



Figure 4: The Transformer-Based Model

$$\hat{Y}^{<t>} = \{Softmax(z_0^{<t>}), Softmax(z_1^{<t>})\} \quad (3)$$

Where  $z_i^{<t>}$  is the logit of the last output layer for the  $t^{th}$  token and  $i^{th}$  class.

In both the Transformers and BiLSTM-ELMo approaches, the Binary Cross-Entropy (BCE) Loss as well as the KL-Divergence (KLD) Loss were used to train the models. The  $Match_m$  score is used as an evaluation for all our models. The equations for both the loss functions are as follows:

$$BCE(y^{<t>}, \hat{y}^{<t>}) = -y^{<t>} \cdot \log(\hat{y}^{<t>}) - (1 - y^{<t>}) \cdot \log(1 - \hat{y}^{<t>}) \quad (4)$$

Where  $y^{<t>} \in \{0, 1\}$  is the true label for emphasis laid on each token and  $\hat{y}^{<t>}$  is the output of the sigmoid activation for each token.

$$KLD(Y^{<t>} || \hat{Y}^{<t>}) = \sum Y^{<t>} \cdot \log\left(\frac{Y^{<t>}}{\hat{Y}^{<t>}}\right) \quad (5)$$

Where  $Y^{<t>}$  is the true probability distribution for the emphasis laid on each token and  $\hat{Y}^{<t>}$  is the output distribution of the softmax activation for each token.

We use the Adam Optimizer for training the models with a learning rate of  $1e-4$  for the BiLSTM-ELMo model for 100 epochs and  $2e-5$  for the Transformer-based models for 100 epochs. The training was performed on 1 NVIDIA Titan X GPU. Our code is available on Github<sup>5</sup>.

## Results

In Table 6 we present scores for both our BiLSTM-ELMo and Transformers approach trained on both BCE Loss and

KLDivergence Loss for LDL. As we can see in the results, LDL as used by (Shirani et. al 2019) doesn't give a huge improvement over results and at times even diminishes the results.

| Model                              | Dev          | Test         |
|------------------------------------|--------------|--------------|
| BiLSTM-ELMo (Baseline)             | -            | 0.475        |
| BiLSTM-ELMo (POS)                  | 0.497        | 0.484        |
| BiLSTM-ELMo (POS) (LDL)            | 0.501        | 0.506        |
| BiLSTM-ELMo (POS, Keyphrase)       | 0.515        | 0.496        |
| BiLSTM-ELMo (POS, Keyphrase) (LDL) | 0.504        | 0.501        |
| XLNet                              | <b>0.536</b> | <b>0.514</b> |
| XLNet (LDL)                        | 0.529        | 0.491        |
| RoBERTa                            | 0.51         | 0.485        |
| RoBERTa (LDL)                      | 0.515        | 0.47         |

Table 6: Performance of BiLSTM-ELMo and Transformers approach on development and Test set. The results are expressed in terms of average  $Match_m$  for  $m \in \{1, 5, 10\}$ . LDL indicates that label distribution learning was employed to train the model with KL-Divergence as the loss function, Binary Cross Entropy otherwise. For BiLSTM-ELMo model the extra features concatenated at the attention layer have been mentioned with each experiment. Baseline. indicates the scores by the baseline model defined by ?

For our final submissions, we tried an ensemble of scores from different models shown in Table 7. Our best scores on the Evaluation leaderboard were obtained using an ensemble of XLNet and RoBERTa with LDL where we stood 3rd. Meanwhile, our best scores on the Post-Evaluation leaderboard were obtained using an ensemble of XLNet and BiLSTM-ELMo approach with POS tags and Keyphrase Feature where we currently stand 1st on the leaderboard.

Additionally, we also ran experiments by dividing the presentations into their constituent sentences in the train and development data. Thus each training instance now corre-

<sup>5</sup><https://github.com/reasonalkumar/CAD21-AAAI21>

| Model                           | Dev         | Test         |
|---------------------------------|-------------|--------------|
| XLNet + RoBERTa (LDL)           | 0.547       | 0.518        |
| XLNet + BiLSTM-ELMo (Keyphrase) | 0.538       | 0.532        |
| XLNet + BiLSTM-ELMo (LDL)       | <b>0.55</b> | <b>0.543</b> |

Table 7: Performance of different ensemble models

sponds to a particular sentence belonging to a presentation slide in the original corpus. The development set results can be found in Table 8. The evaluation scheme used in this experiment uses the same  $Match_m$  as described in the Evaluation Metric section but with  $m = 1, 2, 3, 4$  as used in ?.

| Model             | Dev          |
|-------------------|--------------|
| XLNet             | <b>0.758</b> |
| XLNet (LDL)       | 0.757        |
| RoBERTa           | 0.743        |
| RoBERTa (LDL)     | 0.745        |
| BiLSTM-ELMo       | 0.751        |
| BiLSTM-ELMo (LDL) | 0.752        |

Table 8: Sentence-wise results on the Development set

- i) It is **extremely important** that parents take time to **SLOW DOWN** and give their child their **undivided attention**. The **importance** of that can not be **over-emphasized**.
- ii) It is **extremely important** that parents take time to **SLOW DOWN** and give their child their **undivided attention**. The **importance** of that can not be **over-emphasized**.
- iii) It is **extremely important** that parents take time to **SLOW DOWN** and give their child their **undivided attention**. The **importance** of that can not be **over-emphasized**.
- iv) It is **extremely important** that parents take time to **SLOW DOWN** and give their child their **undivided attention**. The **importance** of that can not be **over-emphasized**.

Figure 5: Emphasis Heatmaps i) Ground Truth ii) BiLSTM-ELMo iii) XLNet iv) Best Ensemble Model

## Analysis

### Length vs Performance

We wanted to understand how the performance of our models was affected by the length of the instances. Table 9 summarizes the performance of our best performing single model, i.e, XLNet on the development set divided into three sets, Short ( $\leq 40$  tokens, 80 samples), Medium (40 to 90 tokens, 262 samples), and Long ( $>90$  tokens, 50 samples). As we can see, the model performance deteriorates with the increasing length of the instances.

|                                | XLNet        |
|--------------------------------|--------------|
| Small ( $\leq 40$ )            | <b>0.648</b> |
| Medium ( $>40$ and $\leq 90$ ) | 0.549        |
| Large ( $>90$ )                | 0.42         |

Table 9: Average  $Match_m$  for best performing XLNet model on different size of instances in the development set

### Emphasis vs Parts of Speech

Table 10 shows POS (Parts of Speech) tags vs. average emphasis on the development dataset. We did this experiment

to understand how our model predictions performed on each POS tag when compared to the actual human-annotated emphasis scores on the development set. We noticed that the original average emphasis scores were highest on Adjectives followed by Noun. On comparing our models, we found that XLNet was able to almost accurately predict the emphasis scores on Adjectives and Noun respectively, and BiLSTM-ELMo also had the highest predictions on Adjectives and Noun respectively. We also noticed that XLNet did a better job on predicting the emphasis score on different POS tags where the predictions were either very close to the human scores or marginally lesser. On the other hand, we noticed that BiLSTM-ELMo’s predictions fell short by bigger margins when compared to XLNet and gave more emphasis to Adverbs than that in the development set.

| POS        | Count | Human | BiLSTM | XLNet |
|------------|-------|-------|--------|-------|
| Noun       | 4719  | 0.169 | 0.134  | 0.168 |
| Verb       | 1420  | 0.118 | 0.083  | 0.113 |
| Adjectives | 982   | 0.186 | 0.140  | 0.181 |
| Det        | 634   | 0.062 | 0.029  | 0.042 |
| Adverbs    | 347   | 0.111 | 0.068  | 0.103 |
| Pronouns   | 165   | 0.040 | 0.068  | 0.022 |
| Punct      | 2082  | 0.034 | 0.015  | 0.025 |

Table 10: POS tags vs. average emphasis on development dataset

## Conclusion

In this paper, we present our approach to AAIL-CAD21 shared task: Predicting Emphasis in Presentation Slides. Our best submission gave us an average  $Match_m$  of 0.518 placing us 3<sup>rd</sup> on the Evaluation phase leaderboard and an average  $Match_m$  of 0.543 placing us 1<sup>st</sup> on the Post-Evaluation leaderboard at the time of writing the paper. Future work includes using a hierarchical approach to emphasis prediction as a sequence labeling task using both sentence-level (individual sentence in a slide) and slide-level representations of a word (?).

## Acknowledgement

Rajiv Ratn Shah is partly supported by the Infosys Center for AI at IIIT Delhi. We also thank Sunny Dsouza and Gautam Maurya for their detailed and valuable feedback. Neque tenetur asperiores repellendus maiores aspernatur molestiae placeat, distinctio quod tempore ab sit eos id repellat enim quis soluta tempora, quisquam praesentium similique ea error, voluptate blanditiis harum architecto aliquid recusandae omnis, nihil eum culpa dicta ut repudiandae facere dolorum obcaecati asperiores sequi est?Quod sapiente laboriosam libero maiores perspiciatis aliquam, omnis dolorem excepturi, esse possimus debitis dolore animi laudantium quaerat dicta excepturi mollitia nam ducimus?Porro minus maxime nisi consequatur ab necessitatibus nesciunt, odit harum commodi, laudantium quas tenetur inventore beatae aperiam laboriosam, quaerat dolor quam modi, quidem porro sequi aliquam