace mucho que no duermo bien...* (I haven't slept well for a long time...) evaluated at distinct times, different representations were obtained, directly influencing the classification result (Fig. 5). A more in-depth analysis in the future will allow us to evaluate the impact of temporal representations and identify critical symptoms linked to diverse pathologies.

4.6 Relevant themes in positive users

A last analysis in our study consisted of obtaining the most relevant themes found in the positive users. Table 3 shows the different themes found in positive

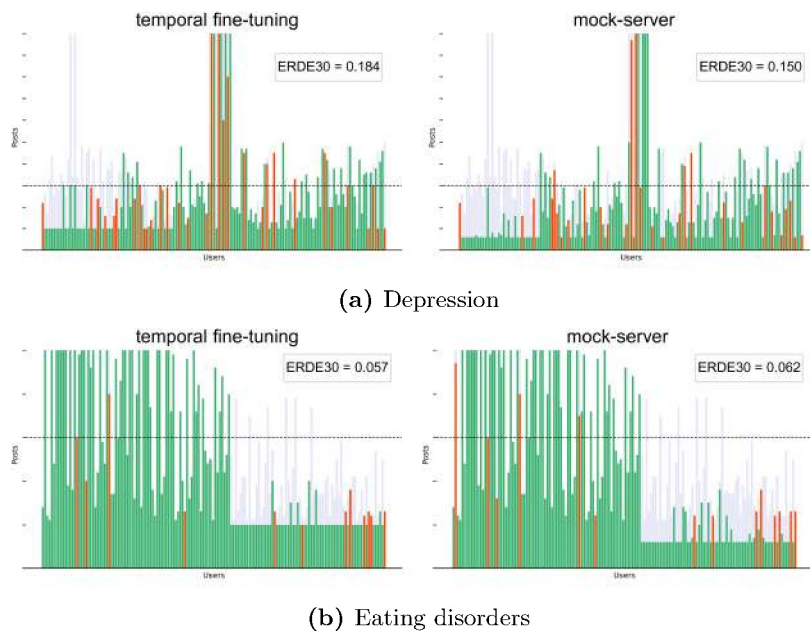


Fig. 4: Evaluation of the models using *temporal fine-tuning* validation and mock-server for both tasks.

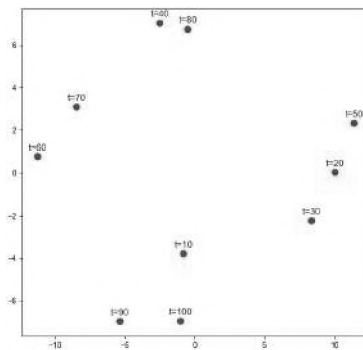


Fig. 5: Temporal representation of a sentence evaluated at times $t=[10, 20, \dots, 100]$. Green dot: positive decision; red dot: negative decision. Embeddings were extracted from the fine-tuned model ([CLS] token), followed by dimensionality reduction through Principal Component Analysis (PCA).

users that were correctly detected. For each TP, we extracted the post window used by the models at the decision moment and used Google’s Gemini model (<https://gemini.google.com/>) via API to capture the themes related to each task.

Table 3: Themes relevant to depression and eating disorders tasks. For each positive user correctly detected, the post window used by the models at the decision moment was extracted. Then, they were analyzed by the Gemini model to capture the relevant themes.

| | |
|--------------------------------|--|
| <i>Depression</i> | Loneliness, Depression, Suicidal ideation, Hopelessness, Anxiety, Isolation, Sadness, Low self-esteem, Loss and grief, Self-harm, Eating disorders, Family conflicts, Emotional dependence, Post-traumatic stress, Sleep disorders, Loss of pleasure, Medication dependency, Stress, Lack of concentration, Fatigue, Failed relationships. |
| <i>Eating disorders</i> | Eating disorders, Body image, Bulimia, Food restriction, Intermittent fasting, Binge eating, Weight control, Advices, Compensatory behaviors, Medical tests, Physical exercise, Parental control, Self-harm, Poor eating habits, Medication dependency, Isolation, Relapses, Vomiting, Concerns, Diets. |

5 Conclusions and Future Work

In this work, we proposed *temporal fine-tuning*, a novel tuning method for pre-trained models based on transformers. We explicitly incorporated time during the learning process through a *delay* scheme and a *loss* function designed considering the ERDE θ metric. We found that our method optimizes the decisions, taking into account different contexts and time progress. It also allows us to evaluate training performance using temporal metrics and finally select the optimal models for an ERD problem. We obtained remarkable results in the depression and eating disorders tasks for the Spanish language, being competitive with the best models of MentalRiskES 2023. In this way, by properly taking advantage of the power of transformers, it is possible to address ERD by combining precision and speed as a *single objective*.

Although we have contributed to a less-explored language, it would be interesting to analyze our proposal in other languages, including English. Besides, new *loss* functions could be designed considering other temporal metrics such as F-latency. On the other hand, it would be important to analyze the impact of the temporal representations obtained with our method, which would contribute to the interpretability of transformers. More research linked to the *combined single-objective* approach is necessary since, in our view, it could be the most appropriate way to address ERD problems.

Acknowledgments. This work was developed at the Laboratorio de Investigación y Desarrollo en Inteligencia Computacional (LIDIC) and was supported by a grant from Universidad Nacional de San Luis (UNSL), Argentina [PROICO 03-0620]. This work is part of the Doctoral thesis of Horacio Thompson.

References

1. Bucur, A.M., Cosma, A., Dinu, L.P.: Early risk detection of pathological gambling, self-harm and depression using bert (2021)
2. Cañete, J., Chaperon, G., Fuentes, R., Ho, J.H., Kang, H., Pérez, J.: Spanish pre-trained bert model and evaluation data (2023)
3. Devaguptam, S., Kogatam, T., Kotian, N., M, A.K.: Early detection of depression using bert and deberta. In: CLEF (Working Notes). pp. 875–882 (2022)
4. Echeverría-Barú, F., Sánchez-Vega, F., López-Monroy, A.P.: Cimat-nlp-gto at mentalrisks 2023: Early detection of mental disorders in spanish messages using style based models and bert models (2023)
5. González-Silot, S., Martínez-Cámara, E., Ureña-López, L.A.: Sinai at mentalrisk: Using emotions for detecting depression. In: IberLEF (Working Notes). CEUR Workshop Proceedings (2023)
6. de Jesús-García-Santiago, M., Sánchez-Vega, F., López-Monroy, A.P.: Improving transformer by instance packaging for mental illnesses identification (2023)
7. Kang, X., Dou, R., Yu, H.: Tual at erisk 2022: Exploring affective memories for early detection of depression. In: CLEF (Working Notes). pp. 1026–1037 (2022)
8. Losada, D.E., Crestani, F.: A test collection for research on depression and language use. In: International Conference of the Cross-Language Evaluation Forum for European Languages. pp. 28–39. Springer (2016)
9. Losada, D.E., Crestani, F., Parapar, J.: erisk 2017: Clef lab on early risk prediction on the internet: experimental foundations. In: Experimental IR Meets Multilinguality, Multimodality, and Interaction: 8th International Conference of the CLEF Association, CLEF 2017, Dublin, Ireland, September 11–14, 2017, Proceedings 8. pp. 346–360. Springer (2017)
10. Losada, D.E., Crestani, F., Parapar, J.: Overview of erisk: early risk prediction on the internet. In: Experimental IR Meets Multilinguality, Multimodality, and Interaction: 9th International Conference of the CLEF Association, CLEF 2018, Avignon, France, September 10–14, 2018, Proceedings 9. pp. 343–361. Springer (2018)
11. Losada, D.E., Crestani, F., Parapar, J.: Overview of erisk 2019 early risk prediction on the internet. In: Experimental IR Meets Multilinguality, Multimodality, and Interaction: 10th International Conference of the CLEF Association, CLEF 2019, Lugano, Switzerland, September 9–12, 2019, Proceedings 10. pp. 340–357. Springer (2019)
12. Losada, D.E., Crestani, F., Parapar, J.: erisk 2020: Self-harm and depression challenges. In: European conference on information retrieval. pp. 557–563. Springer (2020)
13. Loyola, J.M., Burdisso, S., Thompson, H., Cagnina, L.C., Errecalde, M.: Unsl at erisk 2021: A comparison of three early alert policies for early risk detection. In: CLEF (Working Notes). pp. 992–1021 (2021)
14. Loyola, J.M., Thompson, H., Burdisso, S., Errecalde, M.: Unsl at erisk 2022: Decision policies with history for early classification. In: CLEF (Working Notes). pp. 947–960 (2022)

15. Mármol-Romero, A.M., Moreno-Muñoz, A., Plaza-del Arco, F.M., Molina-González, M.D., Martín-Valdivia, M.T., Ureña-López, L.A., Montejo-Raéz, A.: Overview of mentalrisques at iberlef 2023: Early detection of mental disorders risk in spanish. *Procesamiento del Lenguaje Natural* **71**, 329–350 (2023)
16. Mármol-Romero, A.M., Zafra, S.M.J., del Arco, F.M.P., Molina-González, M.D., Valdivia, M.T.M., Montejo-Raéz, A.: Sinai at erisk@ clef 2022: Approaching early detection of gambling and eating disorders with natural language processing. In: *CLEF (Working Notes)*. pp. 961–971 (2022)
17. Martín-Rodilla, P., Losada, D.E., Crestani, F.: Overview of erisk 2022: Early risk prediction on the internet. In: *Experimental IR Meets Multilinguality, Multimodality, and Interaction: 13th International Conference of the CLEF Association, CLEF 2022, Bologna, Italy, September 5–8, 2022, Proceedings*. vol. 13390, p. 233. Springer Nature (2022)
18. Pan, R., Díaz, J., Valencia-García, R.: Umuteam at erisk@ clef 2023 shared task: transformer models for early detection of pathological gambling, depression, and eating disorder. pp. 18–21 (2023)
19. Pan, R., Garcia-Díaz, J., Valencia-García, R.: Umuteam at mental-risques2023@ iberlef: Transformer and ensemble learning models for early detection of eating disorders and depression. In: *IberLEF (Working Notes)*. CEUR Workshop Proceedings (2023)
20. Parapar, J., Martín-Rodilla, P., Losada, D.E., Crestani, F.: Overview of erisk at clef 2021: Early risk prediction on the internet (extended overview). *CLEF (Working Notes)* pp. 864–887 (2021)
21. Parapar, J., Martín-Rodilla, P., Losada, D.E., Crestani, F.: erisk 2023: Depression, pathological gambling, and eating disorder challenges. In: *European Conference on Information Retrieval*. pp. 585–592. Springer (2023)
22. Ramos, L.V., García, C.M., Vázquez, J.M., Pachón, V.: I2c-uhu at mental-risques 2023: Detecting and identifying mental disorder risks in social media using transformer-based models (2023)
23. Rujas, M., Merino-Barbancho, B., Arroyo, P., Fico, G.: Development of a natural language processing-based system for characterizing eating disorders. In: *IberLEF (Working Notes)*. CEUR Workshop Proceedings (2023)
24. Sadeque, F., Xu, D., Bethard, S.: Measuring the latency of depression detection in social media. pp. 495–503 (2018)
25. Stalder, S., Zankov, E.: Zhaw at erisk 2022: Predicting signs of pathological gambling-glove for snowy days. In: *CLEF (Working Notes)*. pp. 987–994 (2022)
26. Sun, C., Li, H., Song, M., Hong, S.: A ranking-based cross-entropy loss for early classification of time series. *IEEE Transactions on Neural Networks and Learning Systems* (2023)
27. Talha, A., Basu, T.: A natural language processing based risk prediction framework for pathological gambling. pp. 18–21 (2023)
28. Thompson, H., Cagnina, L., Errecalde, M.: Strategies to harness the transformers’ potential: Unsl at erisk 2023 (2023)
29. Thompson, H., Errecalde, M.: Early detection of depression and eating disorders in spanish: Unsl at mentalrisques 2023 (2023)
30. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
31. Wu, S.H., Qiu, Z.J.: Cyut at erisk 2022: Early detection of depression based-on concatenating representation of multiple hidden layers of roberta model. In: *CLEF (Working Notes)*. pp. 1014–1025 (2022)