# Fake News Detection Project 2

Marcel Affi, Mateusz Wójcik, Arkadiusz Zalas

# Agenda

- Project Contributions
- Preprocessing Techniques
- LIAR-PLUS dataset
- CT-FAN dataset
- Word2vec with LSTM model
- BERT Experiments
- SentiGAN
- CatGAN
- Improvements & Further Research

# Project Contributions

We continued working on fake news classification task with the following extensions.

- multi-class classification with the degree of certainty that a piece of news is fake (for example, 0-5 scale)
- datasets with extended input, for example justification
- providing analysis how the choice of inputs influence the models' performance
- application of generative models: SentiGAN and CatGAN
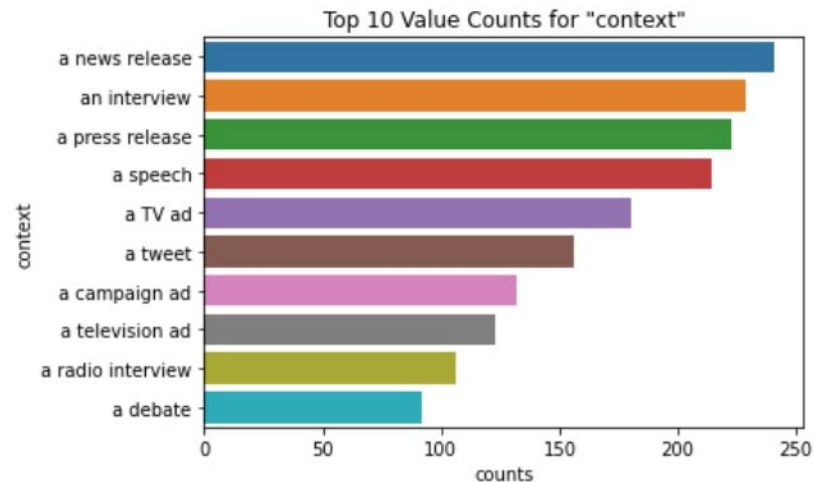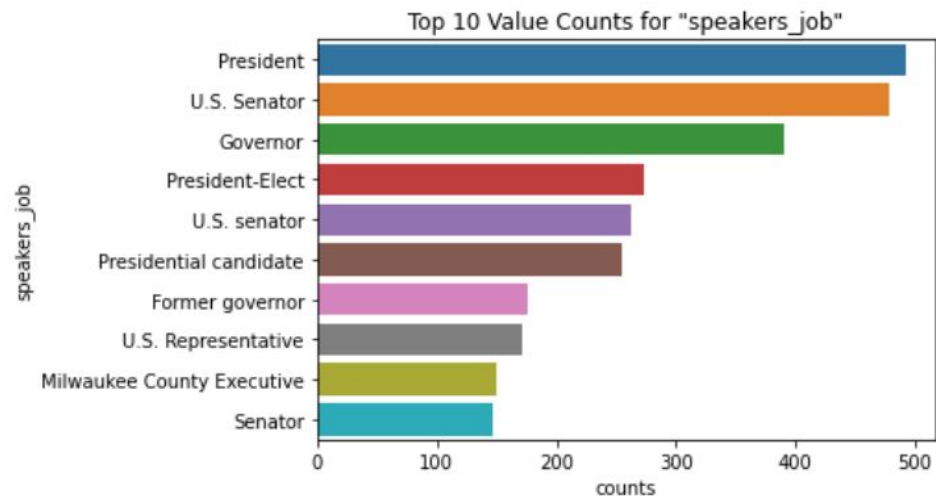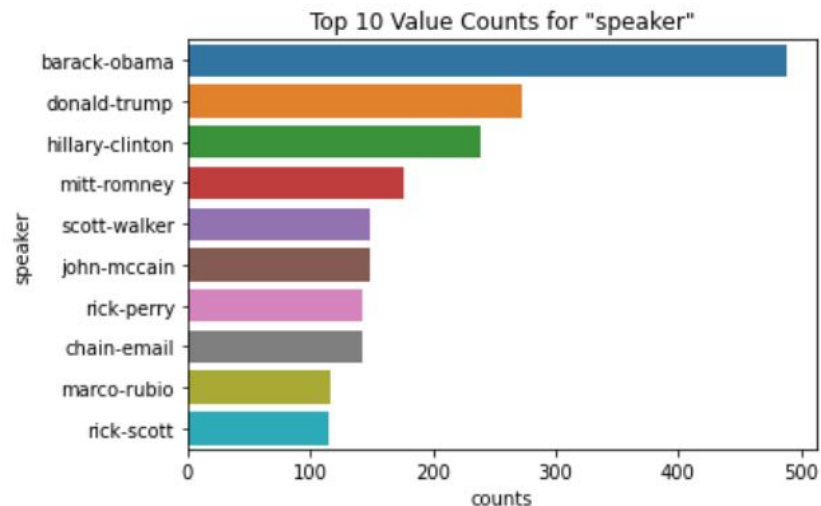
# Preprocessing Techniques

Sets of functions that we applied to prepare data for models
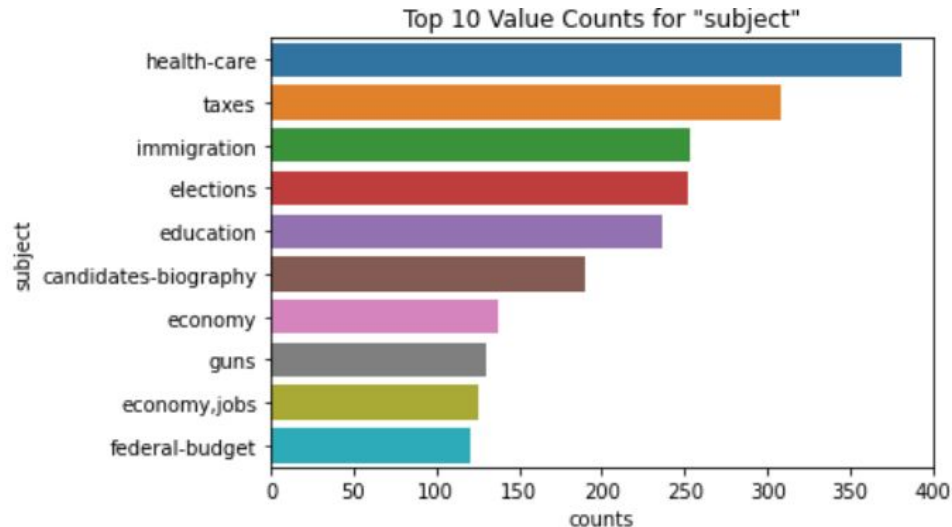
- Using regular expressions to remove:
  - punctuations
  - special symbols like $, #, &
  - digits <- we only replaced digits with words
  - URLs
  - wide spaces
  - single characters
- removing stop words
- lemmatization from different sources (NLTK, SpaCy)
- tokenization suited for a particular architecture
- other techniques required for the models, e.g. padding

# LIAR-PLUS Dataset

The extension of the dataset manually labeled by PolitiFact consisting of short statements, published in **2018**.

- multiple labels: false, pants-fire, barely-true, half-true, mostly-true, true
- around 12800 of observations
- well-balanced set, ranging from ~2000 to ~2600 instances per class
- attributes such as: label, statement, subject, speaker, speaker's job title, party affiliation...
- evidence sentences that have been automatically extracted from full-text verdict reports written by journalists in Politifact

**Top 10 Value Counts for "subject"**

| subject | counts |
| --- | --- |
| health-care | ~385 |
| taxes | ~305 |
| immigration | ~253 |
| elections | ~251 |
| education | ~238 |
| candidates-biography | ~192 |
| economy | ~133 |
| guns | ~128 |
| economy,jobs | ~122 |
| federal-budget | ~118 |

**Top 10 Value Counts for "speaker"**

| speaker | counts |
| --- | --- |
| barack-obama | ~488 |
| donald-trump | ~270 |
| hillary-clinton | ~238 |
| mitt-romney | ~174 |
| scott-walker | ~147 |
| john-mccain | ~145 |
| rick-perry | ~140 |
| chain-email | ~140 |
| marco-rubio | ~110 |
| rick-scott | ~110 |

**Top 10 Value Counts for "speakers_job"**

| speakers_job | counts |
| --- | --- |
| President | ~490 |
| U.S. Senator | ~475 |
| Governor | ~393 |
| President-Elect | ~265 |
| U.S. senator | ~258 |
| Presidential candidate | ~248 |
| Former governor | ~178 |
| U.S. Representative | ~178 |
| Milwaukee County Executive | ~158 |
| Senator | ~155 |

**Top 10 Value Counts for "context"**

| context | counts |
| --- | --- |
| a news release | ~243 |
| an interview | ~231 |
| a press release | ~225 |
| a speech | ~216 |
| a TV ad | ~183 |
| a tweet | ~158 |
| a campaign ad | ~135 |
| a television ad | ~127 |
| a radio interview | ~110 |
| a debate | ~93 |

# Example Justification

**Statement:** Oregon is the only state out of the 50 states in the USA that continues to pay 100% of the medical benefits for its employees and their families.
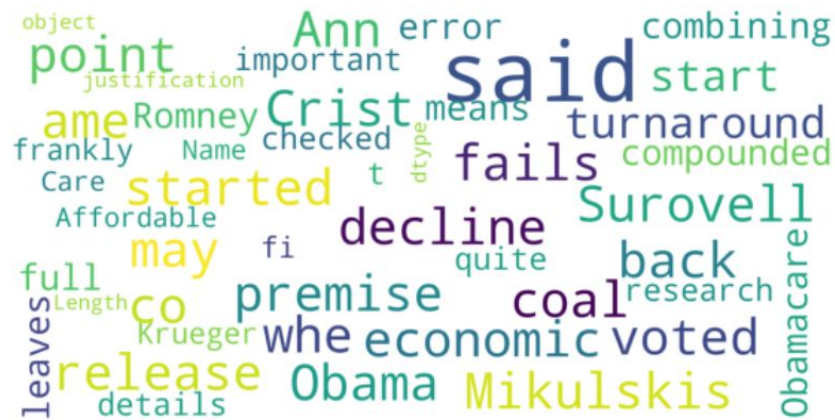
**Label:** false

**Justification:** Oregon is one of two states that covers a full range of benefits -- health, dental and vision -- but it's one of at least four that covers the premiums for its lowest-cost health plan. Richardsons larger point, that Oregon is in a shrinking group of states that do so is certainly a strong and valid one. However, he undercuts his argument by resorting to hyperbole.

"statement" word cloud

"justification" word cloud

"false" word cloud

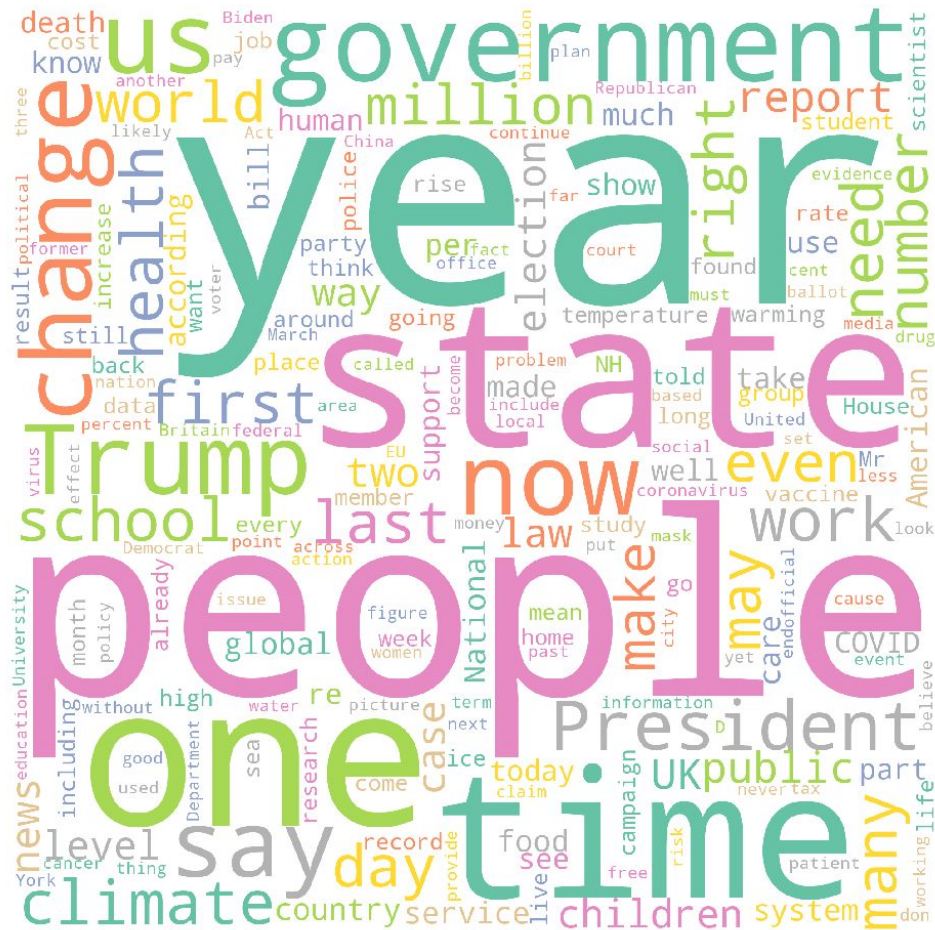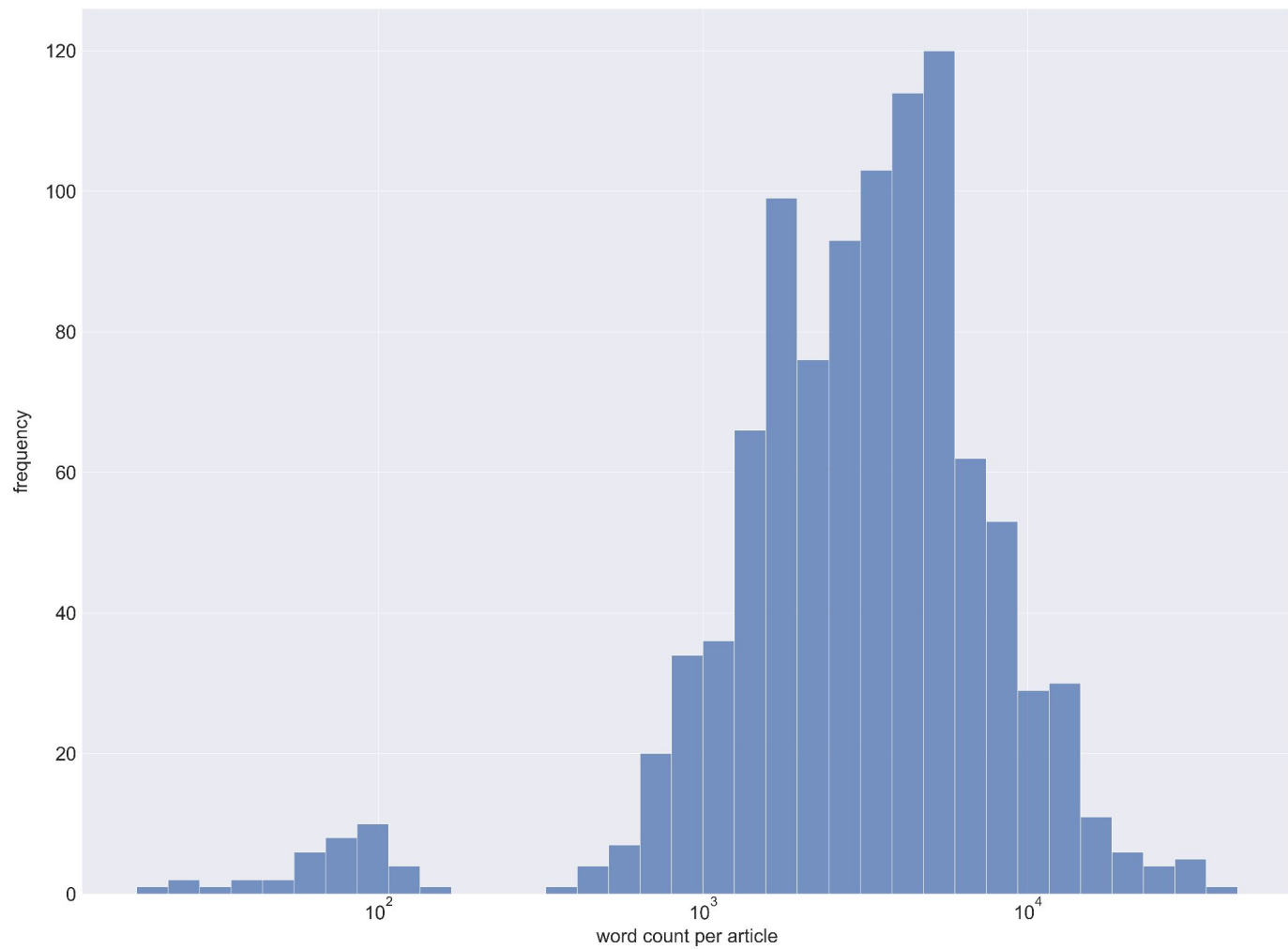Example word clouds for the statement variable, justification variable and statements in false class

# CT-FAN Dataset

The dataset consisting of news articles data, that was published in 2022. The dataset was collected from 2010 to 2022, including various topics related to elections, COVID-19 etc.

- in total **1264** articles
- attributes: id, title, text, rating (label)
- multiple labels: False (46%), Partially False (28%), True (17%) and Other (9%)
- news articles also in German language

# Word2vec embeddings with LSTM model

- applied all preprocessing techniques mentioned earlier

- clipped and padded length of articles to 3000

- 100-dim vectors

- NN model:

  - word2vec based embedding layer

  - LSTM layer with 32 units

  - dropout layers

  - Adam optimizer

# Word2vec embeddings with LSTM model

Obtained results for 10 epochs of training:

| Dataset | Accuracy | Balanced accuracy |
|---|---|---|
| LIAR-PLUS (6 classes): S | 0.19079 | 0.17397 |
| LIAR-PLUS (6 classes): S+J | 0.19205 | 0.18438 |
| CT-FAN (4 classes) | 0.45810 | 0.31610 |

# BERT - Bidirectional Encoder Representations from Transformers

- We used pretrained *BertForSequenceClassification* with *BertTokenizer* implemented in Transformers library in Python
- The architecture consisted of 124 layers, 1024 hidden layers, 16 heads, about 450 MB
- pretrained on uncased vocabulary
- the model was trained on Nvidia RTX 3060

Hyperparameters and other features:

- Maximum length of tokens - 512 (a popular length for pretrained transformers)
- Fixed seed for both dataloaders and the model
- **training-validation-test split proposed by the researchers**
- AdamW optimizer with **0.0001** learning rate with **StepLR scheduler**
- **10–20 epochs of training**

# BERT - Training on LIAR Dataset

As a point of reference, we also trained our BERT model on LIAR-PLUS dataset with using all preprocessing techniques. After 20 epochs, we obtained the following results:

- Accuracy: 0.61510 -> **0.62344**
- F1: 0.62789 -> **0.61682**

| Testing Set | Binary Classification | 6 Class Classification |
|---|---|---|
| LIAR-PLUS Test set | 77.2% | 37.4% |

*https://github.com/manideep2510/siamese-BERT-fake-news-detection-LIAR*
*https://aclanthology.org/W18-5513/*

# Multiclass Classification on LIAR-PLUS Dataset

| Scenario | Balanced Accuracy | Weighted F1-Score | MAE | MSE |
|---|---|---|---|---|
| 80-20 training on Valid: S+J | 0.78586 | 0.79560 | 0.43828 | 1.29297 |
| 80-20 training on Test: S+J | 0.84140 | 0.83385 | 0.36709 | 1.11076 |

# Statement versus Statement & Justification

| Scenario | Balanced Accuracy | Weighted F1-Score | MAE | MSE |
|---|---|---|---|---|
| Train: S | 0.93465 | 0.93677 | 0.11748 | 0.28736 |
| Train: S+J | 0.91823 | 0.92099 | 0.14016 | 0.34311 |
| Validation: S | 0.23441 | 0.24383 | 1.68047 | 4.84922 |
| Validation: S+J | **0.24940** | **0.25722** | 1.61875 | 4.72031 |
| Test: S | **0.23825** | **0.23796** | 1.71851 | 5.05680 |
| Test: S+J | 0.20855 | 0.21413 | 1.79763 | 5.44256 |

# Multiclass Classification on CT-FAN Dataset

| Scenario | Balanced Accuracy | Weighted F1-Score | MAE | MSE |
|---|---|---|---|---|
| Train Dataset | 0.83937 | 0.88907 | 0.15348 | 0.26899 |
| Test Dataset | **0.40765** | **0.53178** | 0.83333 | 1.64103 |

| | Team | True | False | Partially False | Other | Accuracy | Macro-F1 |
|---|---|---|---|---|---|---|---|
| 1 | iCompass [37] | 0.383 | 0.721 | 0.173 | 0.080 | 0.547 | 0.339 |
| 2 | NLP&IR@UNED [38] | 0.446 | 0.729 | 0.097 | 0.057 | 0.541 | 0.332 |
| 3 | Awakened [41] | 0.328 | 0.744 | 0.185 | 0.035 | 0.531 | 0.323 |
| 4 | UNED | 0.346 | 0.725 | 0.191 | 0.000 | 0.544 | 0.315 |

*Overview of the CLEF-2022 CheckThat! Lab: Task 3 on Fake News Detection*
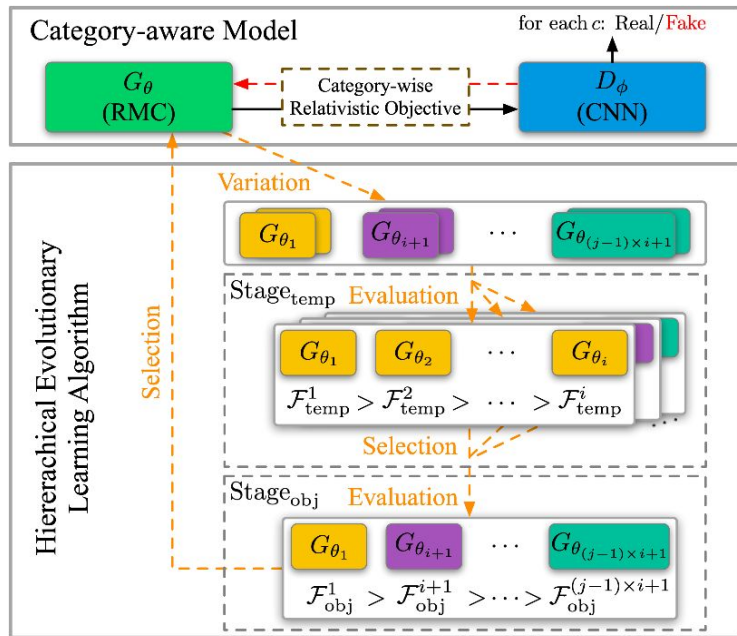
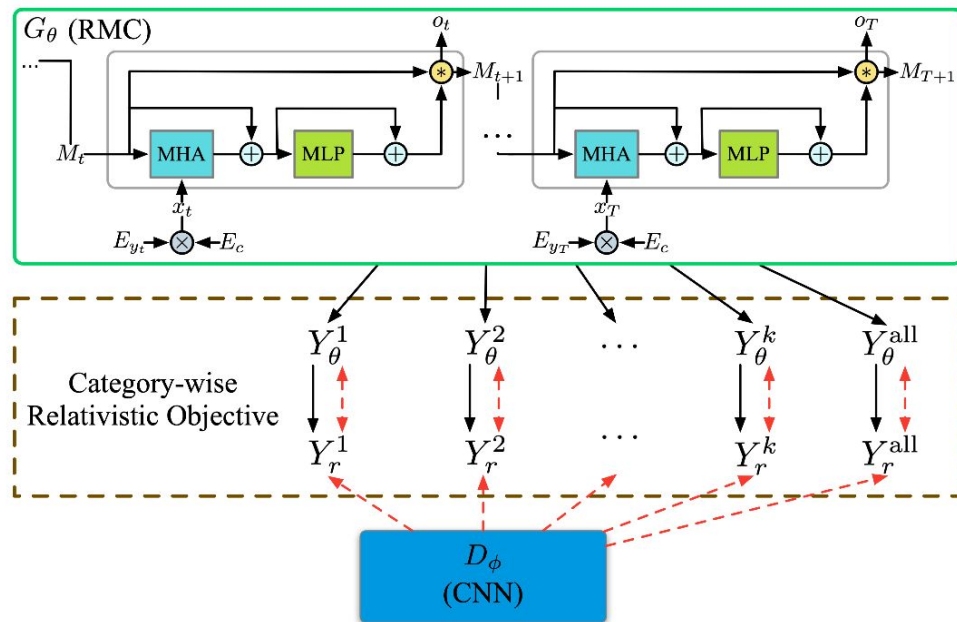# SentiGAN

# SentiGAN - Train Metrics (LIAR)


sentiGAN train accuracy without preprocessing


sentiGAN train accuracy with preprocessing

# SentiGAN - BLEU Score (LIAR)

# SentiGAN - Results (LIAR)

| Label | barely-true | false | half-true | mostly-true | pants-fire | true | AVG |
|---|---|---|---|---|---|---|---|
| Accuracy (no preprocessing) | 79.28 | 12.55 | 2.07 | 0.31 | 24.32 | 15.72 | **22.37** |
| Accuracy (with preprocessing) | 0 | 11.25 | 75.26 | 0 | 0 | 0 | **14.4** |

# CatGAN



(a) CatGAN

(b) Category-aware Model

# CatGAN - Train Metrics (LIAR)



catGAN train accuracy wihtout preprocessing



catGAN train accuracy with preprocessing

# CatGAN - BLEU Score (LIAR)

# CatGAN - Results (LIAR)

| Label | barely-true | false | half-true | mostly-true | pants-fire | true | AVG |
|---|---|---|---|---|---|---|---|
| Accuracy (no preprocessing) | 8.89 | 21.12 | 35.96 | 18.11 | 0 | 19.11 | **17.19** |
| Accuracy (with preprocessing) | 0 | 2.86 | 96.63 | 0 | 0 | 0 | **16.58** |

# Potential Improvements & Further Research

- Make the evaluation and testing **consistent** for all of the architectures used
- Applying **ensemble learning** with models of different complexity and word embeddings
- Using **multiple architectures for each of the inputs** (for example metadata, statement, justification in LIAR-PLUS)
- Training on more **powerful hardware** with **more epochs** and **training optimisation** (early stopping, different schedulers, and hyperparameters)

# Thank you for your attention!

Shall you have any questions, please do not hesitate to ask.