

The 4th International Conference on Natural Language and Speech Processing (ICNLSP 2023)

Kristina Schaaff, Tim Schlippe, Lorenz Mindner

CLASSIFICATION OF HUMAN- AND AI-GENERATED TEXTS FOR ENGLISH, FRENCH, GERMAN, AND SPANISH

Virtual
December 16, 2023

AGENDA

Motivation

1

Human-AI-Generated Text Corpus

2

**Features for the Classification
of Human- and AI-Generated Texts**

3

Human-AI-Generated Text Corpus

4

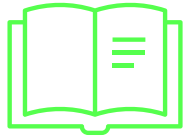
Conclusion and Future Work

5

1

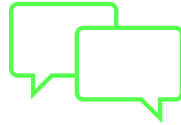
MOTIVATION

TYPICAL USE CASES FOR CHATBOTS



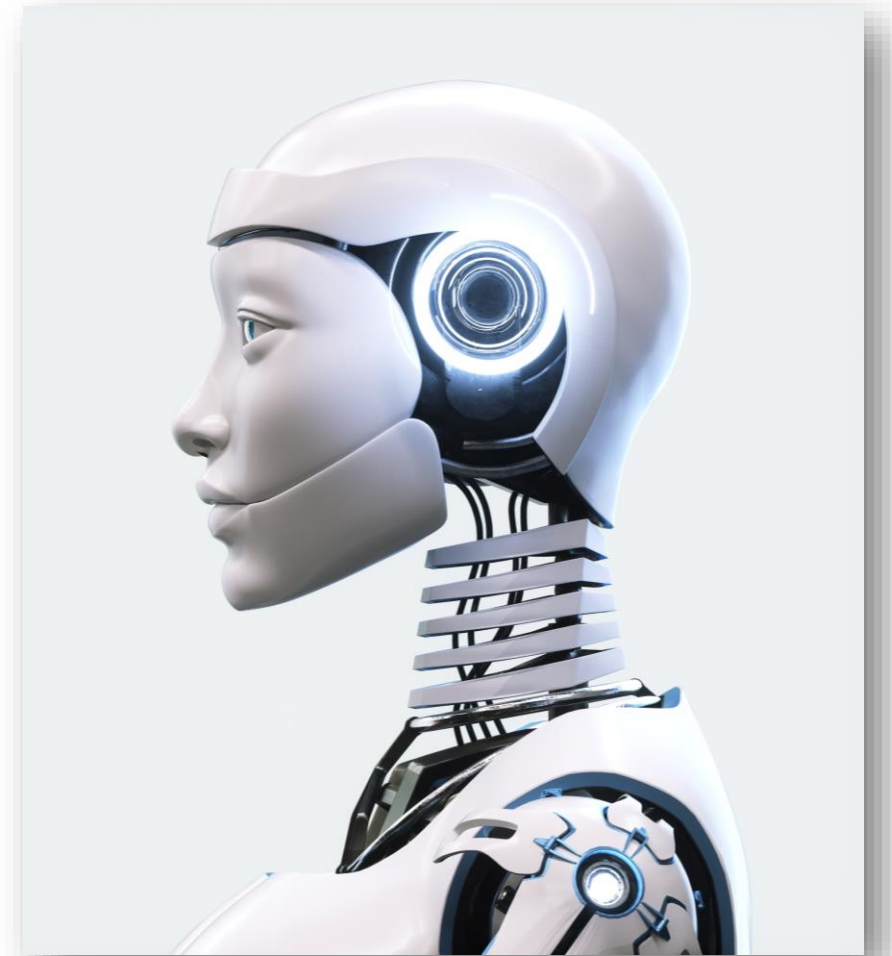
Learning
Support

Chat
Companion

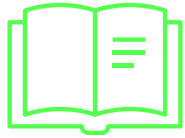


Medical
Advice

Customer
Support

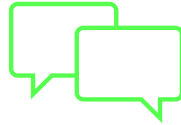


DANGERS



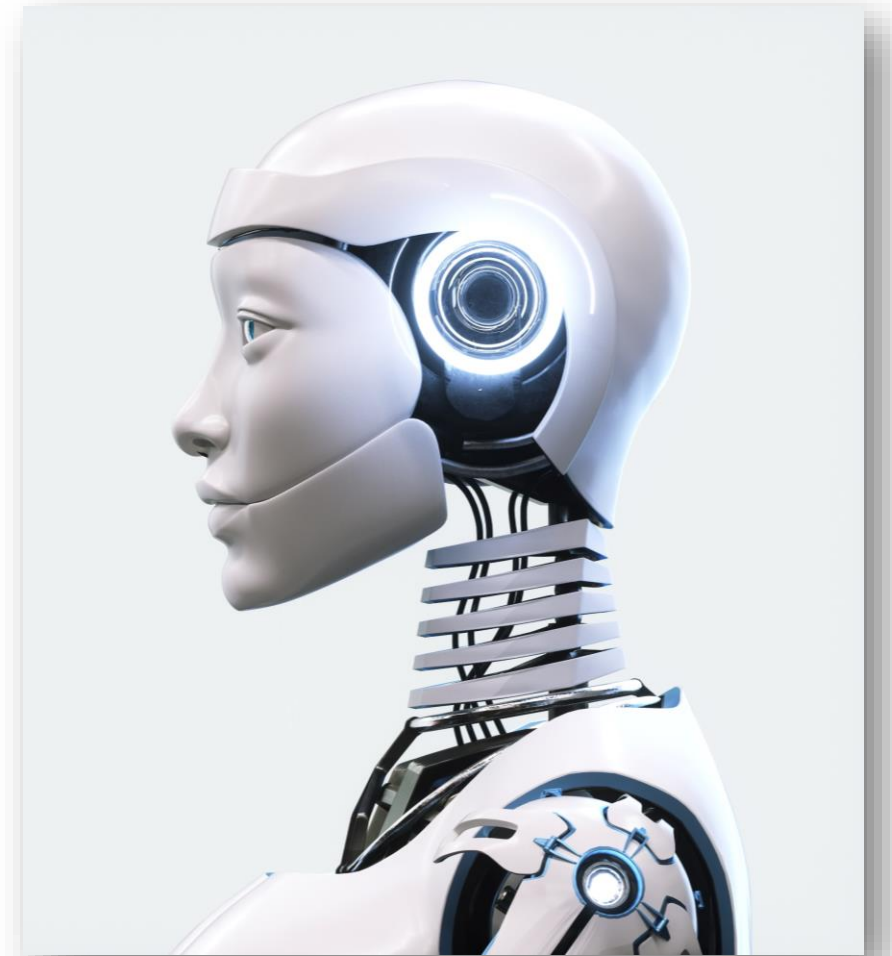
Fake News

Spamming



Plagiarism

...



RQ1

FEATURES

How effective are different features in distinguishing between human- and AI-generated texts from the educational domain?

RQ2

PROMPTS

What impact do we have if we specifically tell the AI to generate text as a human would do it?

RQ3

LANGUAGE

How do the classification results vary across different languages?

RQ4

TEXT DOMAIN

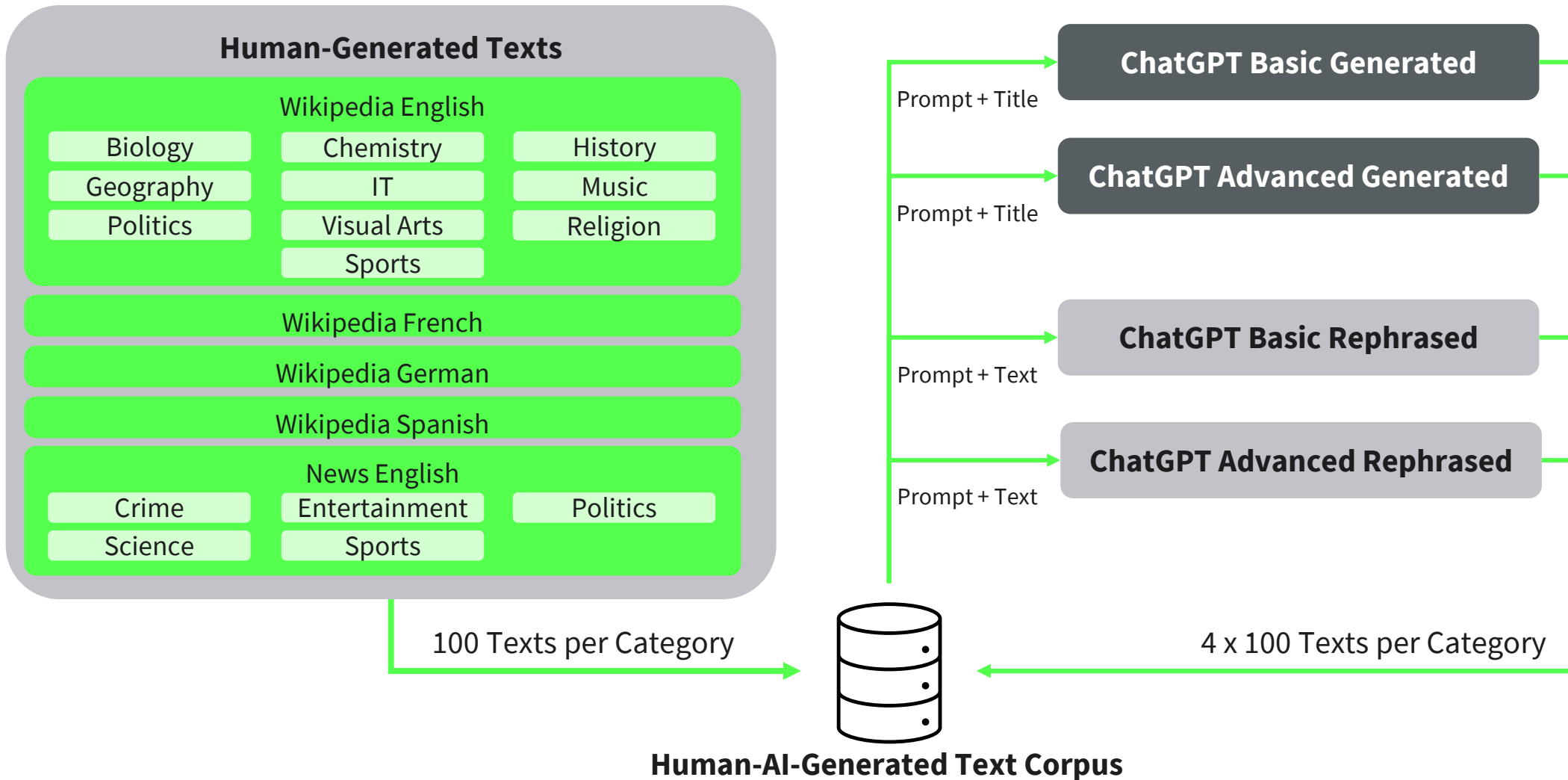
Do the features perform comparably well in the news domain?

2

HUMAN-AI-GENERATED TEXT CORPUS

<https://github.com/LorenzM97/human-AIgeneratedTextCorpus>

HUMAN-AI-GENERATED TEXT CORPUS



HUMAN-AI-GENERATED TEXT CORPUS

Basic Generated

Prompt

Generate a text on the following topic: Australia

Advanced Generated

Prompt

Generate a text on the following topic in a way a human would do it: Australia

Basic Rephrased

Prompt

Rephrase the following text: Australia, officially the Commonwealth of Australia, is a sovereign country [...].

Advanced Rephrased

Prompt

Rephrase the following text in a way a human would do it:
Australia, officially the Commonwealth of Australia, is a sovereign country [...].

3

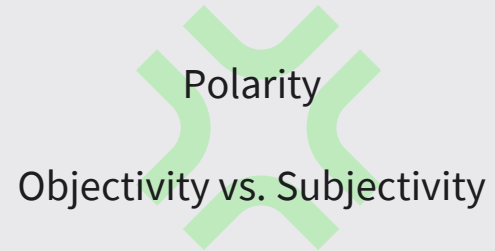
FEATURES FOR THE CLASSIFICATION OF HUMAN- AND AI-GENERATED TEXTS

FEATURES

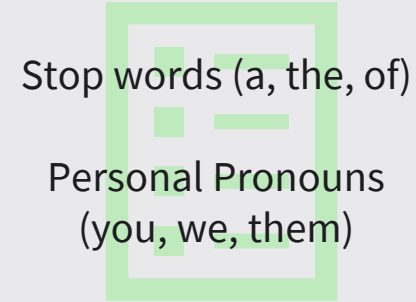
Perplexity Features



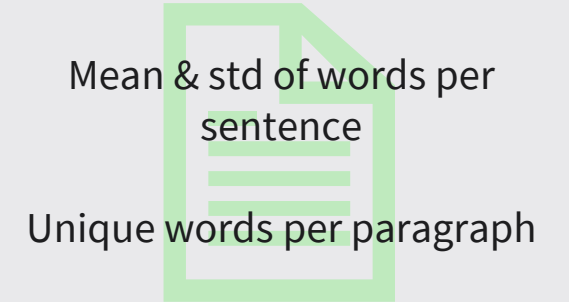
Semantic Features



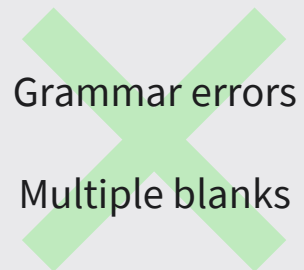
List Lookup Features



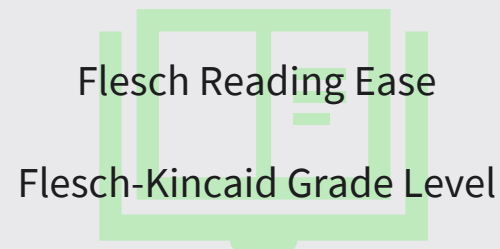
Document Features



Error-Based Features



Readability Features



AI Feedback Features

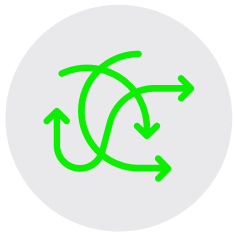


Text Vector Features



FEATURES

8 feature categories, 37 features



Perplexity
Features



Semantic
Features



List Lookup
Features



Document
Features



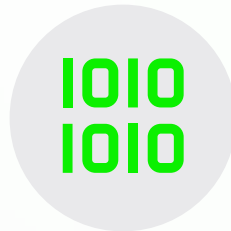
Error-Based
Features



Readability
Features



AI Feedback
Features



Text Vector
Features

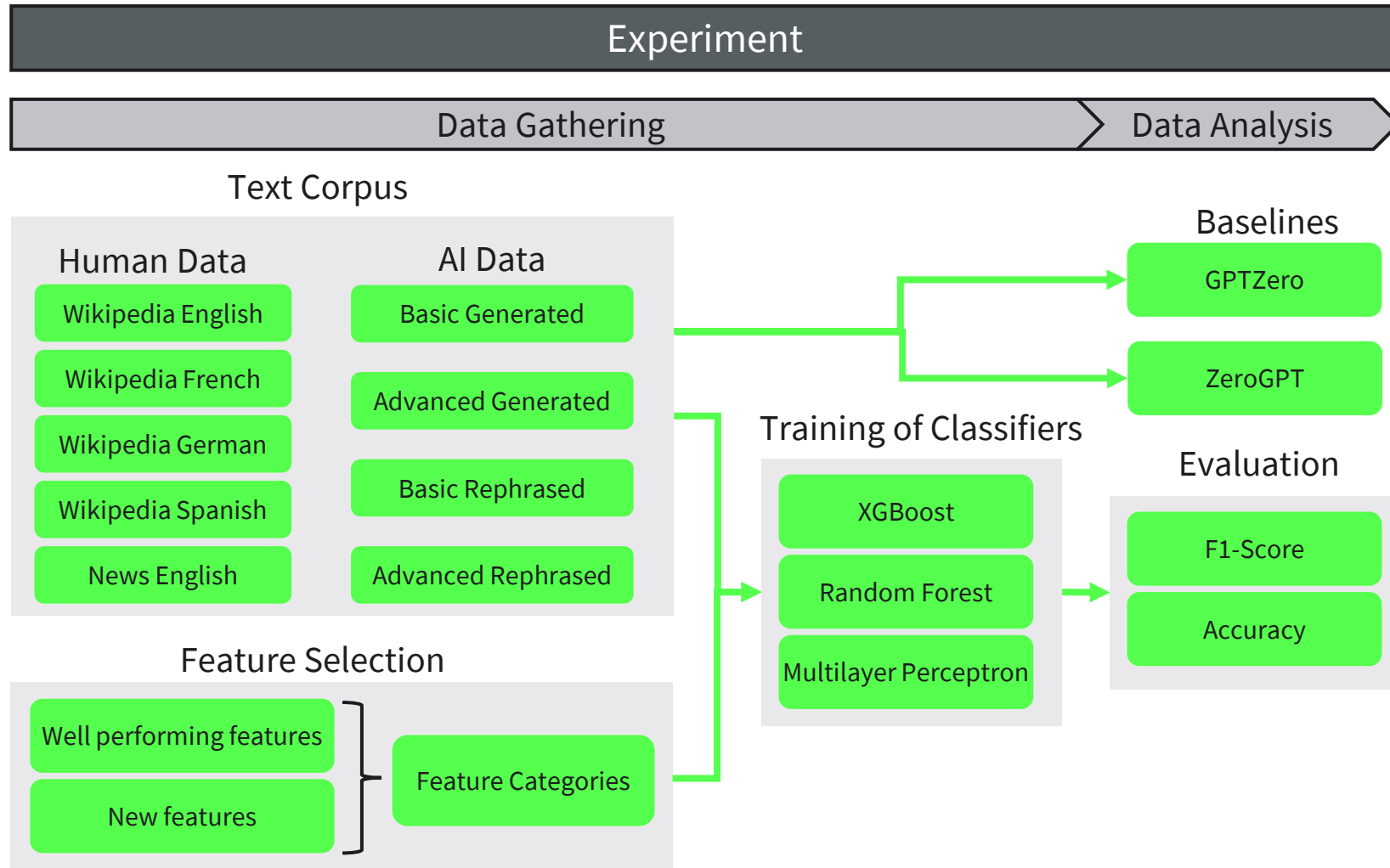
Category	Feature	Description	Reference
Perplexity	PPL_{mean}	mean PPL	[21][14][19]
	PPL_{max}	maximum PPL	[21][14][19]
Semantic	$sentiment_{polarity}$	degree of positivity/negativity [-1,+1]	[14][19]
	$sentiment_{subjectivity}$	degree of subjectivity [0,+1]	new
List Lookup	$stopWord_{count}$	number of stop words	[17]
	$specialChar_{count}$	number of special characters	[28]
	$discourseMarker_{count}$	number of discourse markers	new
	$titleRepetition_{count}$	absolute repetitions of title	new
	$titleRepetition_{relative}$	relative repetitions of title	new
Document	$wordsPerParagraph_{mean}$	\emptyset number of words per paragraph	[28]
	$wordsPerParagraph_{stdev}$	stdev of $wordsPerParagraph$	[28]
	$sentencesPerParagraph_{mean}$	\emptyset number of sentences per paragraph	[28]
	$sentencesPerParagraph_{stdev}$	stdev of $sentencesPerParagraph$	[28]
	$wordsPerSentence_{mean}$	\emptyset number of words per sentence	[28]
	$wordsPerSentence_{stdev}$	stdev of $wordsPerSentence$	[28]
	$uniqWordsPerSentence_{mean}$	\emptyset number of unique words per sentence	[17]
	$uniqWordsPerSentence_{stdev}$	stdev of $uniqWordsPerSentence$	new
	$words_{count}$	number of running words	[19][17][28]
	$uniqWords_{count}$	number of unique words	[28]
	$uniqWords_{relative}$	relative number of unique words	[28]
	$paragraph_{count}$	number of paragraphs	[28]
	$sentence_{count}$	number of sentences	[28]
	$punctuation_{count}$	number of punctuation marks	[28]
	$quotation_{count}$	number of quotation marks	new
$character_{count}$	number of characters	[28]	
$uppercaseWords_{relative}$	relative number of words in uppercase	[17]	
$personalPronoun_{count}$	absolute number of personal pronouns	[14]	
$personalPronoun_{relative}$	relative number of personal pronouns	[14]	
$POSPerSentence_{mean}$	\emptyset number of unique POS-tags/sentence	[19][28][18]	
Error-Based	$grammarError_{count}$	number of spelling/grammar errors	new
	$multiBlank_{count}$	number of multiple blanks	new
Readability	$fleschReadingEase$	Flesch Reading Ease score [0-100]	[17][29]
	$fleschKincaidGradeLevel$	Readability as U.S. grade level [0-100]	[17][30]
AI Feedback	$AIFeedback$	Ask AI if text was generated by AI	new
Text Vector	$TF-IDF$	500-dim TF-IDF vector of 1-/2-grams	[17][31]
	$Sentence-BERT$	\emptyset Sentence-BERT vector	[32]
	$Sentence-BERT-dist$	\emptyset distance of Sentence-BERT vectors	new

Table 3: Summary of our Features for the Classification of Generated Texts.

4

EXPERIMENTS AND RESULTS

EXPERIMENTAL SETUP



Training for:

- Human vs. basic AI-generated
- Human vs. advanced AI-generated
- Human vs. basic AI-rephrased
- Human vs. advanced AI-rephrased

RESULTS

RQ1: FEATURES

Text Corpus
RQ1
“How effective are different features in distinguishing between human- and AI-generated texts from the educational domain?”
Human Data
AI Data
Wikipedia English
Basic Generated
Basic Rephrased

Comparison of Feature Categories



F1-Scores Basic Generated



1	94.9% TextVector _{traditional+new}	Readability _{traditional} 59.3%	1
1	94.9% Document _{traditional}	ErrorBased _{new} 63.9%	2
3	85.3% Perplexity _{traditional}	AIFeedback _{new} 68.1%	3



F1-Scores Basic Rephrased



1	79.7% Document _{traditional}	AIFeedback _{new} 50.9%	1
2	78.2% TextVector _{traditional+new}	Readability _{traditional} 51.1%	2
3	73.9% ListLookup _{traditional+new}	Semantic _{traditional+new} 64.4%	3

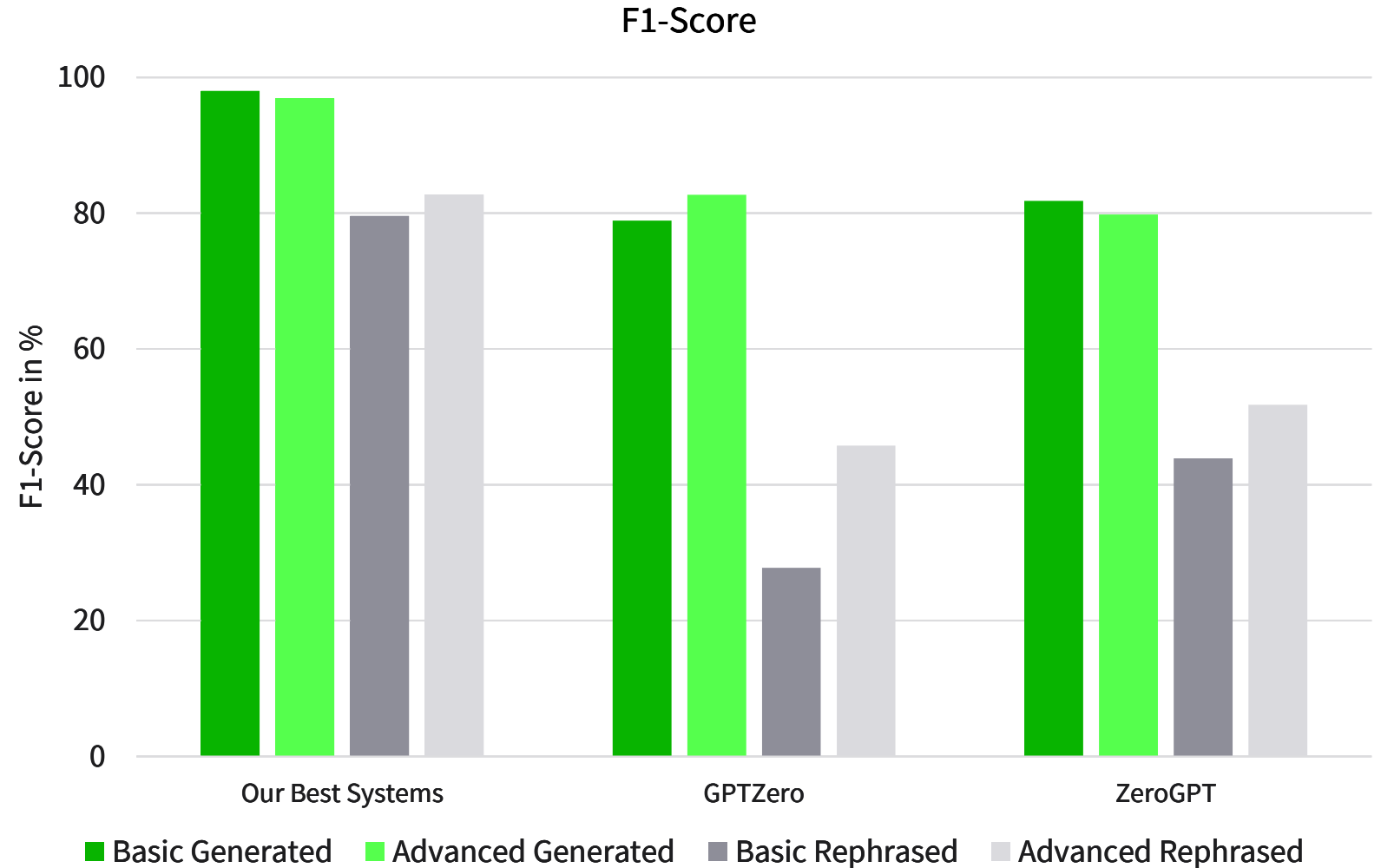
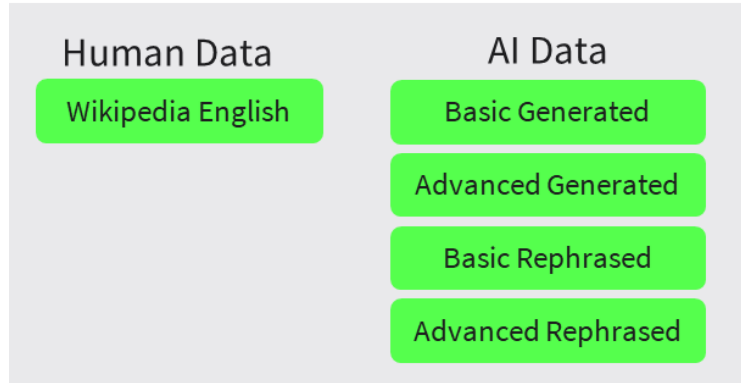
RESULTS

RQ2: PROMPT

RQ2

“What impact do we have if we specifically tell the AI to generate text as a human would do it?”

Text Corpus



RESULTS

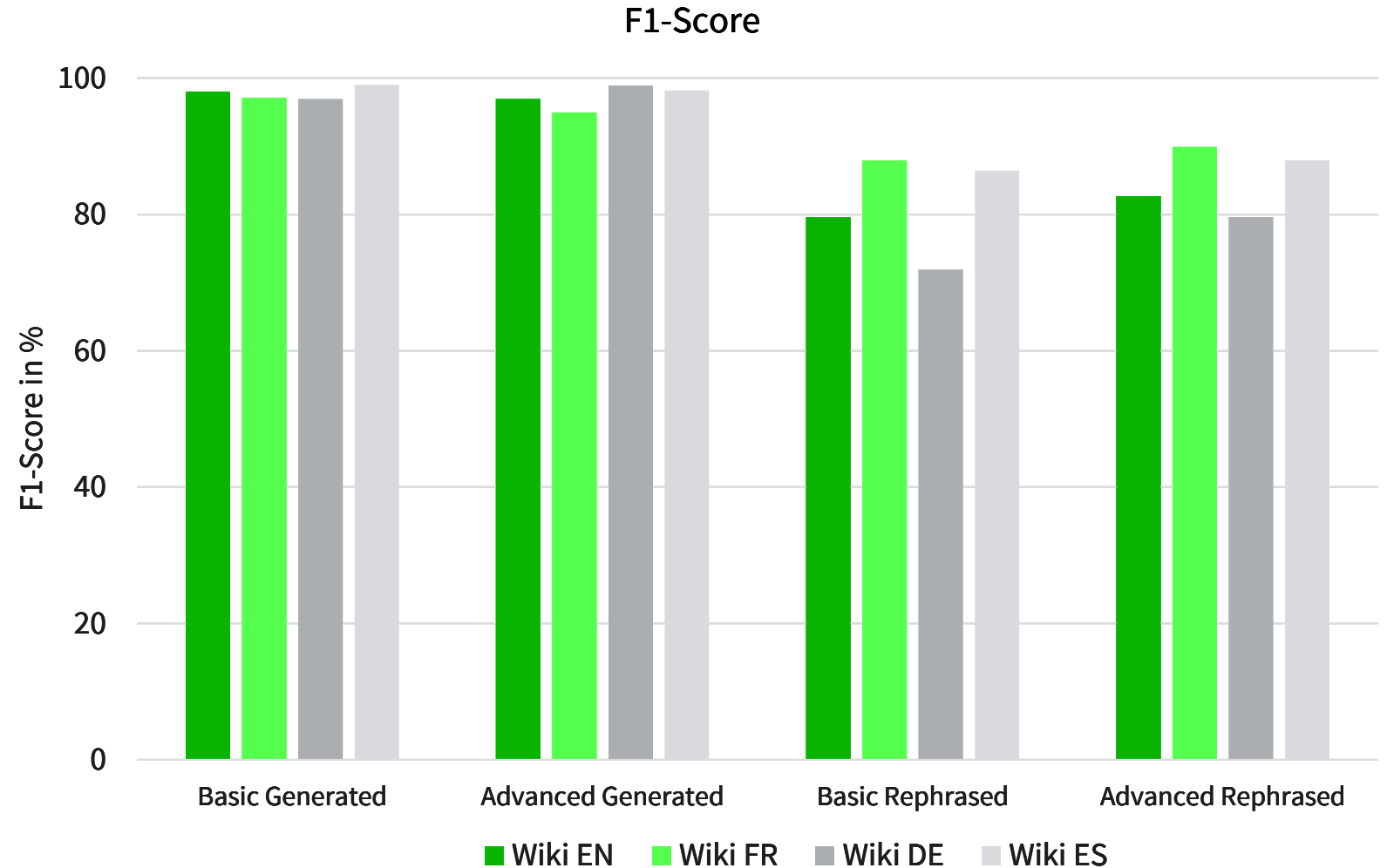
RQ3: LANGUAGE

RQ3

“How do the classification results vary across different languages?”

Text Corpus

Human Data	AI Data
Wikipedia English	Basic Generated
Wikipedia French	Advanced Generated
Wikipedia German	Basic Rephrased
Wikipedia Spanish	Advanced Rephrased

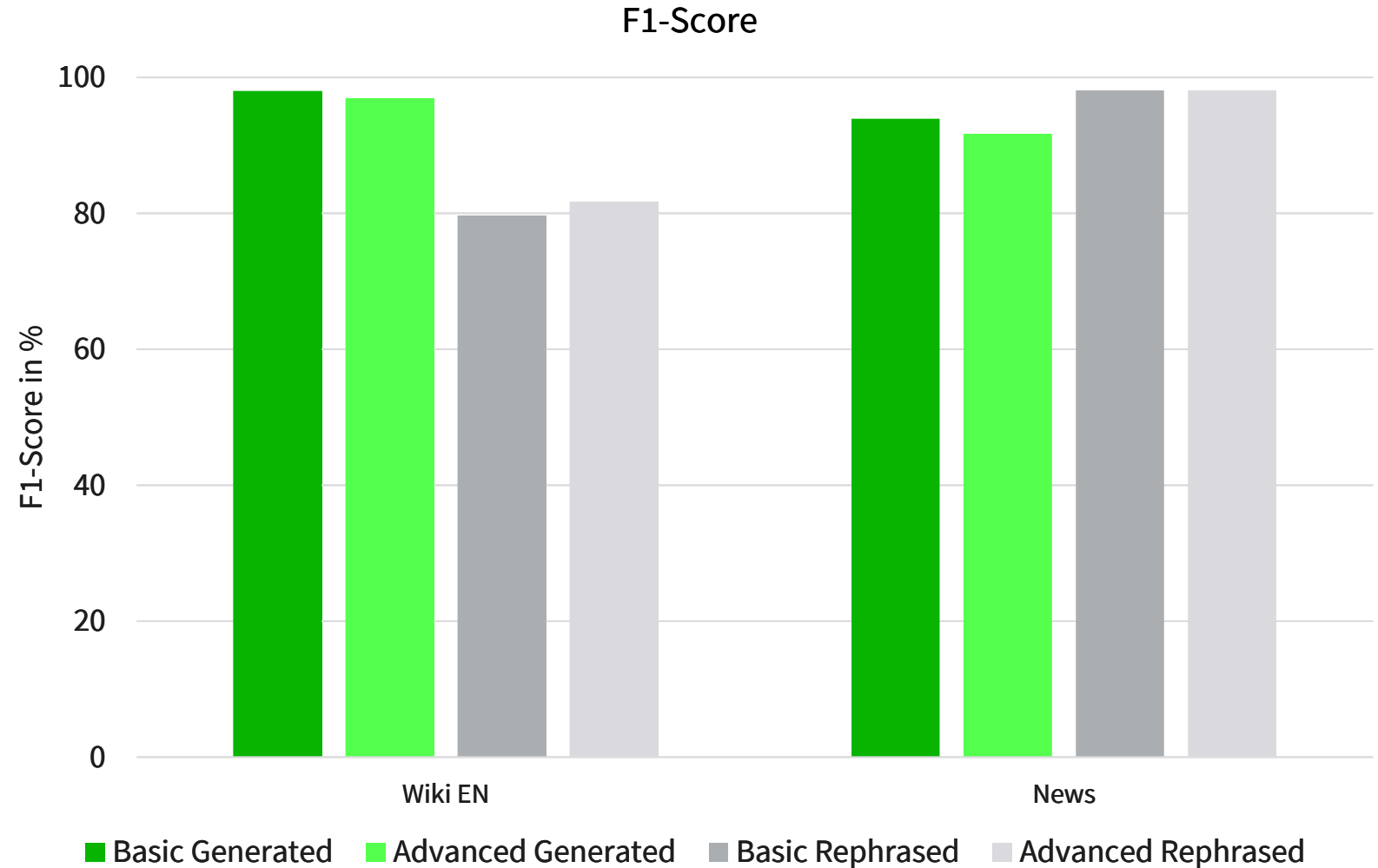
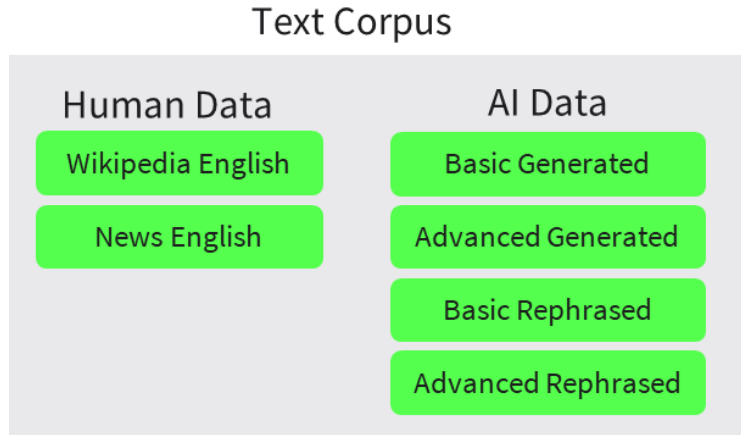


RESULTS

RQ4: TEXT DOMAIN

RQ4

“Do the features perform comparably well in the news domain?”



COMPARISON TO BASELINES

F1-SCORES OF BEST SYSTEMS

Wikipedia English

	GPTZero	ZeroGPT
Basic Generated	+24.2%	+19.8%
Basic Rephrased	+186.7%	+81.5%
Advanced Generated	+17.2%	+21.4%
Advanced Rephrased	+78.4%	+57.7%

News English

	GPTZero	ZeroGPT
Basic Generated	+25.9%	+21.0%
Basic Rephrased	+145.2%	+98.2%
Advanced Generated	+22.9%	+18.1%
Advanced Rephrased	+181.9%	+83.0%

French German Spanish

	ZeroGPT	ZeroGPT	ZeroGPT
Basic Generated	+33.7%	+36.7%	+38.5%
Basic Rephrased	+30.4%	+45.3%	+35.6%
Advanced Generated	+28.9%	+51.2%	+37.3%
Advanced Rephrased	+37.7%	+61.5%	+40.6%

5

CONCLUSION AND FUTURE WORK

CONCLUSION & FUTURE WORK



Generated > Rephrased

Minor influence of prompt

Comparable Performance Across Languages

F1-scores for AI-generated texts between **97%** and **99%**

F1-scores for AI-rephrased texts between **78%** and **89%**



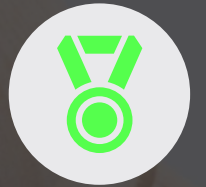
Educational Domain < News Domain

Comparable performance for AI-generated texts

Better results for AI-rephrased texts

Baselines < Our Systems

Our best basic text rephrasing detection system performs almost twice as good



Future Work



Improvement of text generation

Investigation of other prompts and text domains

Comparison to transformer-based models

THANK YOU

Tim Schlippe

 tim.schlippe@iu.org