# Automatic Metrics in NLG:
# A Survey of Current Evaluation Practices

Patrícia Schmidtová, Saad Mahamood, Simone Balloccu, Ondřej Dušek, Albert Gatt, Dimitra Gkatzia, David M. Howcroft, Ondřej Plátek, Adarsa Sivaprasad

- Automatic metrics are quick proxies, but…

  - Some have a poor correlation with human judgment

  - Many cannot capture factuality or faithfulness issues in text

  - Different implementations make results hard to interpret and reproduce

  - They can be over-reported without adding any informational value

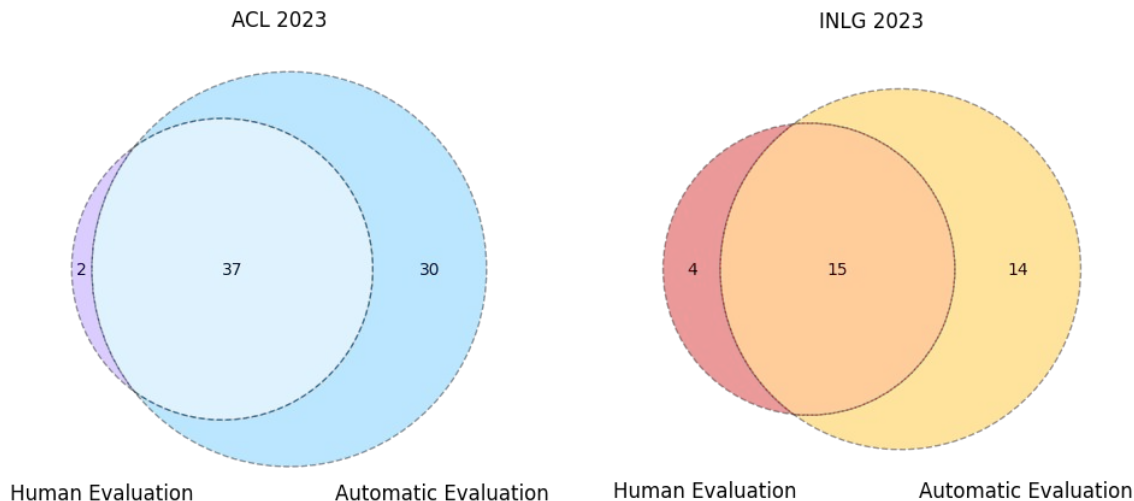# How are automatic metrics used in  NLG?

# Method

We collected papers from INLG 2023 and ACL 2023 Generation track and annotated the following information:

- **Name** of the evaluation method

- Was the method **newly introduced**?

- Which **task**(s) was this metric used to evaluate?

- Did the authors comment on any **correlation between automatic** and **human evaluation**?

- Did the authors provide **implementation details** for the metric?

- Was the metric **only** reported in the **Appendix**?

- Did the authors explain the **rationale** for the metric?

# Overview of Results

- **110** papers total (**36** from **INLG** and **74** from **ACL**)
- **102** papers **included** any **evaluation**
- **57%** use **human** evaluation
- **94%** use **automatic** evaluation
- **51%** use **both**
- **634 counts** of automatic metrics (**283 unique**)

ACL 2023

INLG 2023



Human Evaluation    2    37    30    Automatic Evaluation

Human Evaluation    4    15    14    Automatic Evaluation

# Metric Families & Categories

| Metric Task Name | INLG | ACL | Total |
|---|---|---|---|
| Overlap | 71 | 201 | 272 |
| Semantic Similarity | 20 | 59 | 79 |
| Match | 15 | 61 | 76 |
| Text Properties | 12 | 63 | 75 |
| Text Classifier | 17 | 57 | 74 |
| Factuality | 49 | 21 | 70 |
| Perplexity | 3 | 37 | 40 |
| Distance-based | 1 | 15 | 16 |
| Combination | 0 | 14 | 14 |
| Inference Speed | 0 | 4 | 4 |

| Metric Family Name | INLG | ACL | Total |
|---|---|---|---|
| BLEU | 26 | 69 | 95 |
| ROUGE | 27 | 65 | 92 |
| N-gram diversity | 6 | 49 | 55 |
| Style Classifier | 5 | 37 | 42 |
| BERTScore | 8 | 32 | 40 |
| Perplexity | 3 | 29 | 32 |
| METEOR | 6 | 21 | 27 |
| Semantic Similarity | 9 | 12 | 21 |
| Overlap | 6 | 21 | 27 |
| Factuality | 5 | 13 | 18 |
| Accuracy | 8 | 8 | 16 |
| Quality Estimation | 7 | 7 | 14 |
| Combination | 0 | 14 | 14 |
| BARTScore | 2 | 10 | 12 |

. . .

| | INLG | ACL | Total |
|---|---|---|---|
| Recall | 2 | 44 | 6 |
| Edit Distance | 1 | 5 | 6 |
| Flesch Readability | 1 | 3 | 4 |
| Inference Speed | 0 | 4 | 4 |
| Precision | 1 | 2 | 3 |
| loss/error | 0 | 3 | 3 |
| chrF++ | 1 | 1 | 2 |

# What kinds of metrics were used?



Metric family use per venue

# Correlation with Human Evaluation

# Results per Task

# The Lack of Implementation Details

# Recommendations - Evaluation Quality

- Rationalize your selection of metrics

- Comment on metric combinations

- Do not copy-paste widely used metrics

- Respect the intended use of metrics

- Discuss (dis)agreements between human and automatic evaluation

# Recommendations - Evaluation Reproducibility

- Share evaluation details

- Share data samples

- Release code

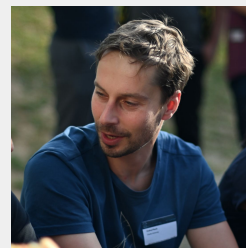# Thank You! Questions?

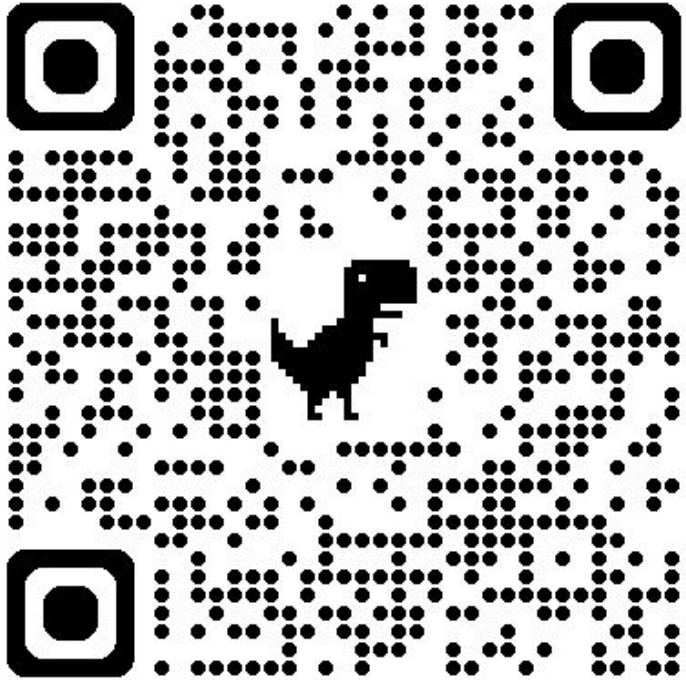Patricia  Saad  Simone  Ondřej  Albert

Dimitra  Dave  Ondřej  Adarsa

Correspondence to: schmidtova@ufal.mff.cuni.cz
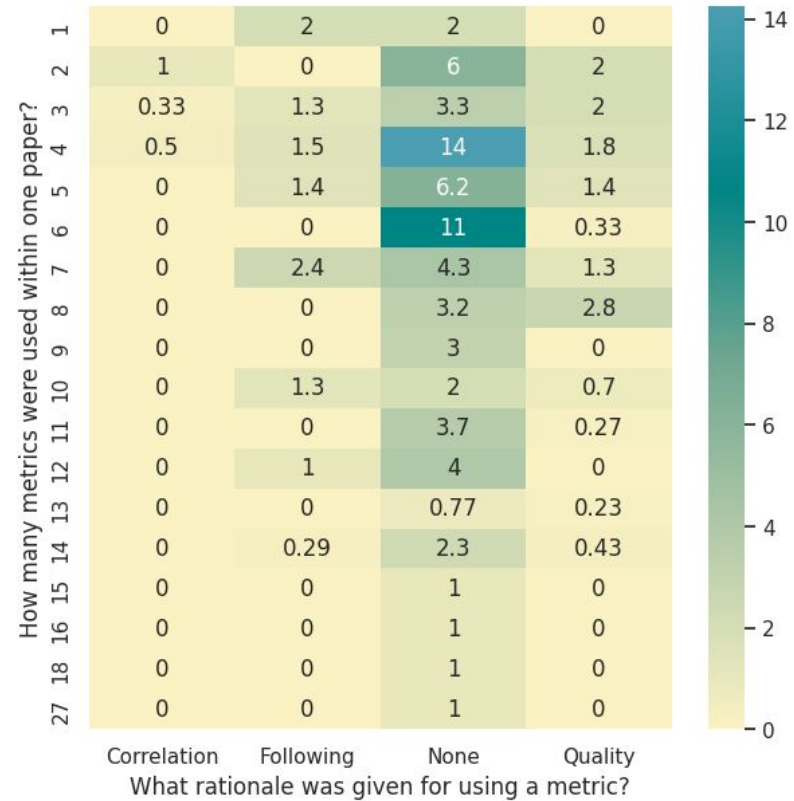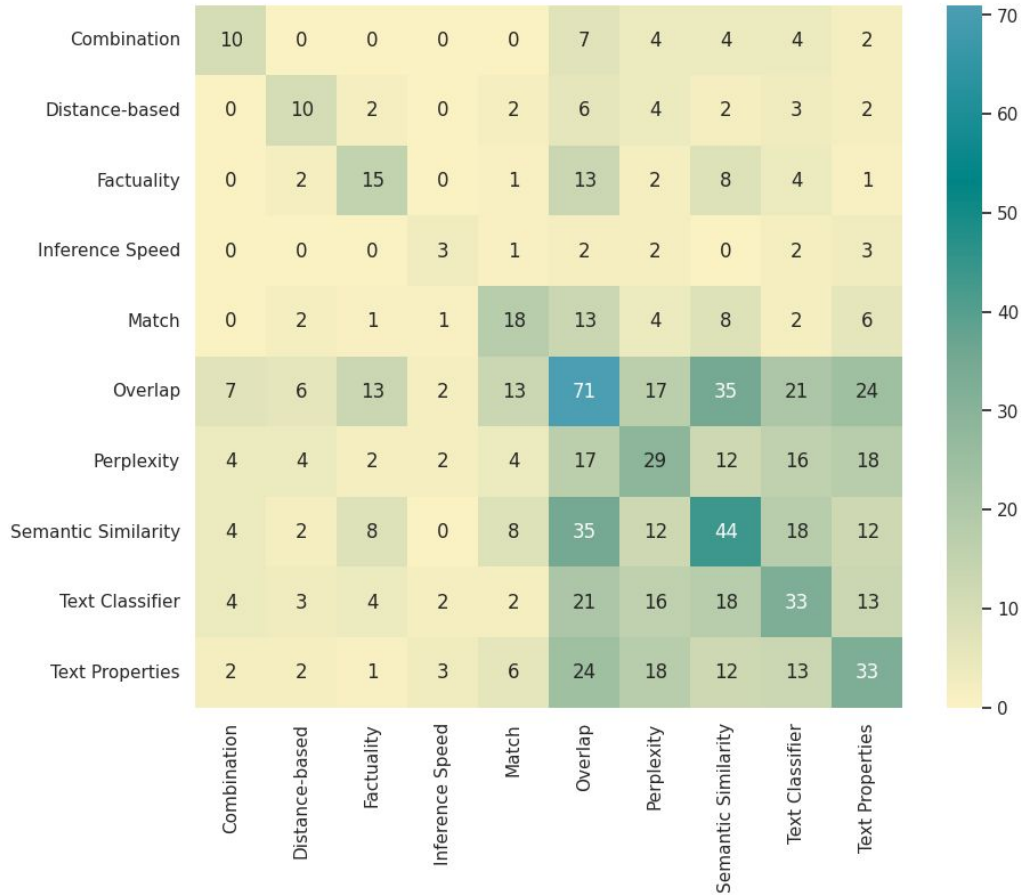
European Research Council
Established by the European Commission
erc

**Link to the paper**



**Link to GitHub**

# Backup Slides

# Correlations

# Variants of BLEU and ROUGE