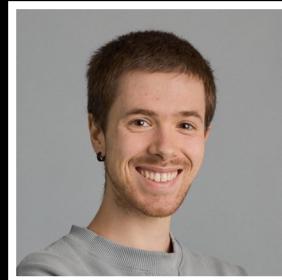


Homoglyph-Based Attacks: Circumventing LLM Detectors

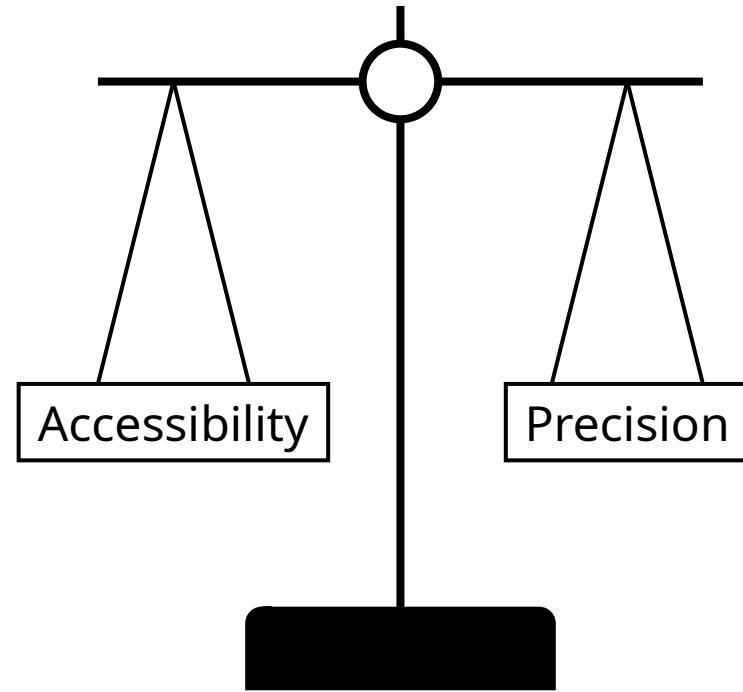


Aldan Creo

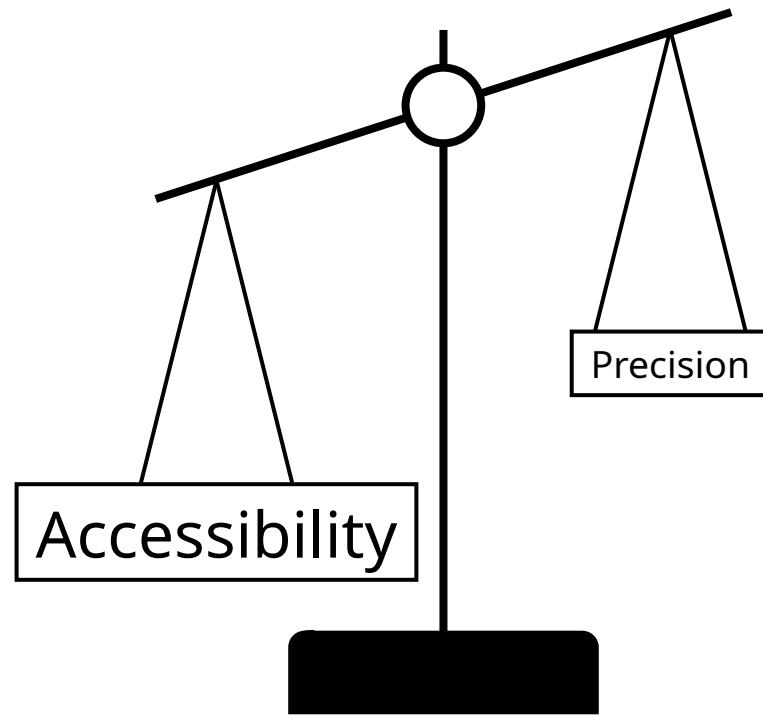
Technology Research Specialist

Some notes...

1

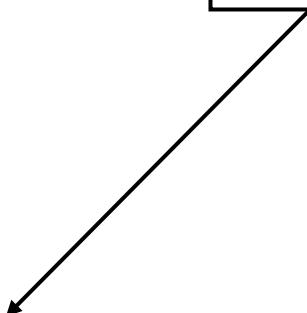


1



2

Precision



(Reference)

3

Takeaway

Bingo cards – explanation
later

Let's get
started!

Who am I?

research



The 31st International
Conference on Computational
Linguistics

Evading AI- Generated Content Detectors using Homoglyphs

Aldan Creo, Shushanta Pudasaini

Structure

What are homoglyphs?

AI-generated text detectors

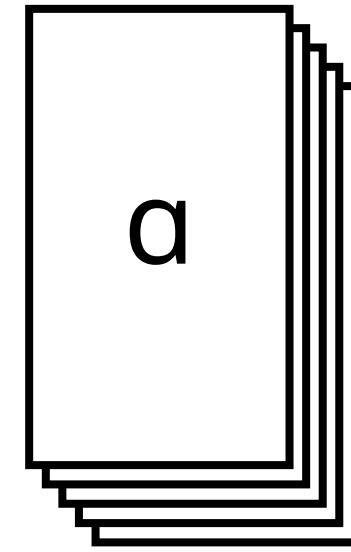
Homoglyph-based attacks

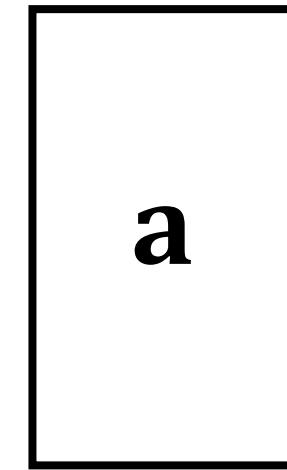
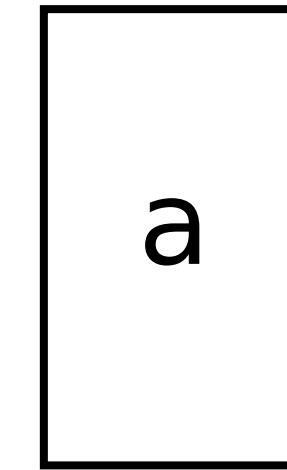
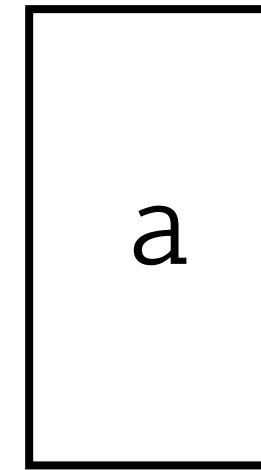
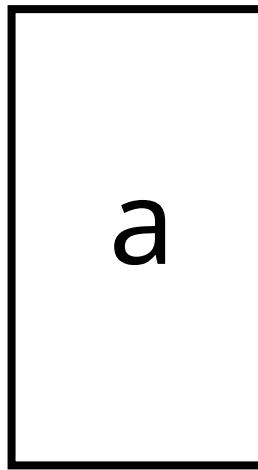
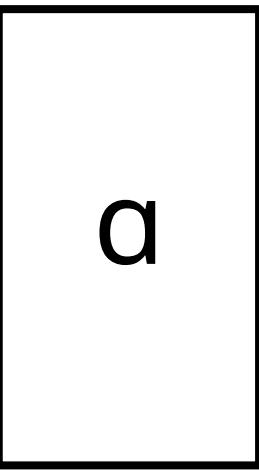
Technical analysis

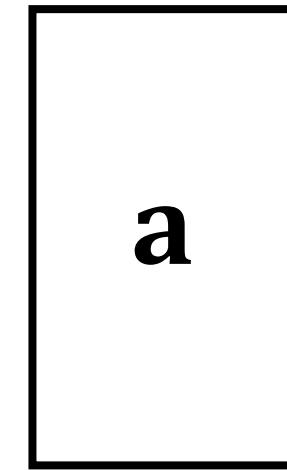
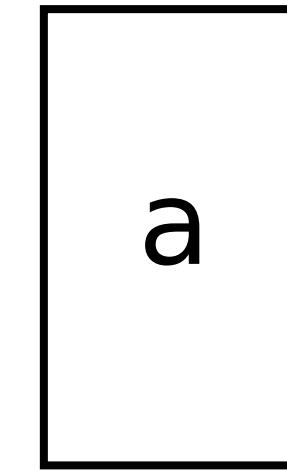
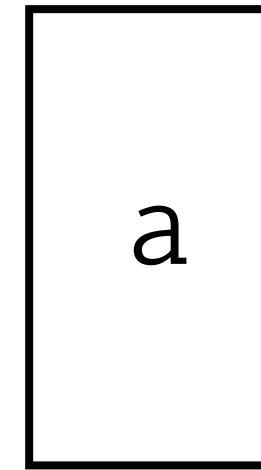
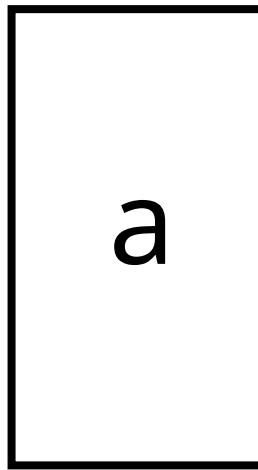
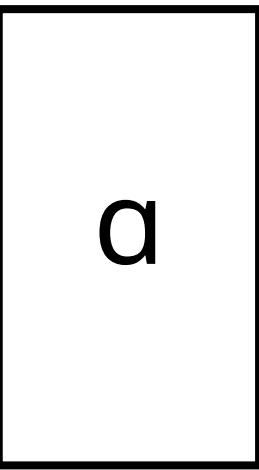
Effectiveness

Implications

What are homoglyphs?







a
0251
LATIN SMALL LETTER ALPHA

a
0061
LATIN SMALL LETTER A

a
1D68A
MATHEMATICAL
MONOSPACE
SMALL A

a
1D5BA
MATHEMATICAL
SF SMALL A

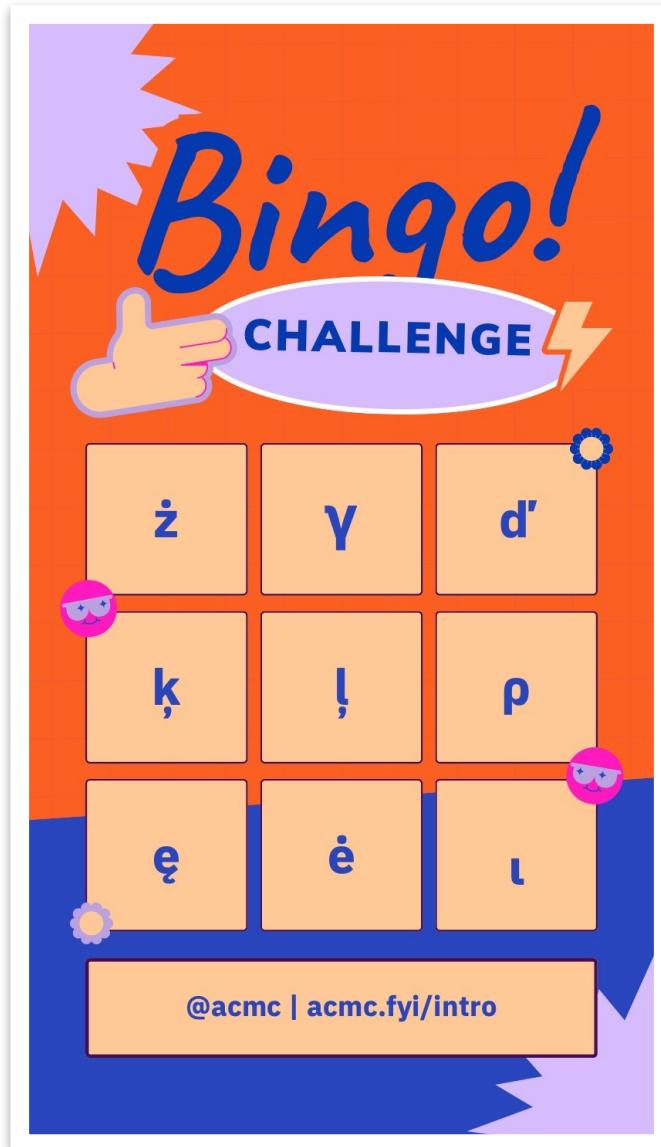
a
1D41A
MATHEMATICAL
BOLD SMALL A

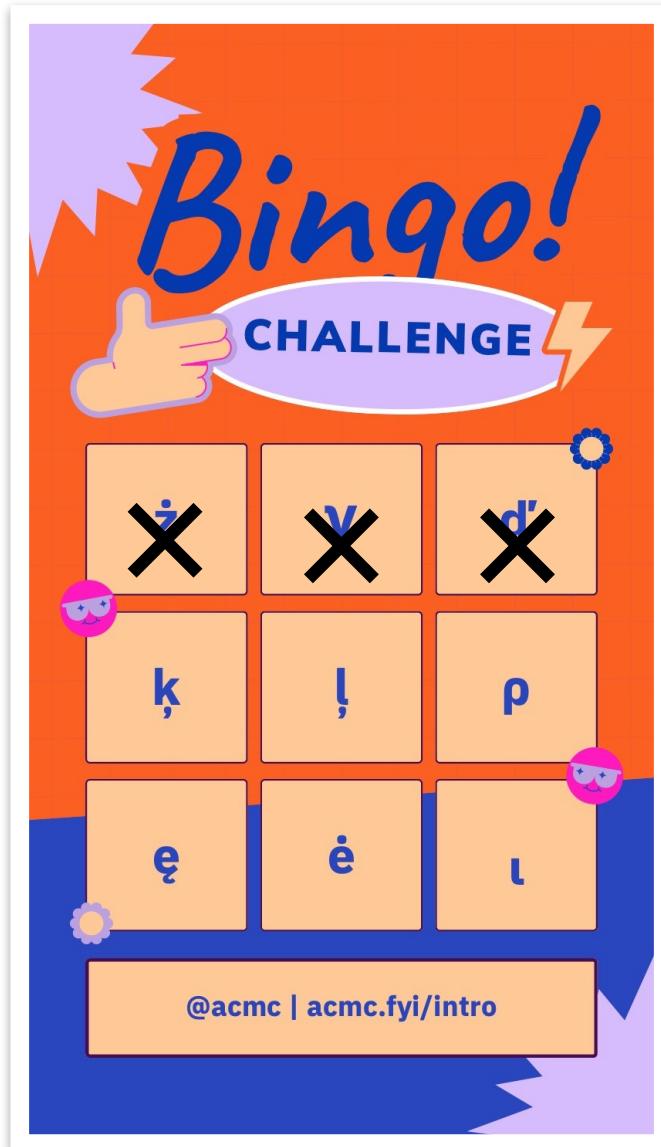
Homoglyphs

Characters that seem identical
but have different meanings

Homoglyphs

Characters that seem identical
but have different encodings





⚠ not all homoglyphs are in the
slides!

Past applications

URL attacks

URL attacks

Click here to reset your password: <https://paypal.com>

URL attacks

Click here to reset your password: <https://xn--paypa-n6a.com/>

Hiding information

Steganography

Hiding data in plain sight
to avoid detection

"If the first 'a' is an 'a', trigger the exploit"

"If the first 'a' is an 'a', trigger the exploit"

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi ultricies quam quis est tempus maximus. Praesent bibendum tortor nec neque sagittis, ut mollis turpis fringilla. Vestibulum velit orci, malesuada sit amet odio et, varius imperdiet eros. Interdum et malesuada fames ac ante ipsum primis in faucibus. Aenean finibus diam sed venenatis ullamcorper. Mauris hendrerit ligula est, nec varius sem hendrerit in. Vestibulum ut lobortis elit, at lobortis risus.

Homoglyphs: characters we can't discern

AI-generated text detectors

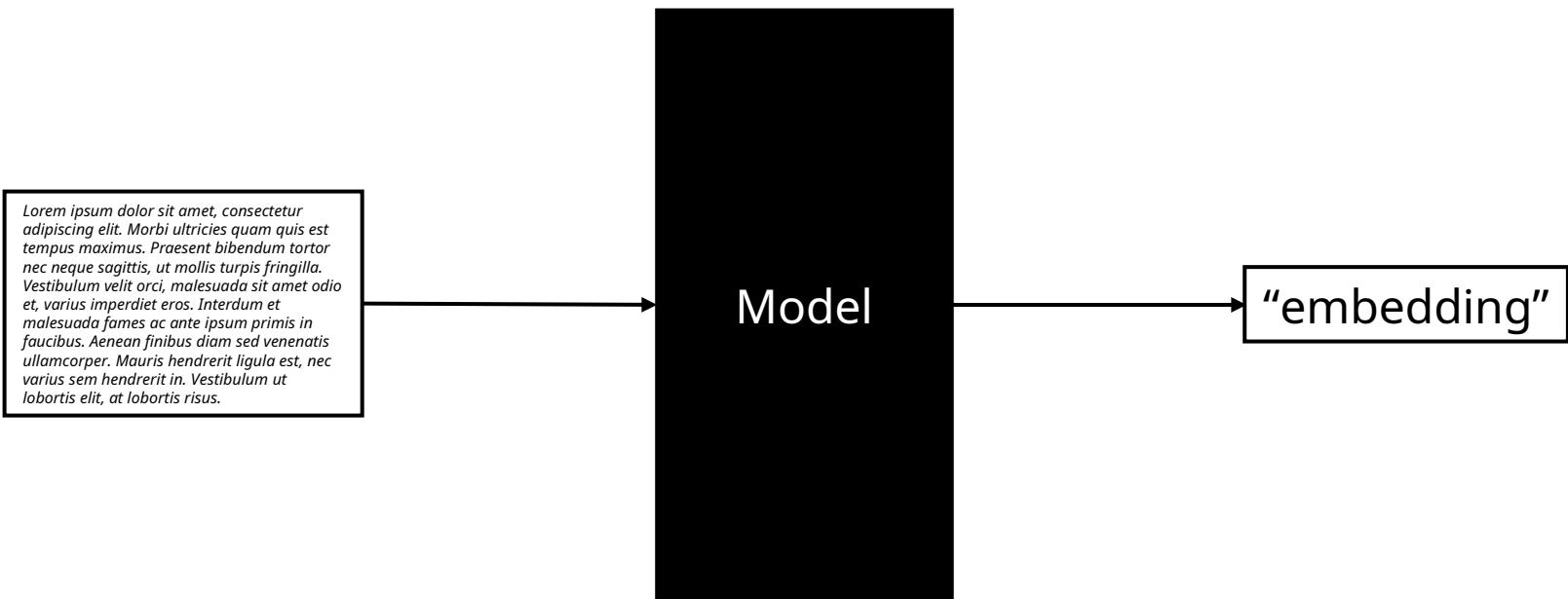
Classifiers

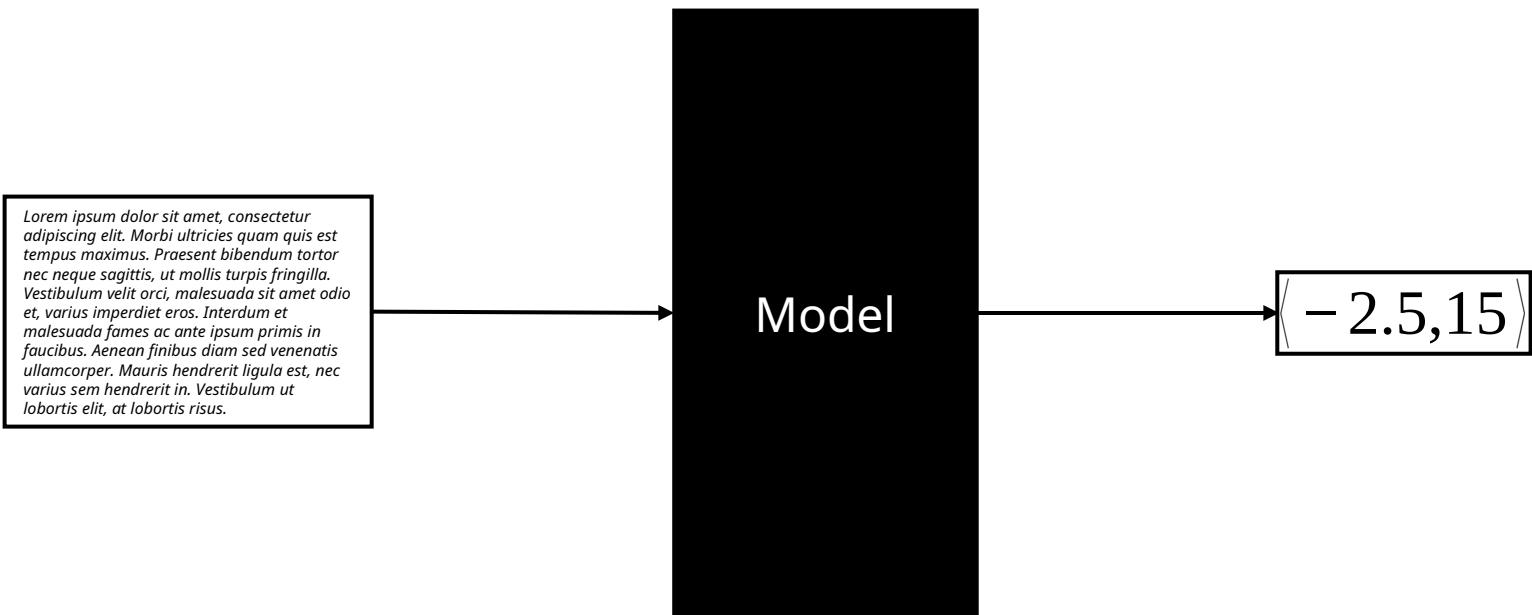
Lorem ipsum dolor sit amet, consectetur adipiscing elit. Morbi ultricies quam quis est tempus maximus. Praesent bibendum tortor nec neque sagittis, ut mollis turpis fringilla. Vestibulum velit orci, malesuada sit amet odio et, varius imperdiet eros. Interdum et malesuada fames ac ante ipsum primis in faucibus. Aenean finibus diam sed venenatis ullamcorper. Mauris hendrerit ligula est, nec varius sem hendrerit in. Vestibulum ut lobortis elit, at lobortis risus.

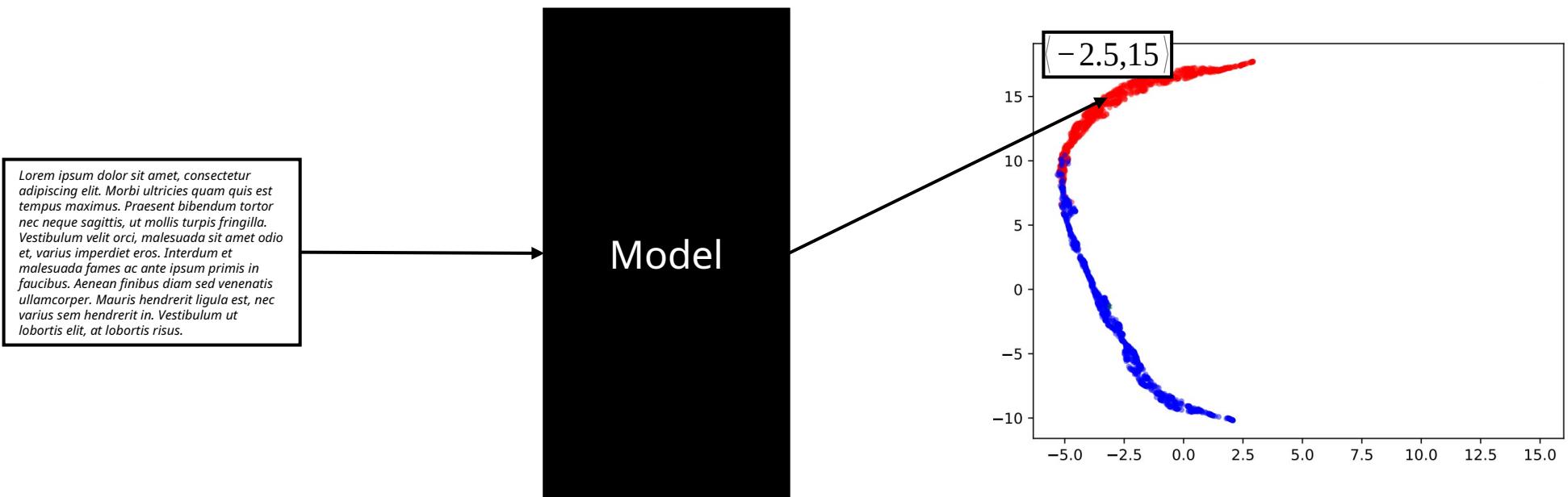
*Lorem ipsum dolor sit amet, consectetur
 adipiscing elit. Morbi ultricies quam quis est
 tempus maximus. Praesent bibendum tortor
 nec neque sagittis, ut mollis turpis fringilla.
 Vestibulum velit orci, malesuada sit amet odio
 et, varius imperdier eros. Interdum et
 malesuada fames ac ante ipsum primis in
 faucibus. Aenean finibus diam sed venenatis
 ullamcorper. Mauris hendrerit ligula est, nec
 variis sem hendrerit in. Vestibulum ut
 lobortis elit, at lobortis risus.*

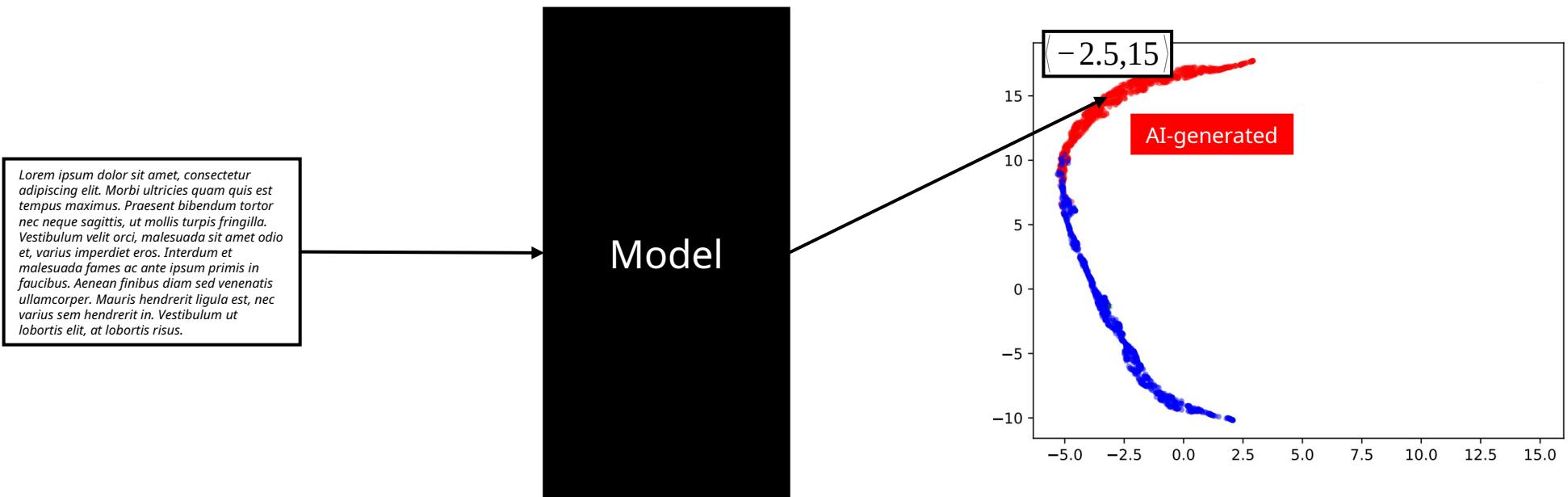


Model







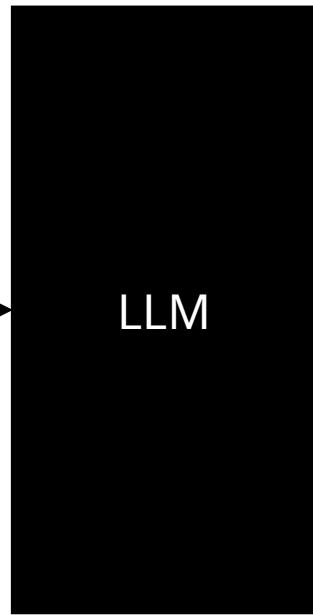


Perplexity-based detectors

How do LLMs
work?

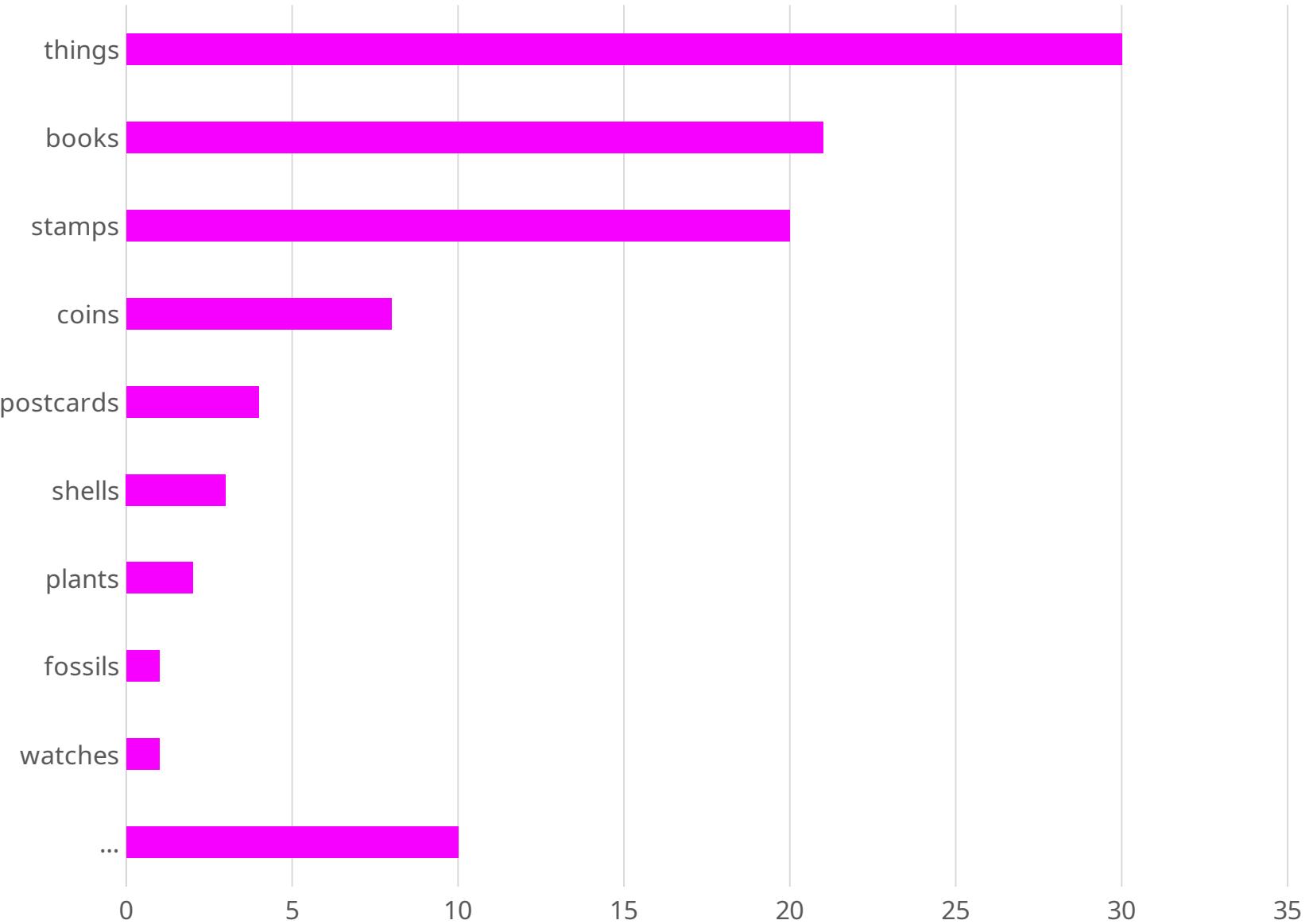
As a hobby, I like to collect _____

As a hobby, I like to collect

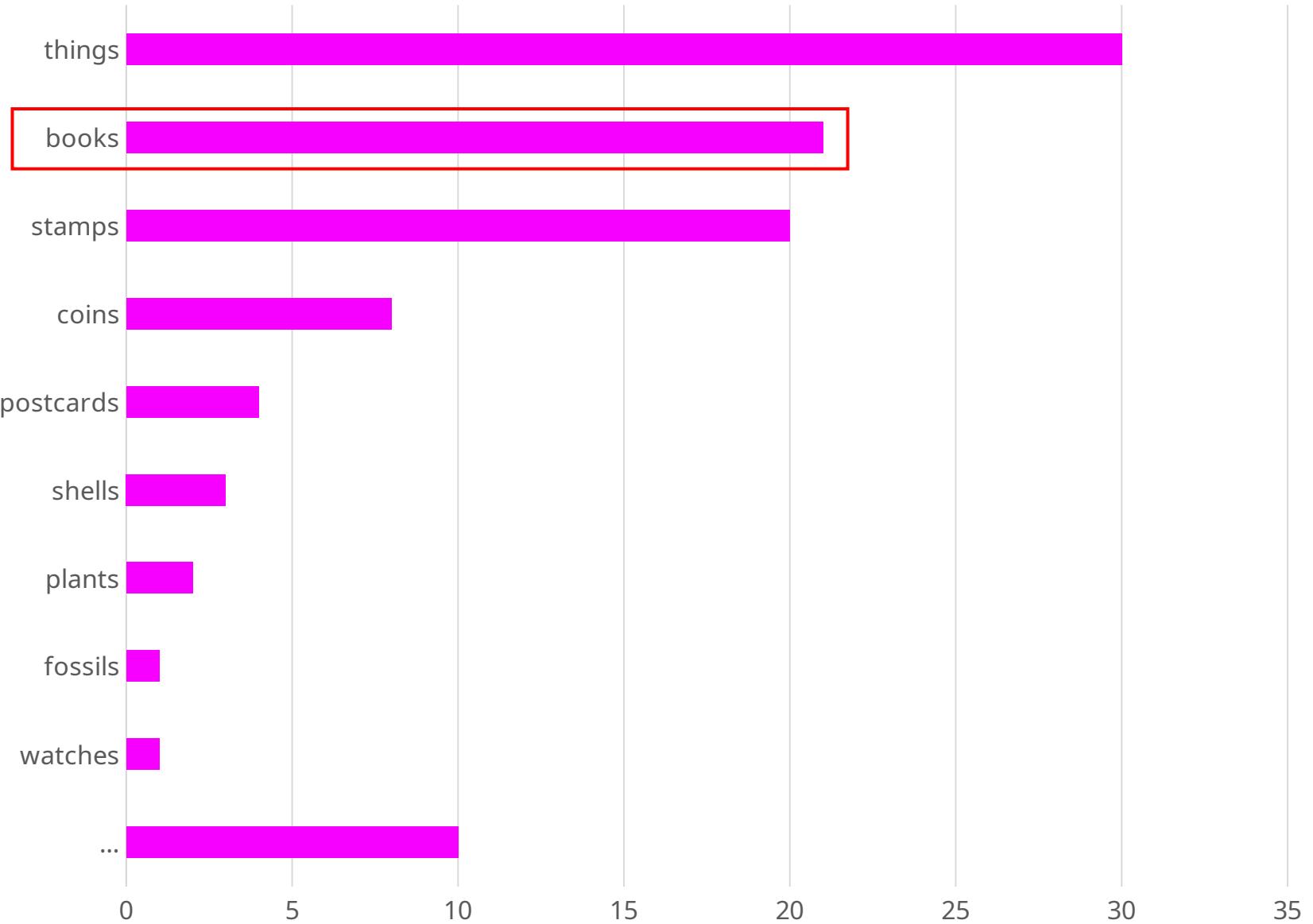


LLM

As a hobby, I like to collect



As a hobby, I like to collect



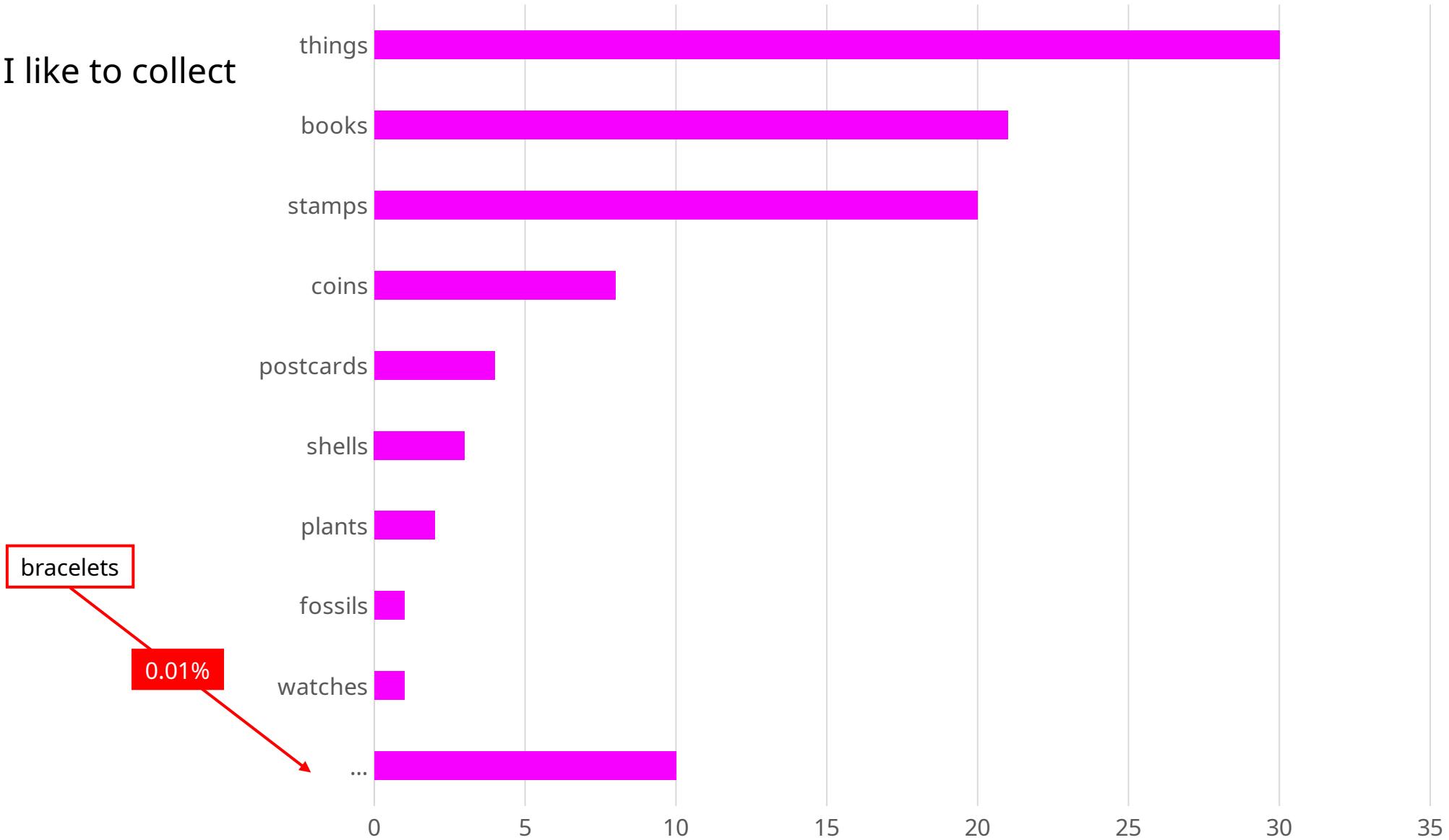
As a hobby, I like to collectbooks

AI-generated

As a hobby, I like to collectbracelets

Human-written

As a hobby, I like to collect



As a hobby, I like to collect

books

↑ probability

bracelets

↓ probability

As a hobby, I like to collect

books

↑ probability

↓ perplexity

bracelets

↓ probability

↑ perplexity

As a hobby, I like to collect

books

↑ probability

↓ perplexity

AI-generated

bracelets

↓ probability

↑ perplexity

human-written

Repeat for the whole text 
prediction

Watermarks

**OpenAI has built a
text watermarking
method to detect
ChatGPT-written
content**

COMPANY HAS MULLED ITS
RELEASE OVER THE PAST YEAR

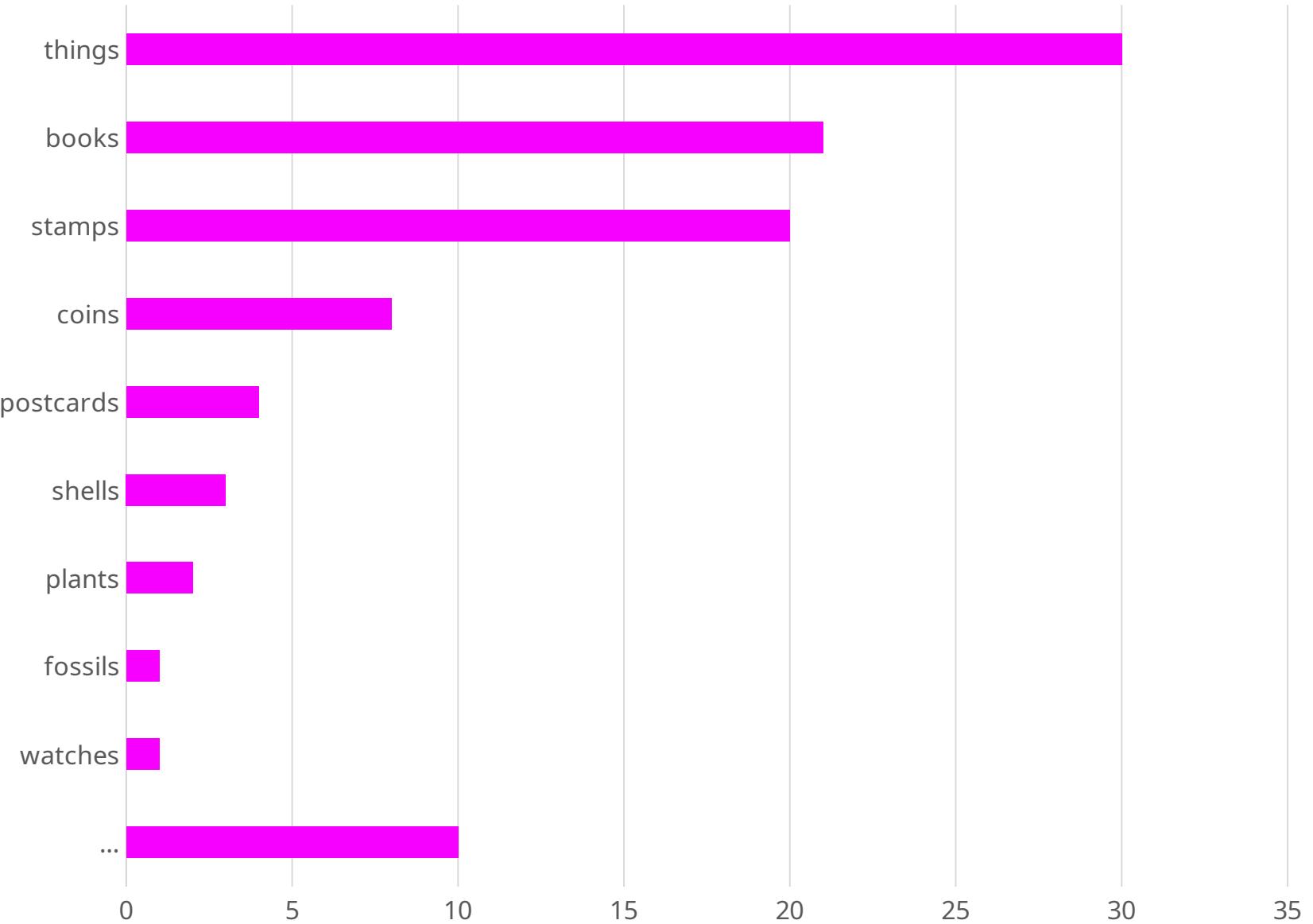
**OpenAI has the tech to watermark
ChatGPT text—it just won't release it**

Some say watermarking is the responsible thing to do...

**OpenAI sitting on tool to
watermark AI-generated content**

The company has yet to release it, because it fears losing users

As a hobby, I like to collect



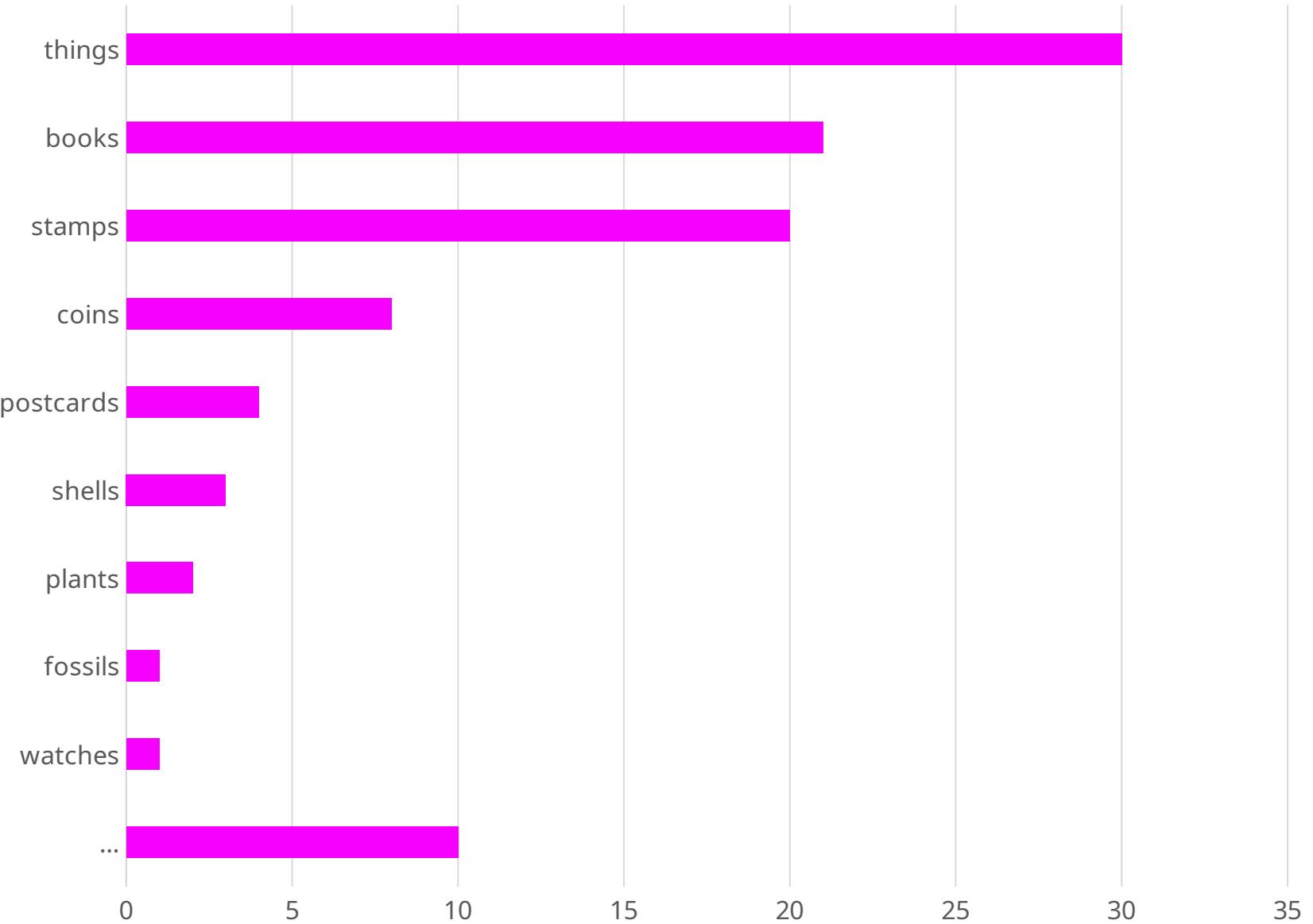
“Pick the N-th word. N is the number of letters in the previous word”

“Pick the N-th word. N is the number of letters in the previous word”

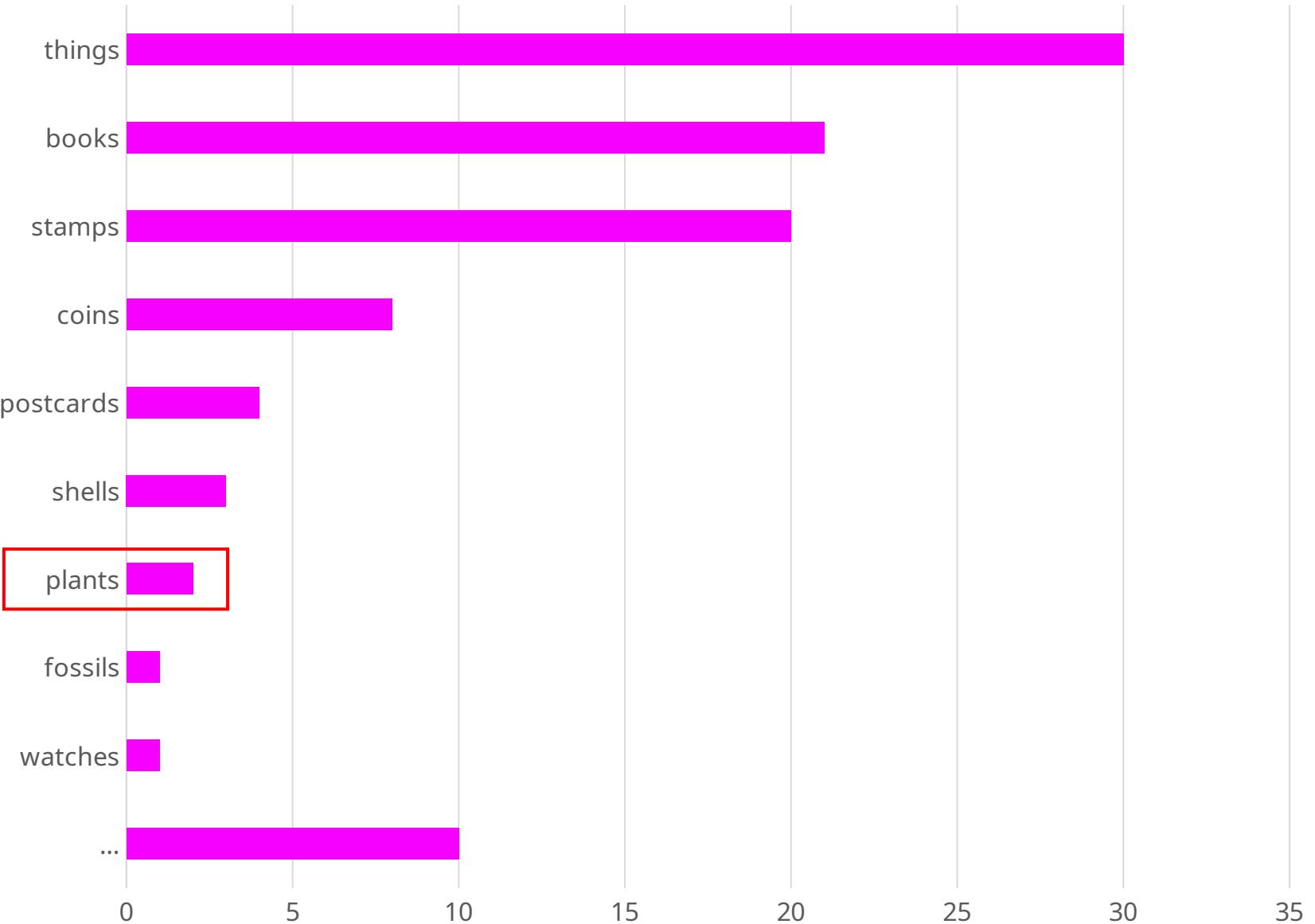
As a hobby, I like to collect

7 letters

As a hobby, I like to collect



As a hobby, I like to collect



As a hobby, I like to collect plants

Watermarked

Repeat for all of the words that we
write

If a lot of words follow the rule ↗ AI-generated

* side
note

Detection of AI-generated texts: different techniques

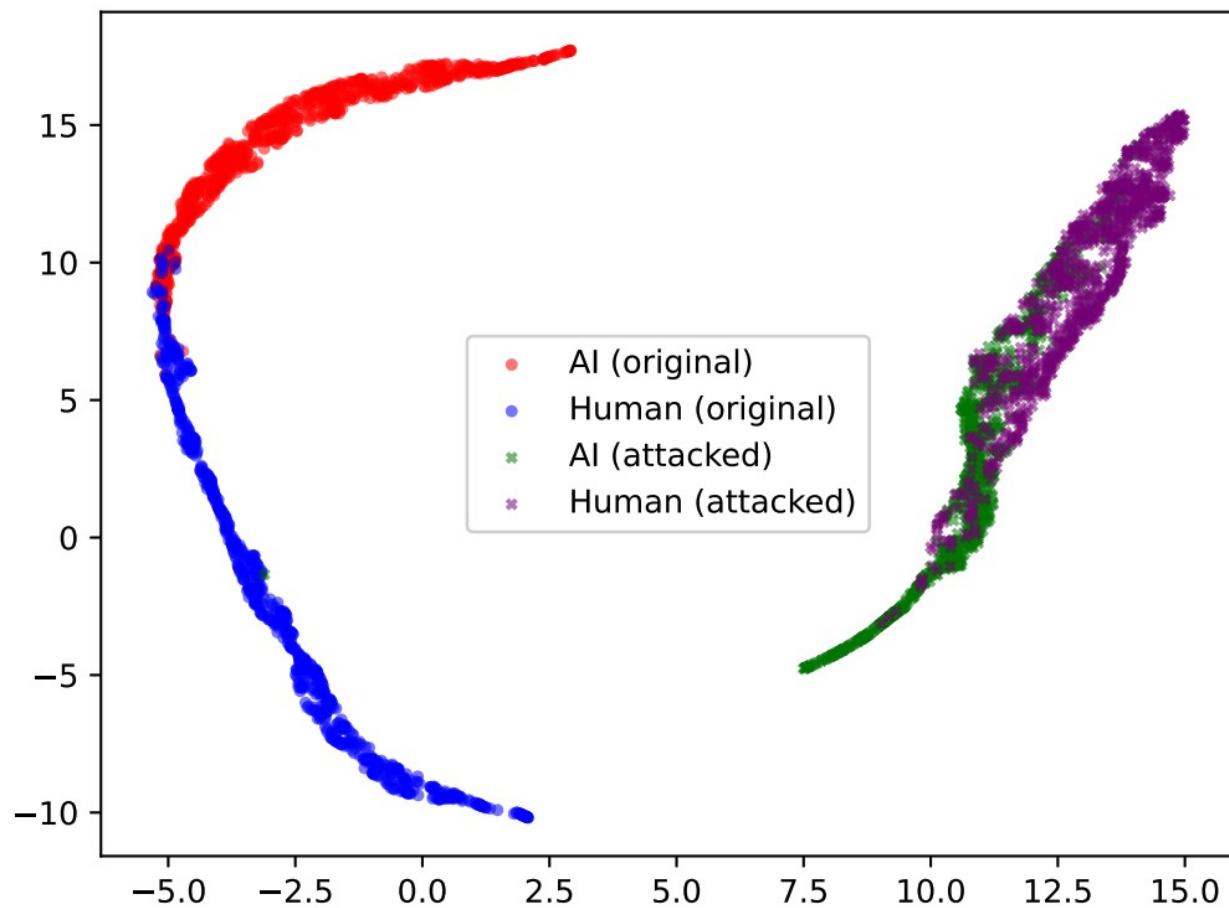
Homoglyph-based attacks

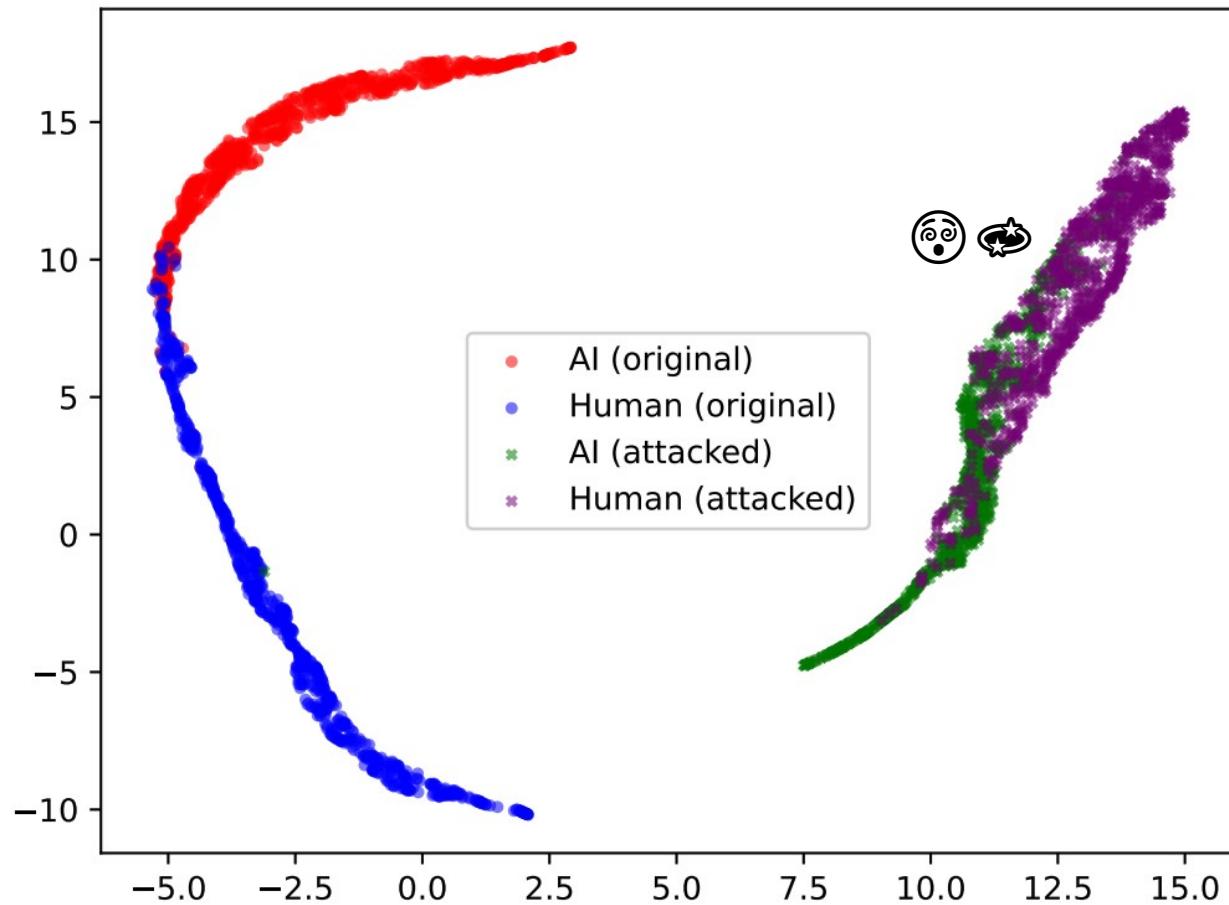
Dem
o

Homoglyph-based attacks: text becomes **unrecognizable**

Technical analysis

Classifiers



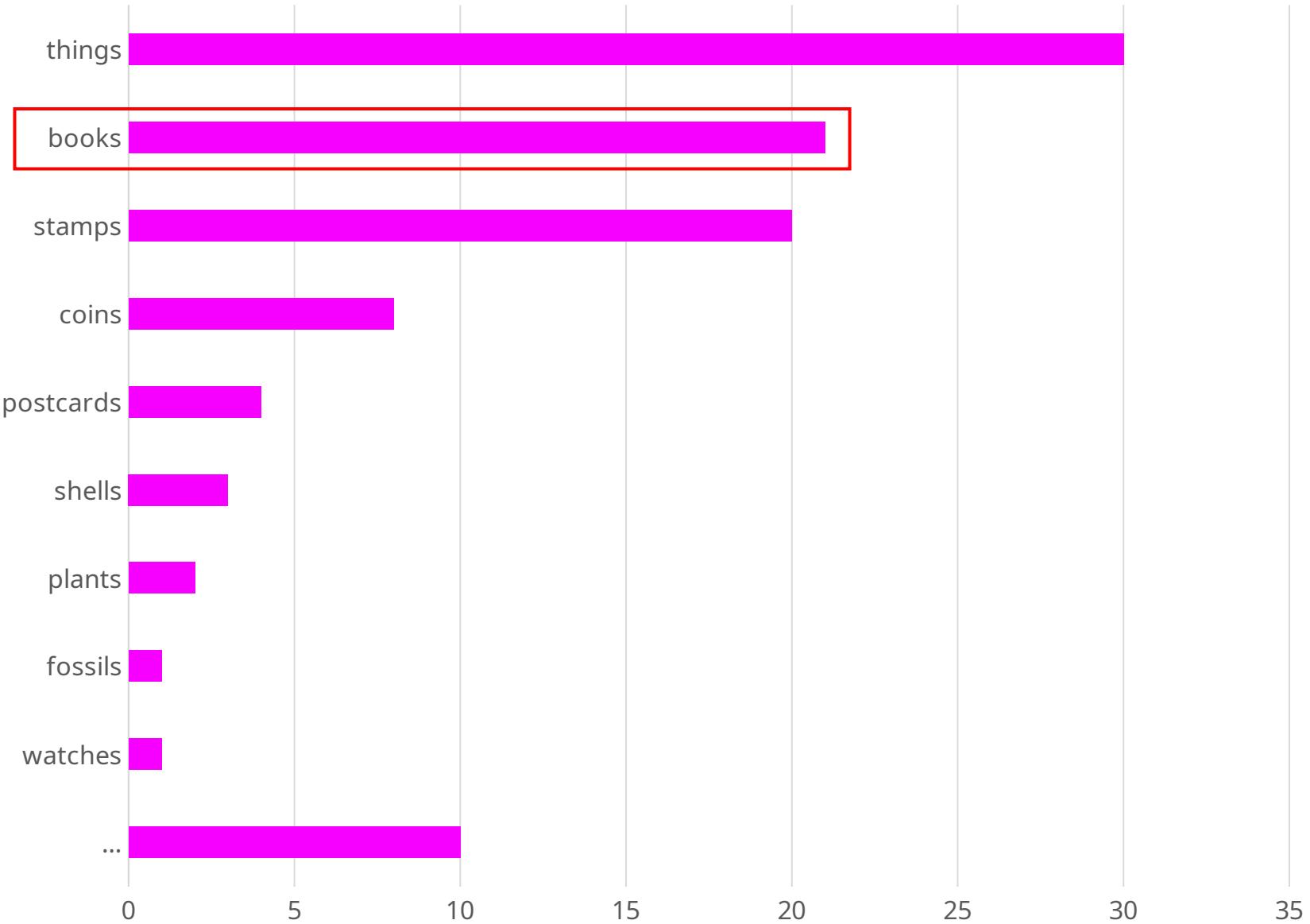


Perplexity-based detectors

As a hobby, I like to collectbooks

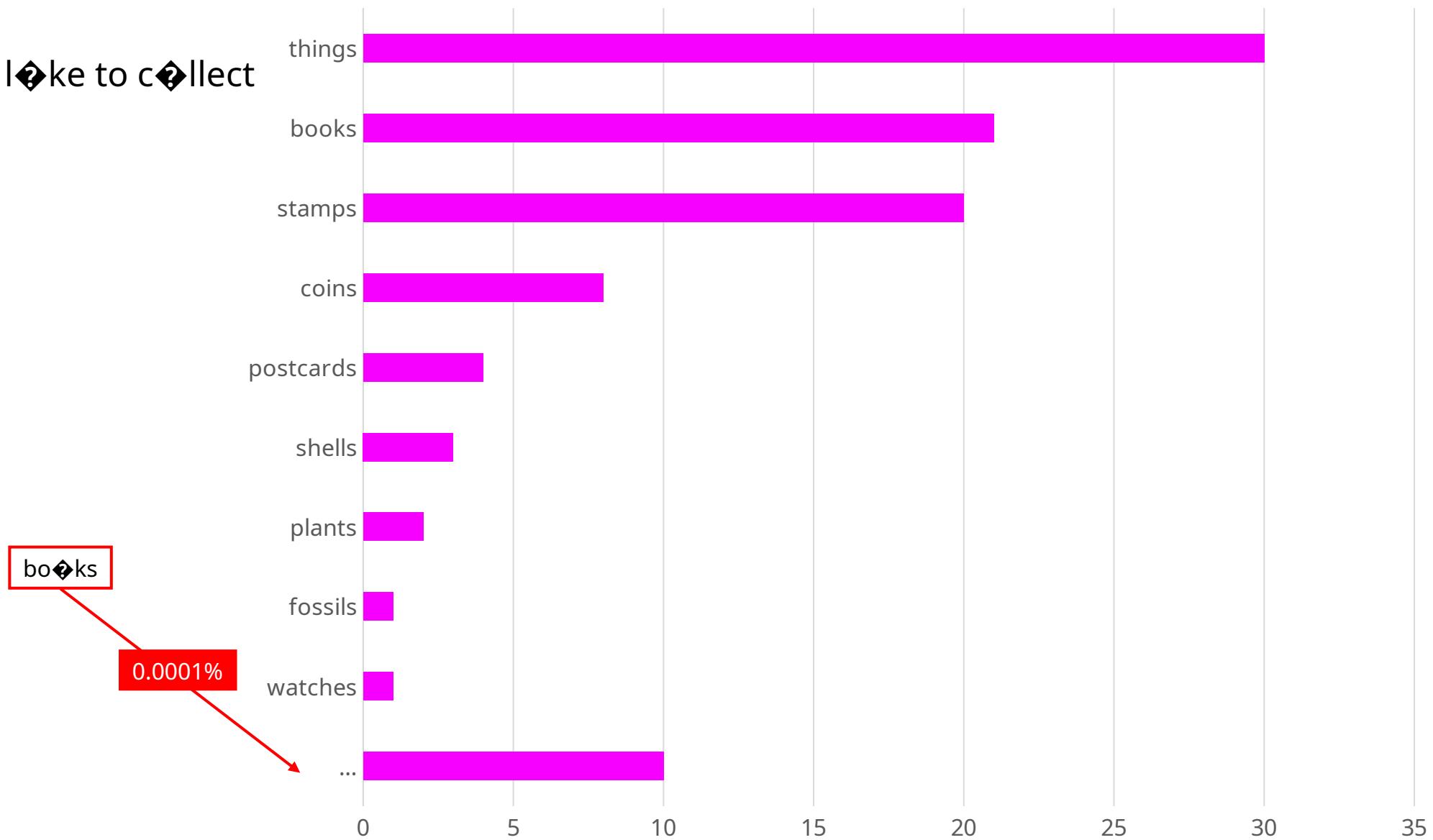
AI-generated

As a hobby, I like to collect



As a hobb?, I l?ke to c?llct?ks

As a hobby, I like to collect



As a hobb?, I l?ke to c?llct?ks

Human-written

Watermarks

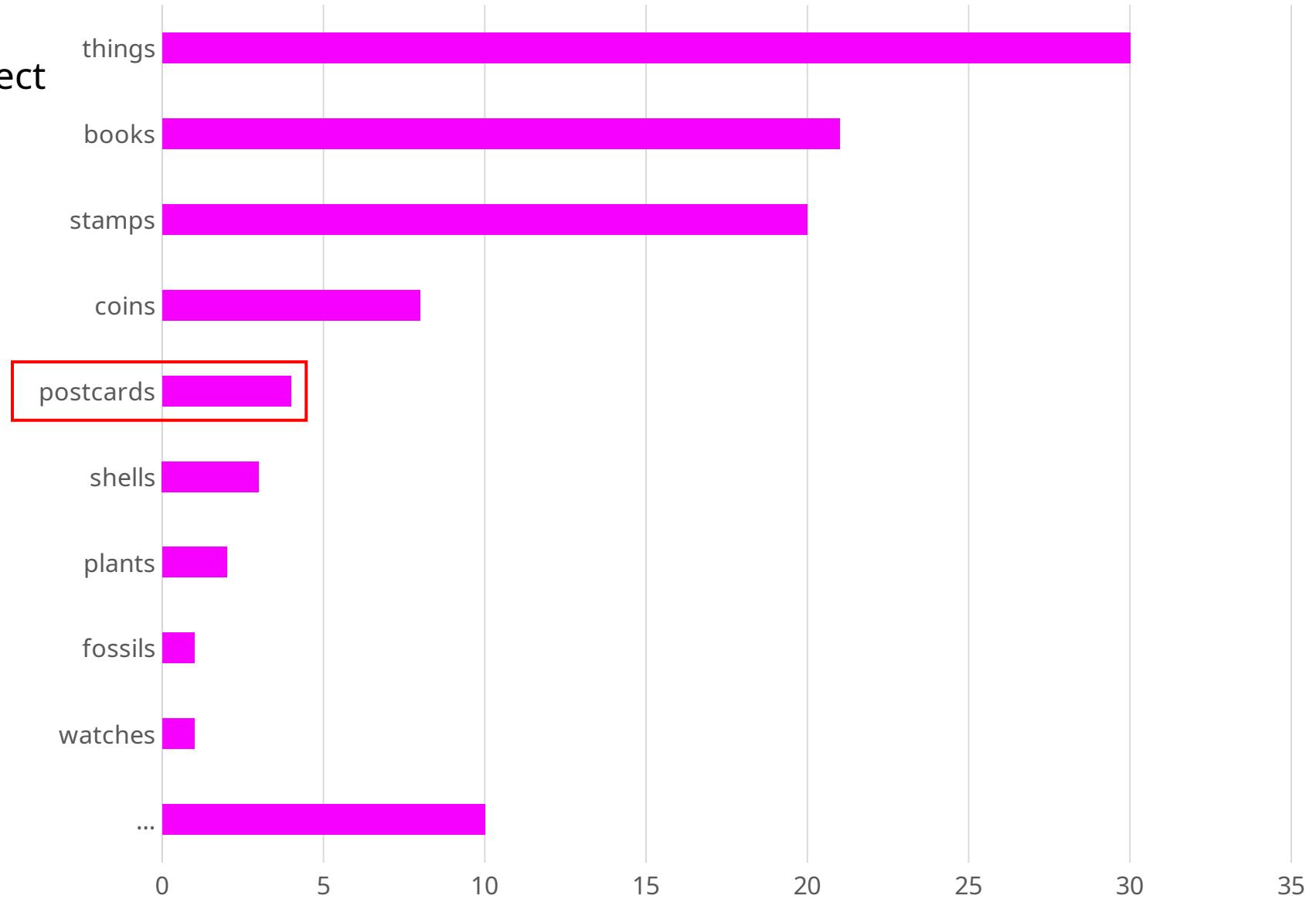
As a hobby, I like to collect plants

As a hobb♦, I l♦ke to c♦llect plants

“Pick the N-th word. N is the number of letters in the previous word.”

As a hobb?, I l?ke to c?llct plants

As a hobby, I like to collect

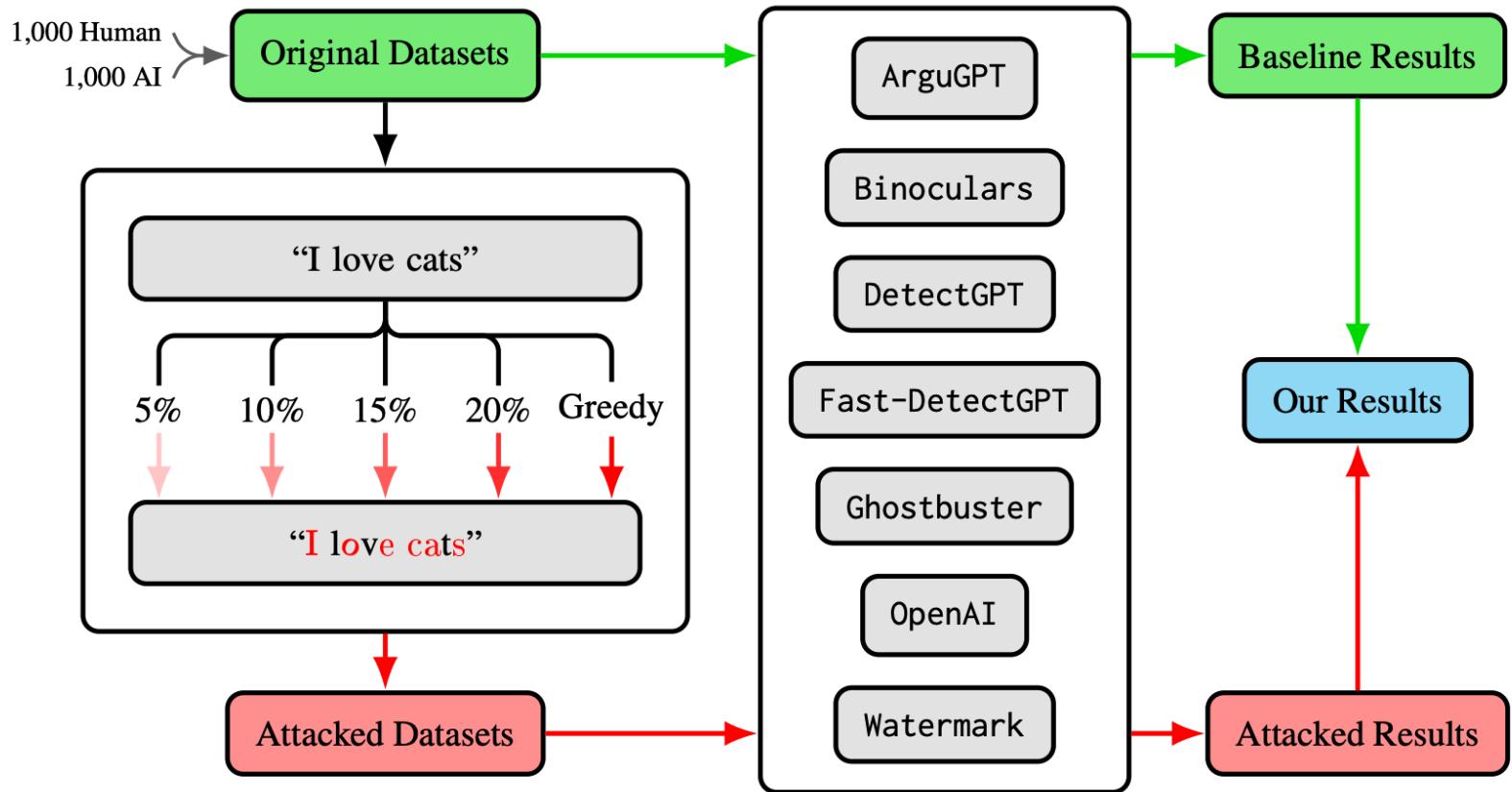


Doesn't follow the rule ↗ human-written

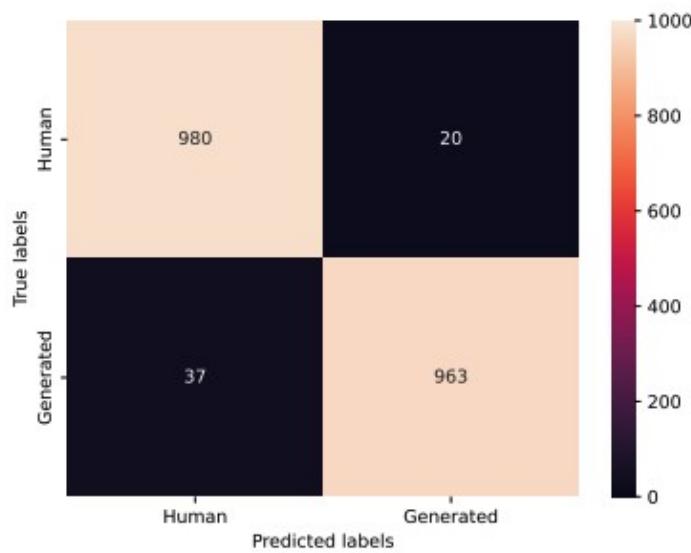
Different mechanisms of action – all exploit “confusion”

Effectiveness

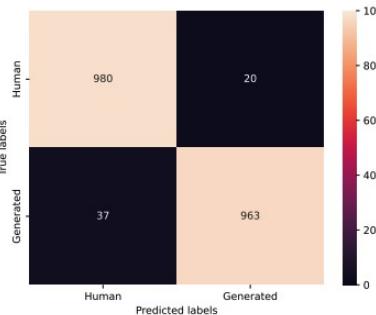
Experimental approach



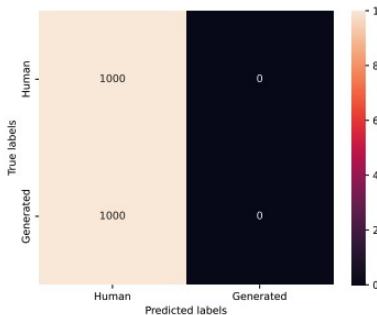
Results



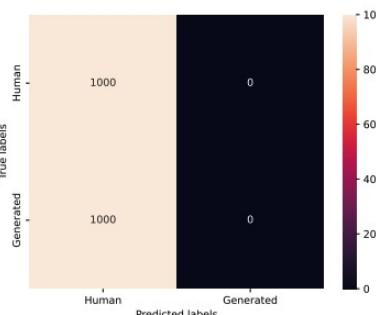
(a) No attack



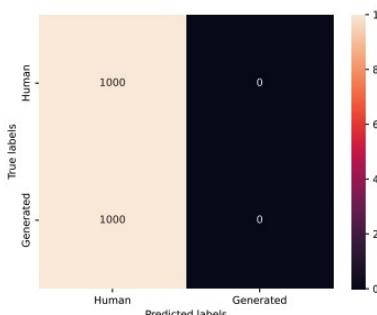
(a) No attack



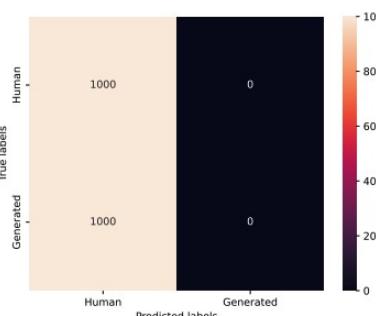
(b) Random attack (5%)



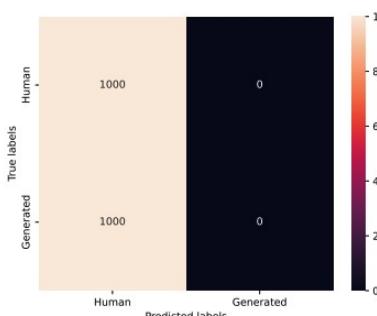
(c) Random attack (10%)



(d) Random attack (15%)



(e) Random attack (20%)



(f) Greedy attack

Dataset	Detector	Original	5%	10%	15%	20%	Greedy
	ArguGPT	0.94	0.0	0.0	0.0	0.0	0.0

Dataset	Detector	Original	5%	10%	15%	20%	Greedy
<i>CHEAT</i>	ArguGPT	0.94	0.0	0.0	0.0	0.0	0.0
	Binoculars	0.93					
	DetectGPT	0.14					
	Fast-DetectGPT	0.9					
	Ghostbuster	0.64					
	OpenAI	0.47					
<i>essay</i>	ArguGPT	0.92					
	Binoculars	0.91					
	DetectGPT	0.24					
	Fast-DetectGPT	0.88					
	Ghostbuster	0.92					
<i>reuter</i>	ArguGPT	0.92					
	Binoculars	0.8					
	DetectGPT	0.23					
	Fast-DetectGPT	0.92					
	Ghostbuster	0.93					
	OpenAI	0.27					
<i>writing prompts</i>	ArguGPT	0.39					
	Binoculars	0.85					
	DetectGPT	0.44					
	Fast-DetectGPT	0.79					
	Ghostbuster	0.88					
<i>realnewslike</i>	OpenAI	-0.05					
	Watermark	0.92					
Average		0.64					
Standard deviation		0.36					

Dataset	Detector	Original	5%	10%	15%	20%	Greedy
<i>CHEAT</i>	ArguGPT	0.94	0.0	0.0	0.0	0.0	0.0
	Binoculars	0.93	0.37	0.11	0.04	0.02	0.13
	DetectGPT	0.14	-0.02	0.03	0.13	0.06	0.0
	Fast-DetectGPT	0.9	0.23	0.04	0.02	0.0	-0.01
	Ghostbuster	0.64	0.41	0.32	0.12	0.06	0.02
	OpenAI	0.47	0.0	0.0	0.0	-0.02	0.0
<i>essay</i>	ArguGPT	0.92	0.0	0.0	0.0	0.0	0.0
	Binoculars	0.91	0.22	0.05	0.0	0.0	0.05
	DetectGPT	0.24	-0.01	0.11	0.21	0.08	0.0
	Fast-DetectGPT	0.88	0.22	0.04	0.0	0.0	-0.08
	Ghostbuster	0.92	0.73	0.51	0.13	0.0	0.0
	OpenAI	-0.21	0.0	0.0	0.0	0.0	0.03
<i>reuter</i>	ArguGPT	0.92	0.0	0.0	0.0	0.0	0.0
	Binoculars	0.8	0.22	0.07	0.03	0.02	0.08
	DetectGPT	0.23	0.0	0.03	0.34	0.14	0.0
	Fast-DetectGPT	0.92	0.28	0.1	0.02	0.0	0.04
	Ghostbuster	0.93	0.61	0.51	0.16	0.04	0.0
	OpenAI	0.27	0.0	-0.04	-0.09	-0.11	-0.06
<i>writing prompts</i>	ArguGPT	0.39	0.0	0.0	0.0	0.0	0.0
	Binoculars	0.85	0.2	0.0	0.0	0.0	-0.04
	DetectGPT	0.44	0.04	0.01	0.02	0.02	0.0
	Fast-DetectGPT	0.79	0.3	0.05	-0.03	0.0	-0.33
	Ghostbuster	0.88	0.42	0.64	0.33	0.09	0.0
	OpenAI	-0.05	-0.04	-0.05	-0.13	-0.11	0.01
<i>realnewslike</i>	Watermark	0.92	0.18	-0.01	0.0	-0.03	0.0
Average		0.64	0.17	0.1	0.05	0.01	-0.01
Standard deviation		0.36	0.21	0.19	0.11	0.05	0.08

Highly effective, renders all detectors **ineffective**

Implications

Twitter
bots

Academic
misconduct

Fraud, social
engineering

Identity
theft

Fake news, loss of trust in
media

Lower access barrier: increased risk

Safeguards

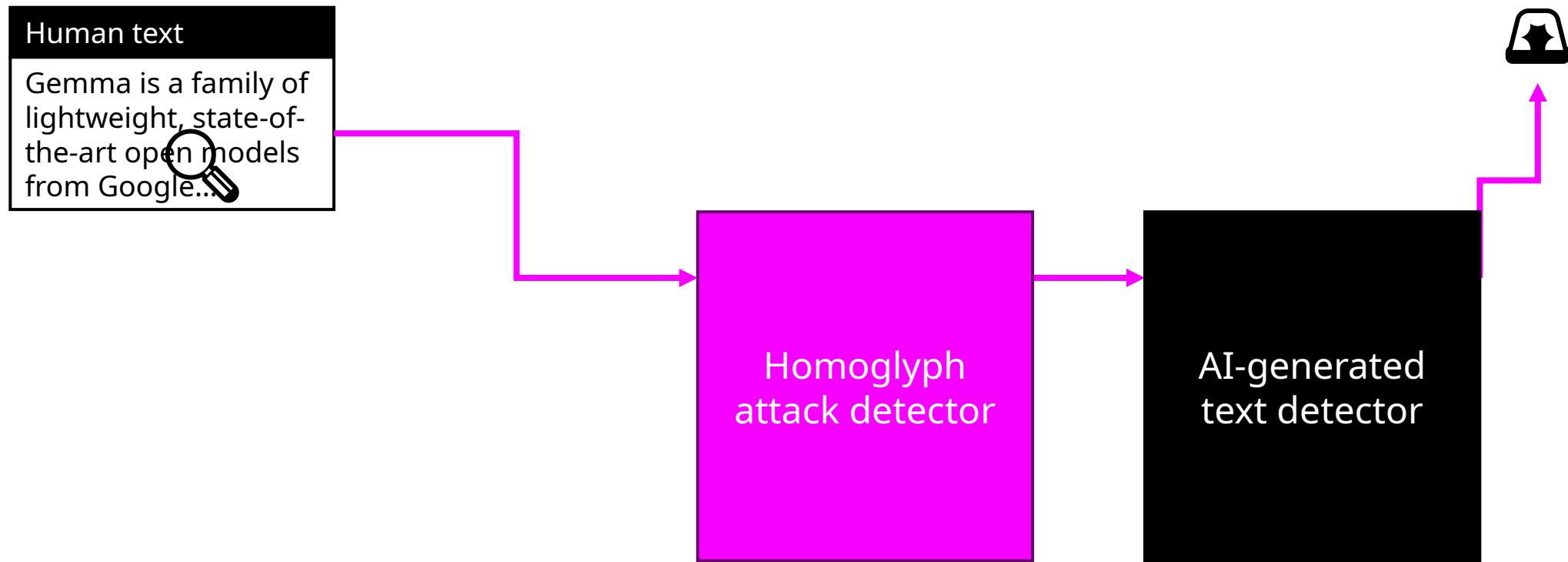
Safeguards

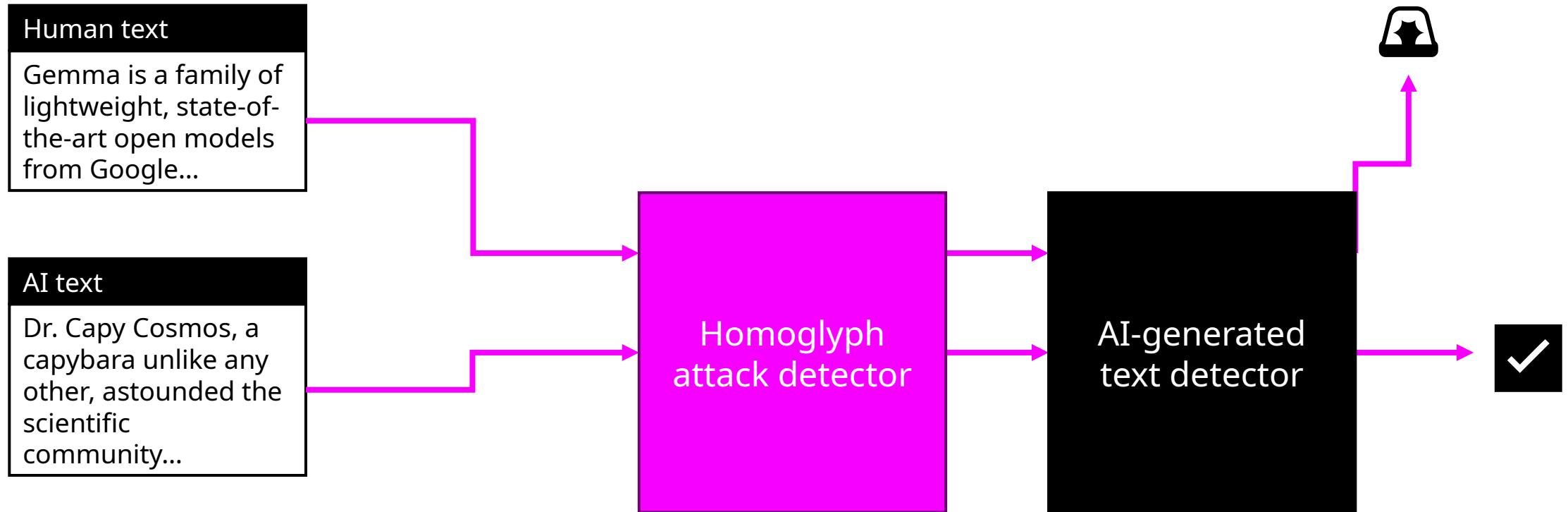
Input constraints

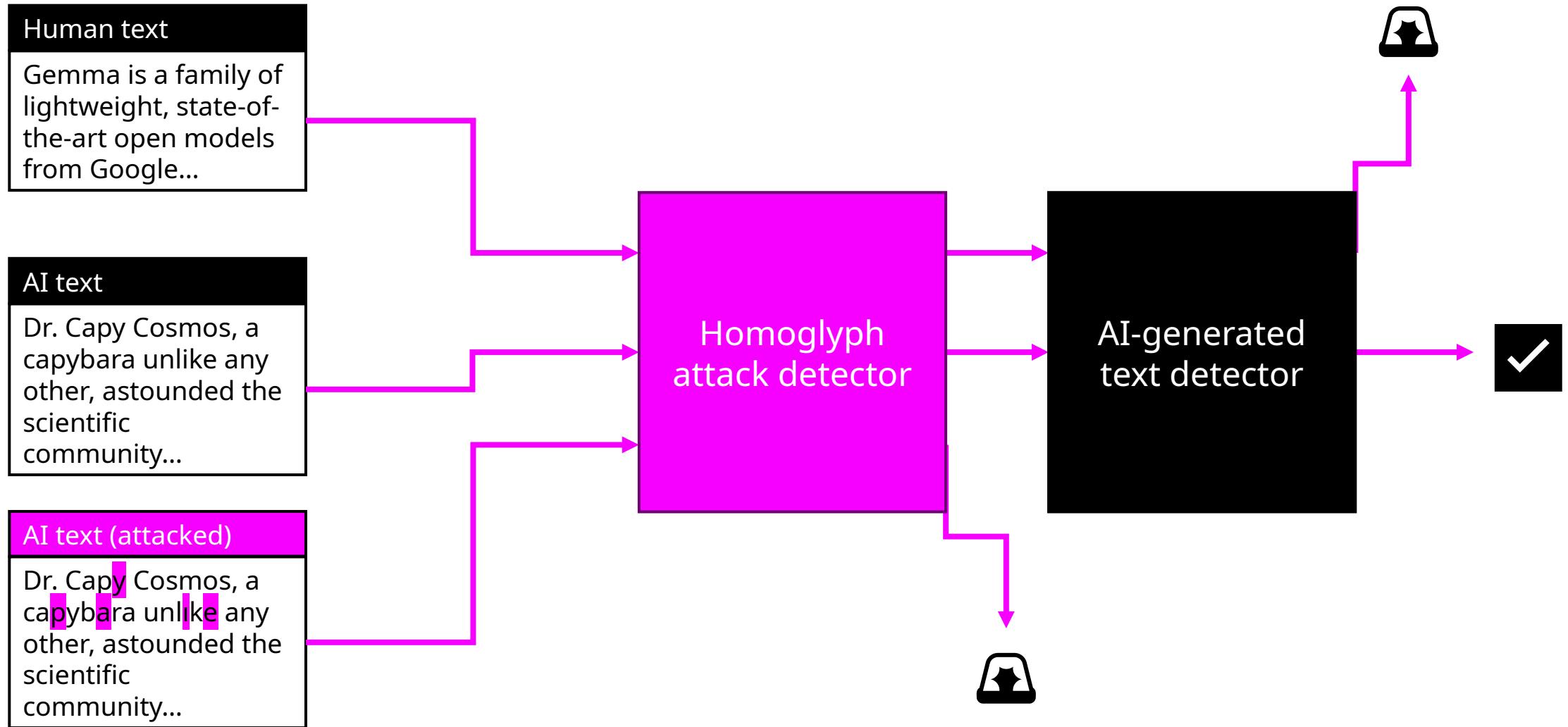
AI-generated
text detector

Homoglyph
attack detector

AI-generated
text detector







Maybe not
enough...

Maybe not
enough...

Specifically, considering that LLMs utilize decoder-only structures, the distribution for the i -th token at the final layer L , denoted as p_i^L , depends on the past i tokens. This represents the probability of any token in the vocabulary becoming the $(i + 1)$ -th token. The distribution difference between the target token and the alternative token, $p_i^L(w_t) - p_i^L(w_a)$, can be attributed to the cumulative contribution of the first i tokens. Similarly, $p_{i-1}^L(w_t) - p_{i-1}^L(w_a)$ reflects the collective contribution of the initial $i - 1$ tokens.

Maybe not
enough...

Specifically, considering that LLMs utilize decoder-only structures, the distribution for the i -th token at the final layer L , denoted as p_i^L , depends on the past i tokens. This represents the probability of any token in the vocabulary becoming the $(i + 1)$ -th token. The distribution difference between the target token and the alternative token, $p_i^L(w_t) - p_i^L(w_a)$, can be attributed to the cumulative contribution of the first i tokens. Similarly, $p_{i-1}^L(w_t) - p_{i-1}^L(w_a)$ reflects the collective contribution of the initial $i - 1$ tokens.

Maybe not
enough...

Specifically, considering that LLMs utilize decoder-only structures, the distribution for the i -th token at the final layer L , denoted as p_i^L , depends on the past i tokens. This represents the probability of any token in the vocabulary becoming the $(i + 1)$ -th token. The distribution difference between the target token and the alternative token, $p_i^L(w_t) - p_i^L(w_a)$, can be attributed to the cumulative contribution of the first i tokens. Similarly, $p_{i-1}^L(w_t) - p_{i-1}^L(w_a)$ reflects the collective contribution of the initial $i - 1$ tokens.



False positive

We need to get
smarter

Rule-based approach

Specifically, considering that LLMs utilize decoder-only structures, the distribution for the i -th token at the final layer L , denoted as p_i^L , depends on the past i tokens. This represents the probability of any token in the vocabulary becoming the $(i + 1)$ -th token. The distribution difference between the target token and the alternative token, $p_i^L(w_t) - p_i^L(w_a)$, can be attributed to the cumulative contribution of the first i tokens. Similarly, $p_{i-1}^L(w_t) - p_{i-1}^L(w_a)$ reflects the collective contribution of the initial $i - 1$ tokens.

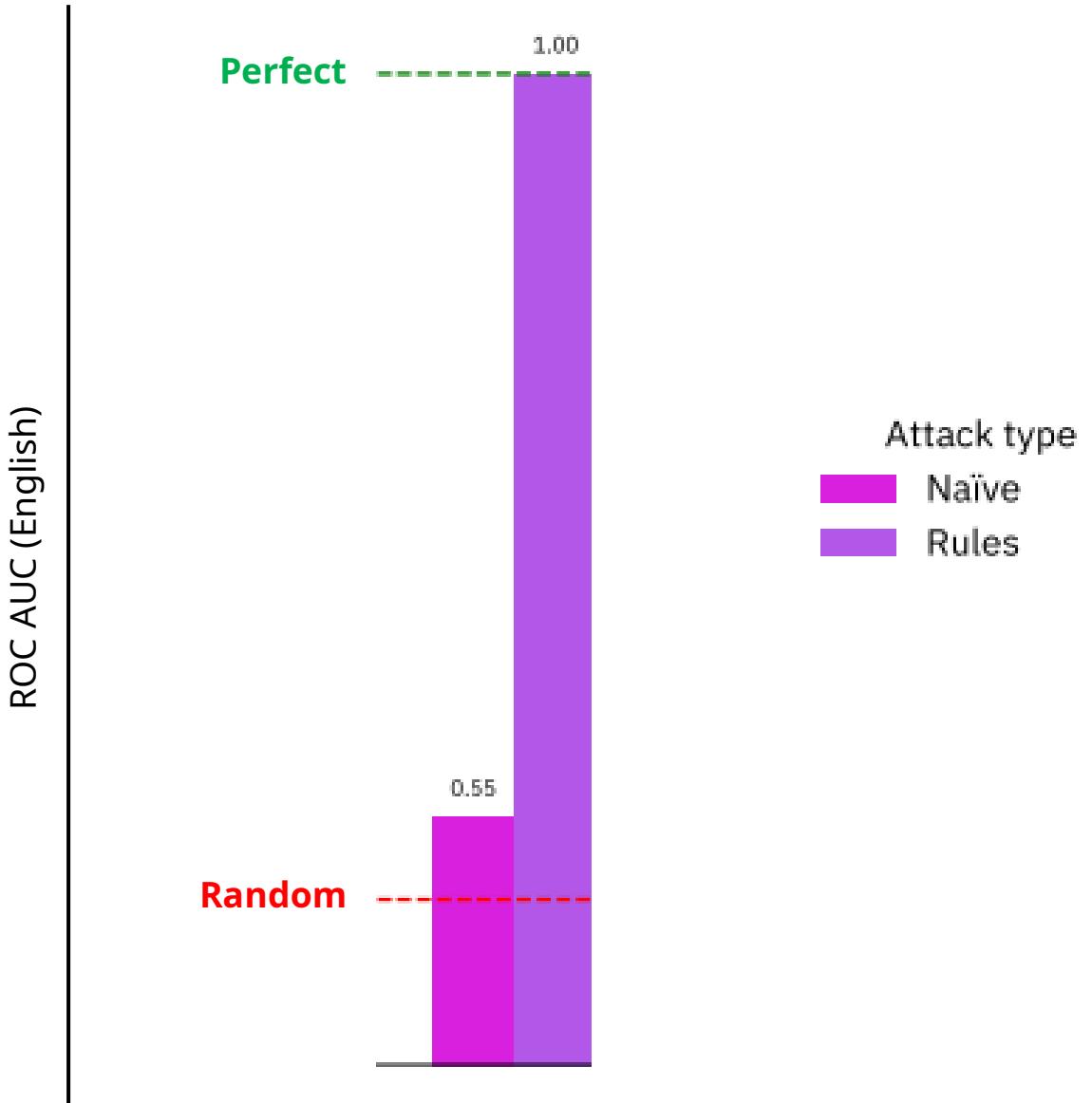
Percentage of words that mix homoglyphs

Rule-based approach

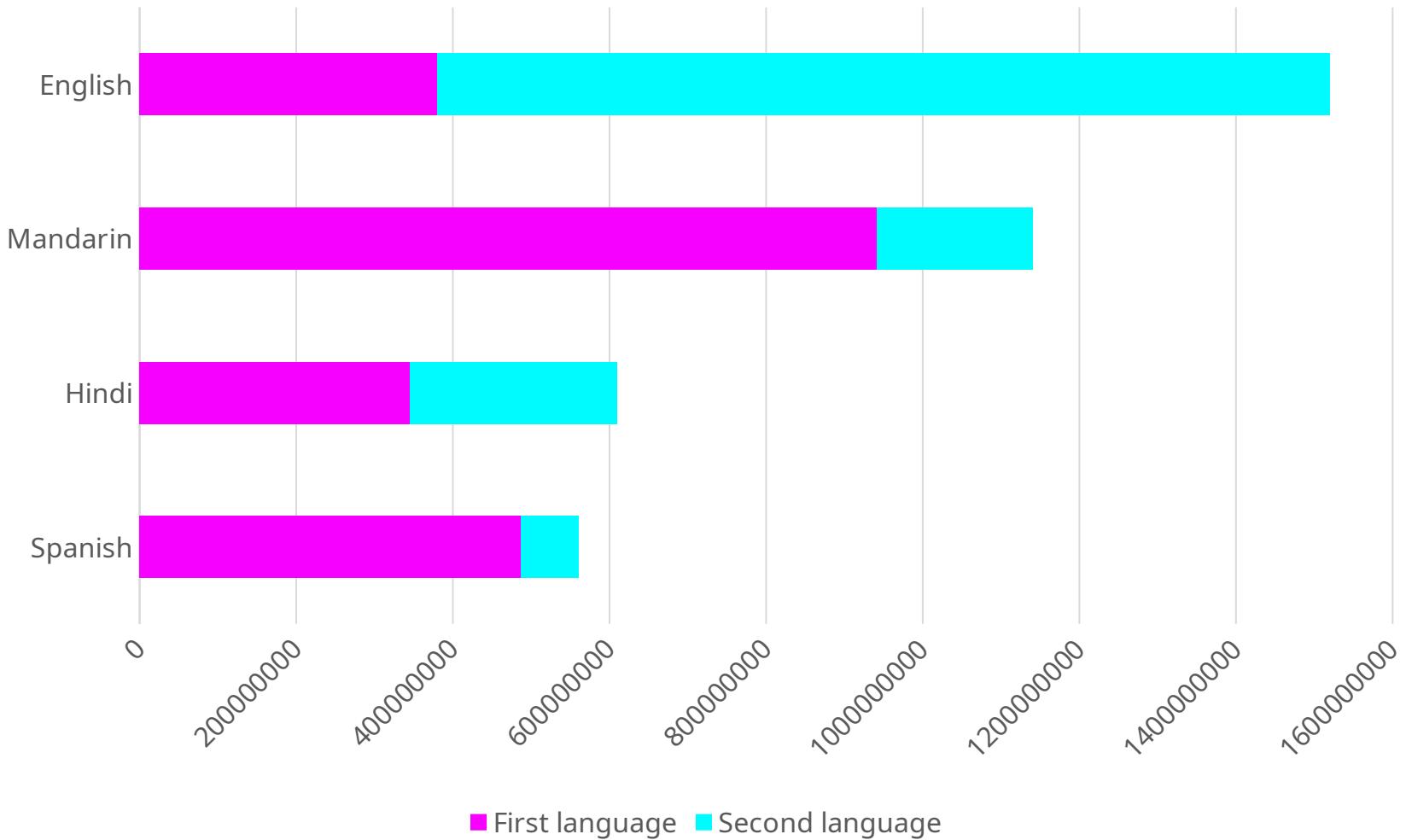
Specifically, considering that LLMs utilize decoder-only structures, the distribution for the i -th token at the final layer L , denoted as p_i^L , depends on the past i tokens. This represents the probability of any token in the vocabulary becoming the $(i + 1)$ -th token. The distribution difference between the target token and the alternative token, $p_i^L(w_t) - p_i^L(w_a)$, can be attributed to the cumulative contribution of the first i tokens. Similarly, $p_{i-1}^L(w_t) - p_{i-1}^L(w_a)$ reflects the collective contribution of the initial $i - 1$ tokens.

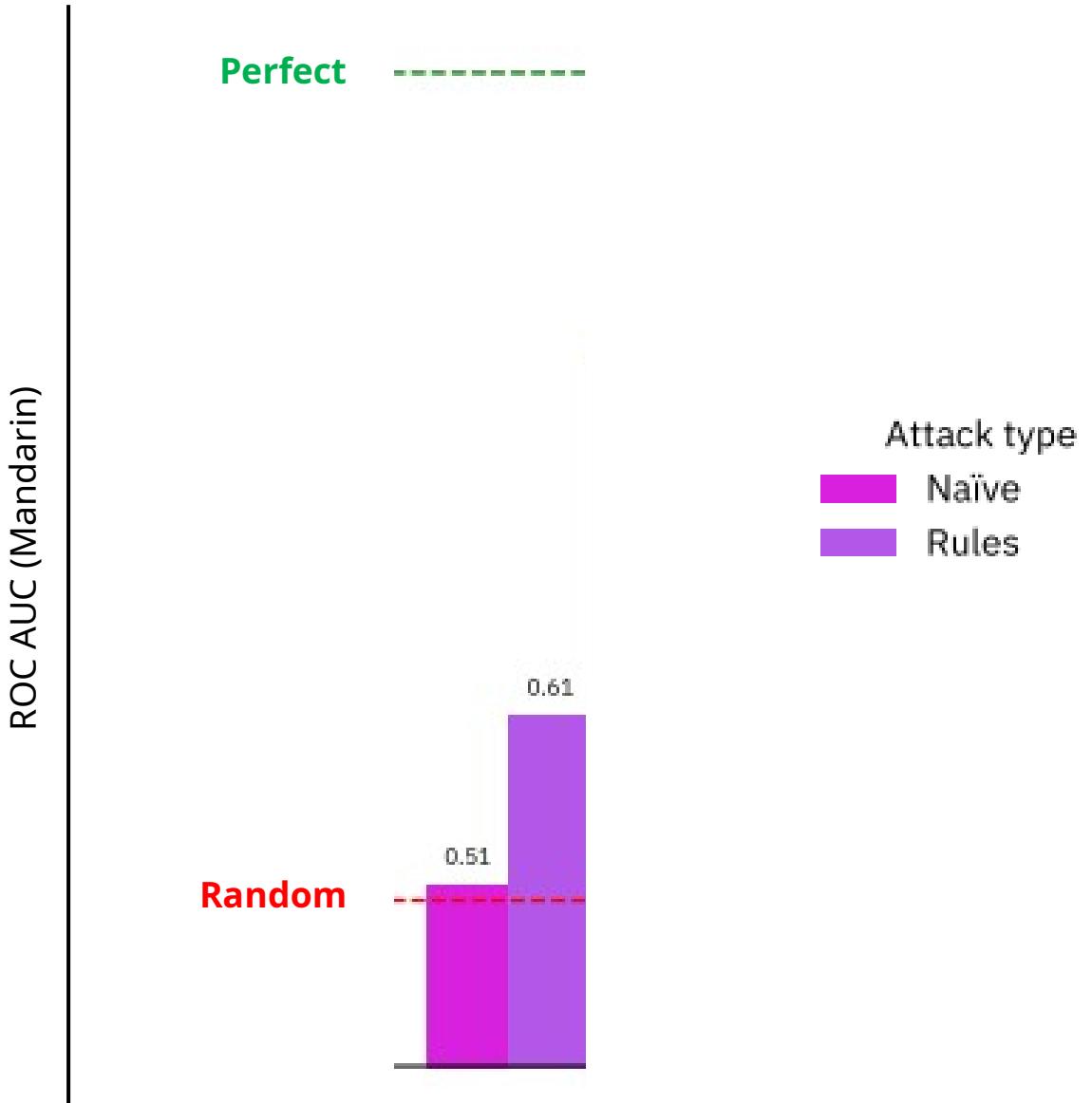


All good



Number of speakers per language in the world





What else can we
do?

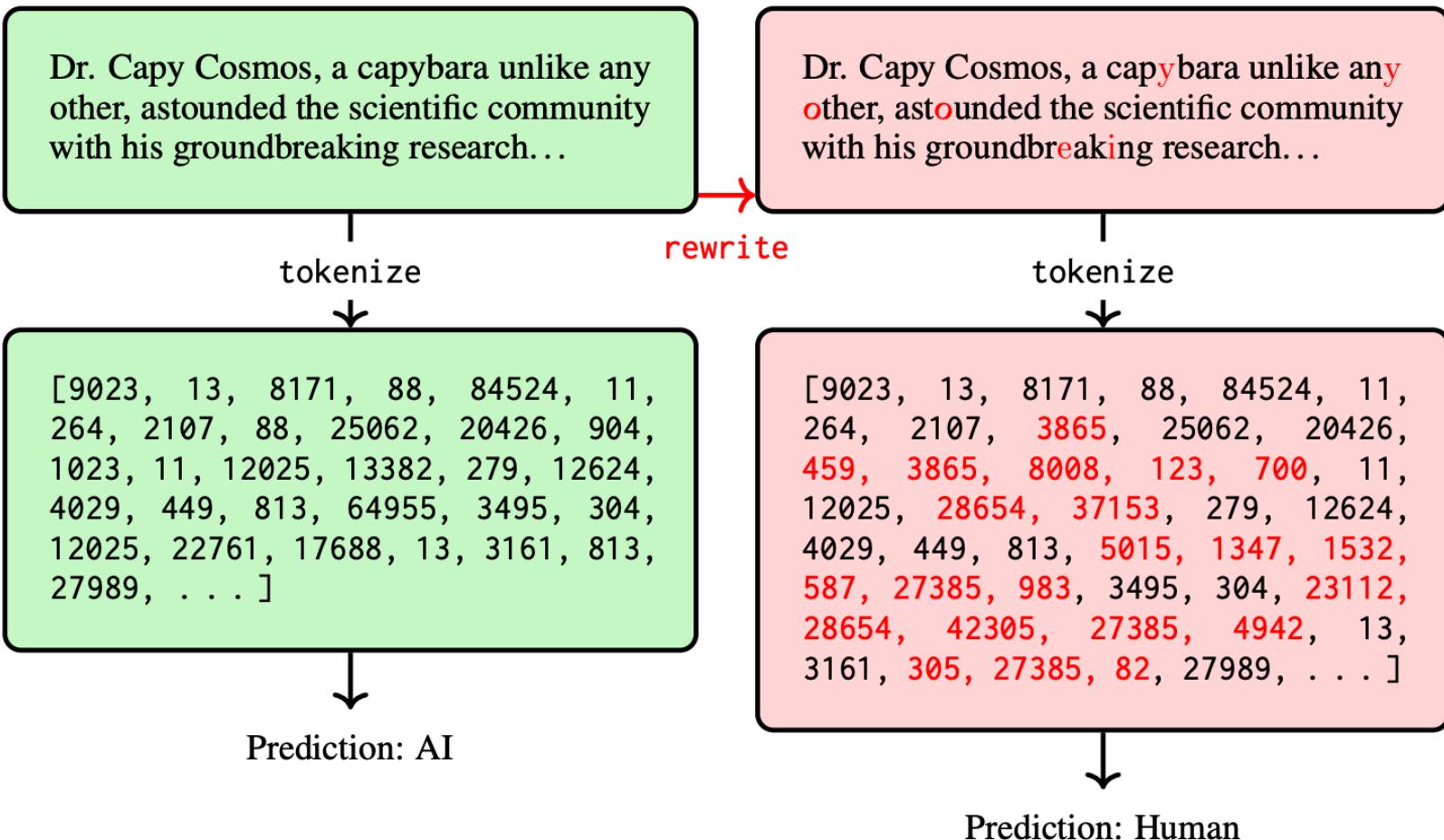
What else can we
do?

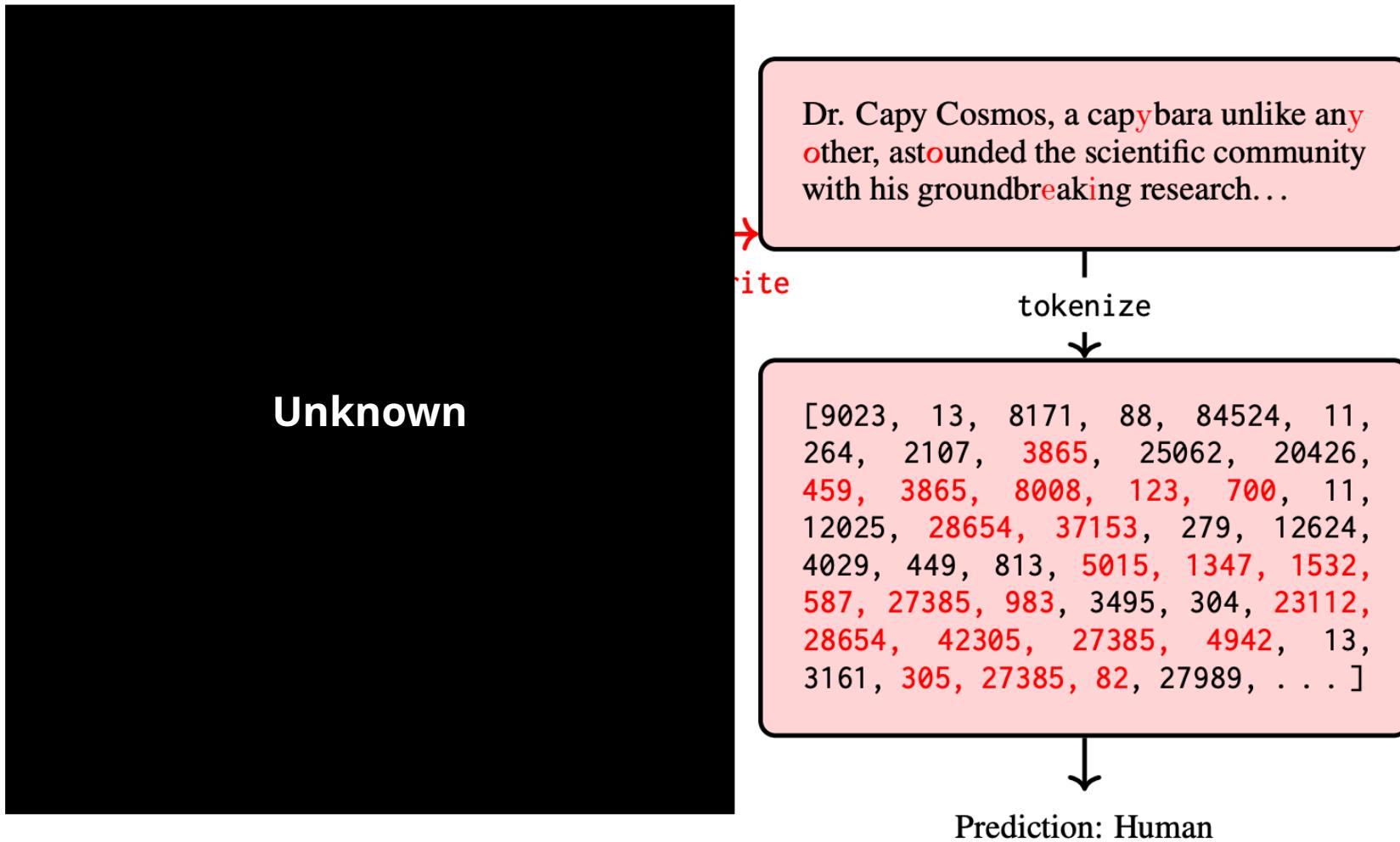
2 idéas

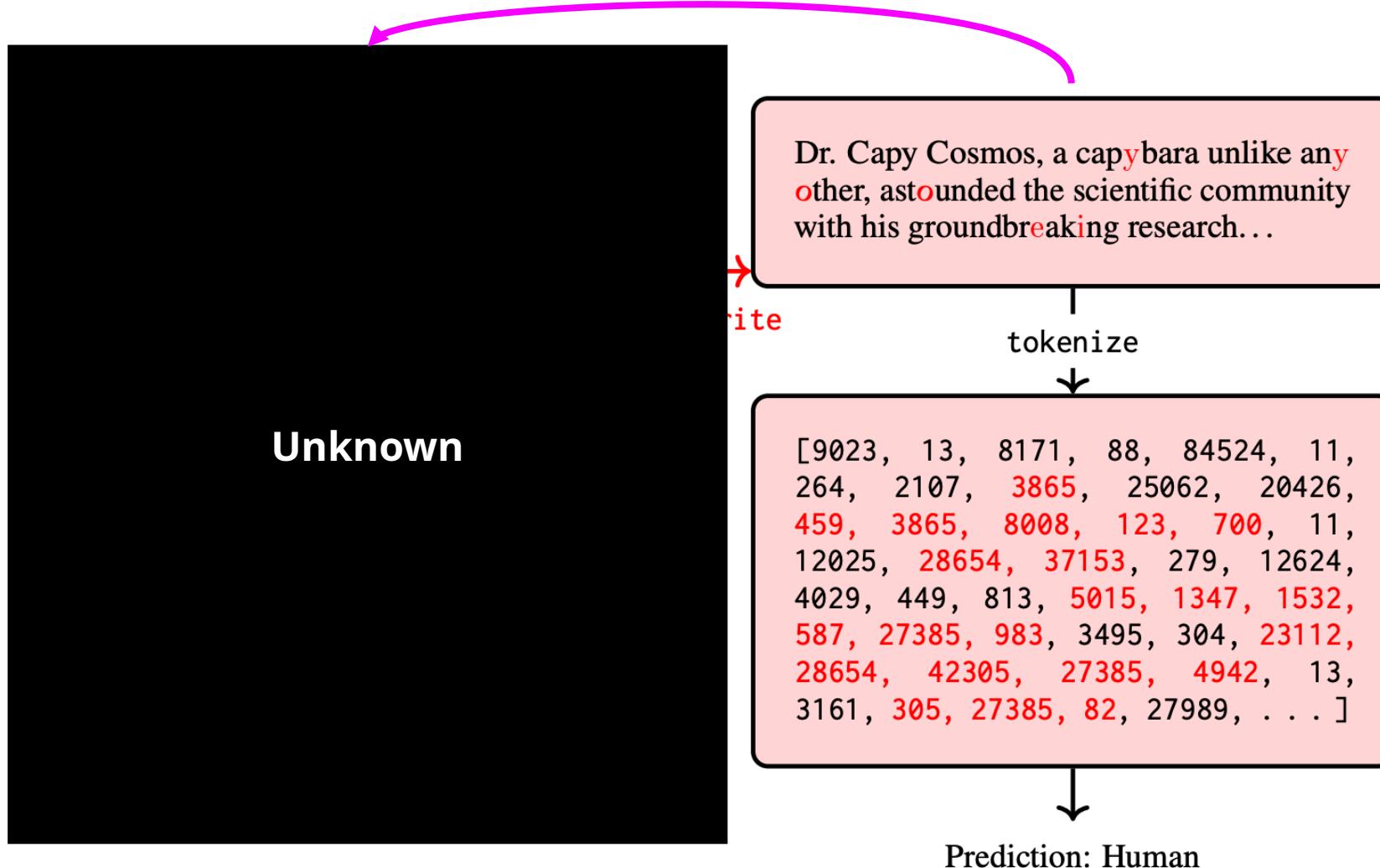
Idea 1: Tokenization

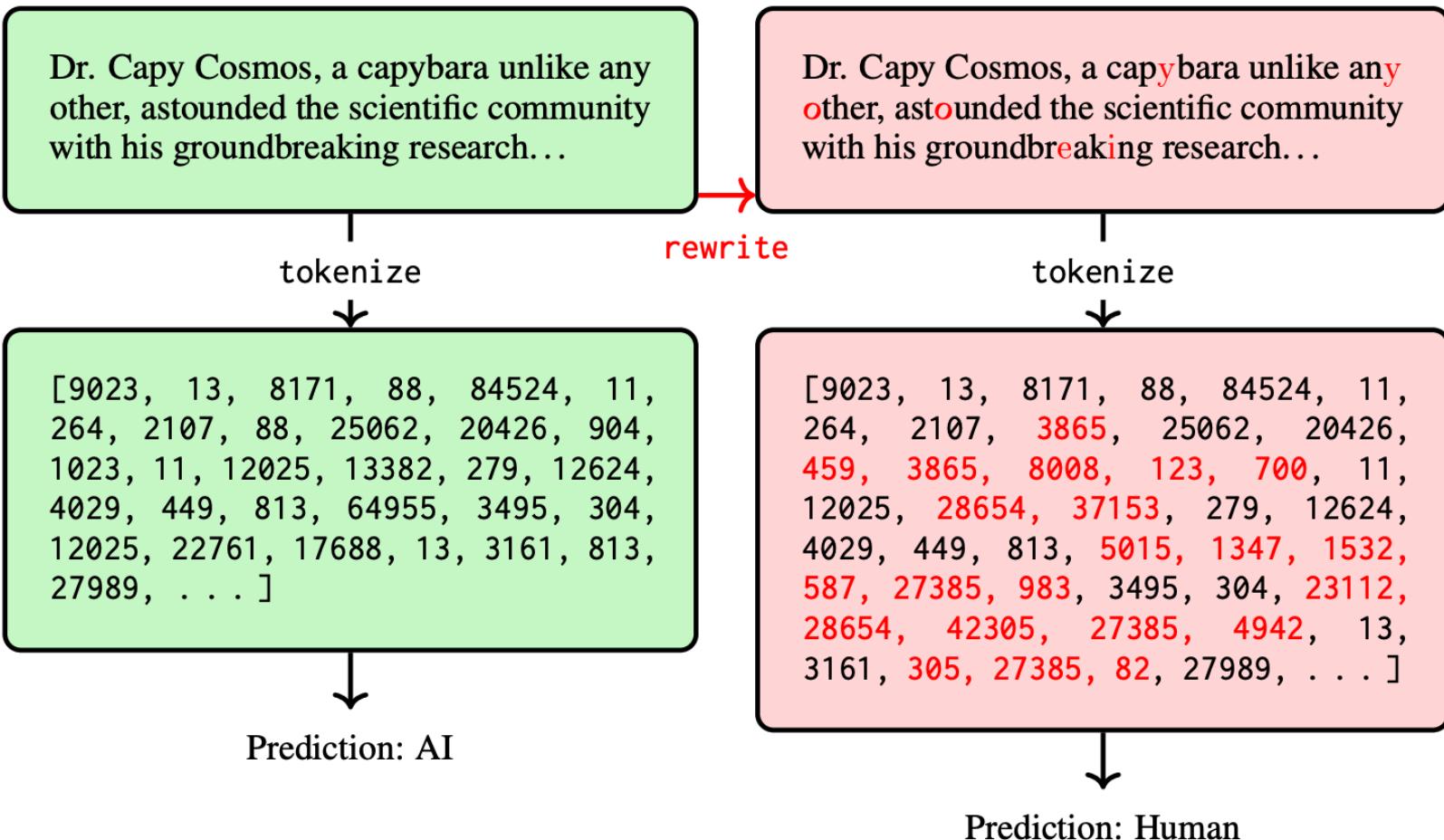
Dr. Capy Cosmos, a capybara unlike any other, astounded the scientific community with his groundbreaking research...

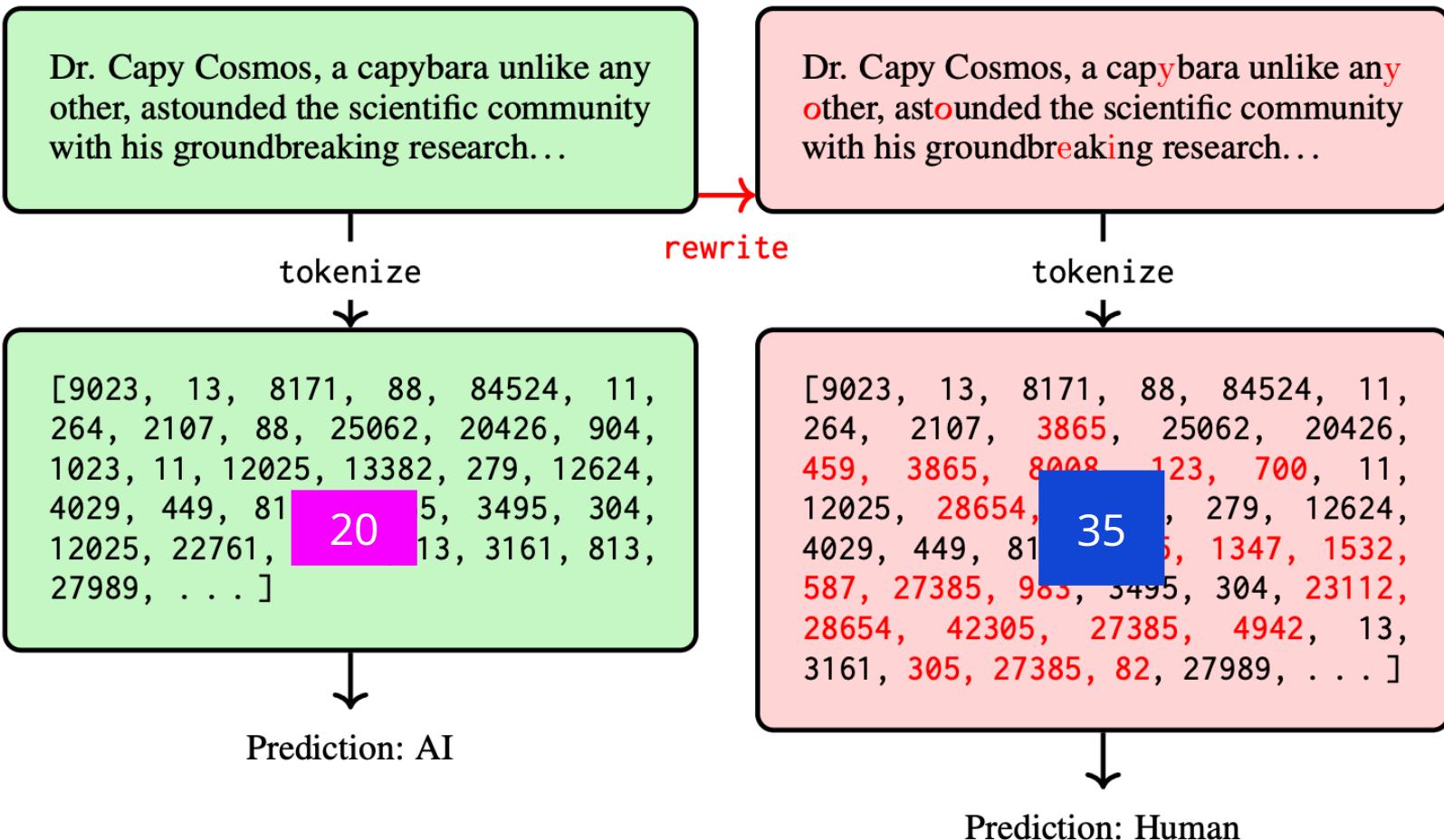
Dr. Capy Cosmos, a capybara unlike any other, astounded the scientific community with his groundbreaking research...

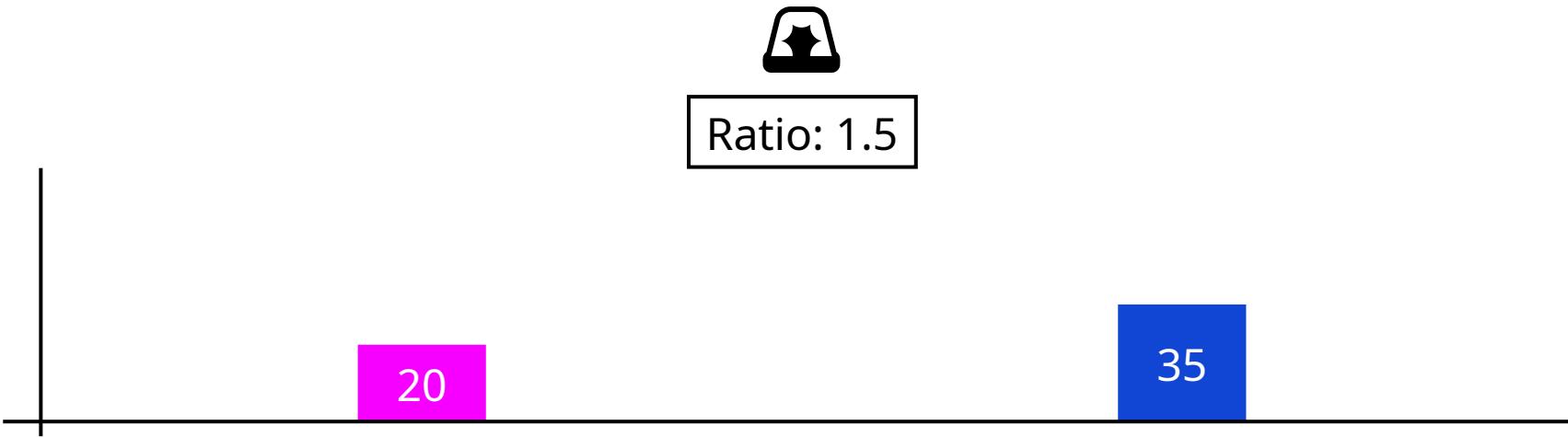






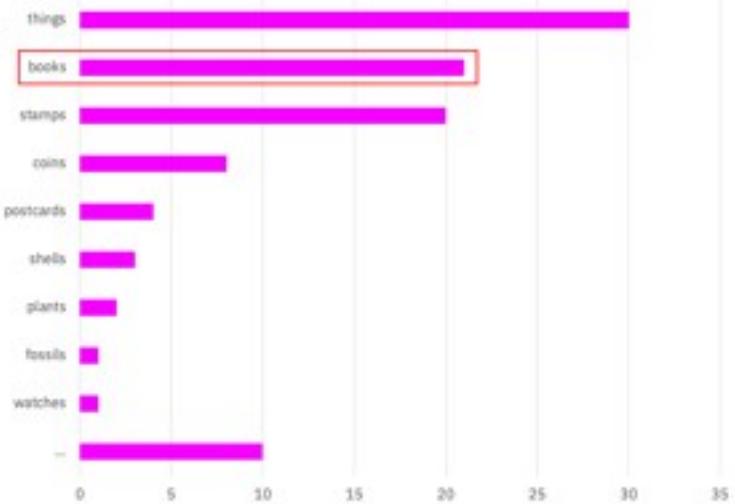




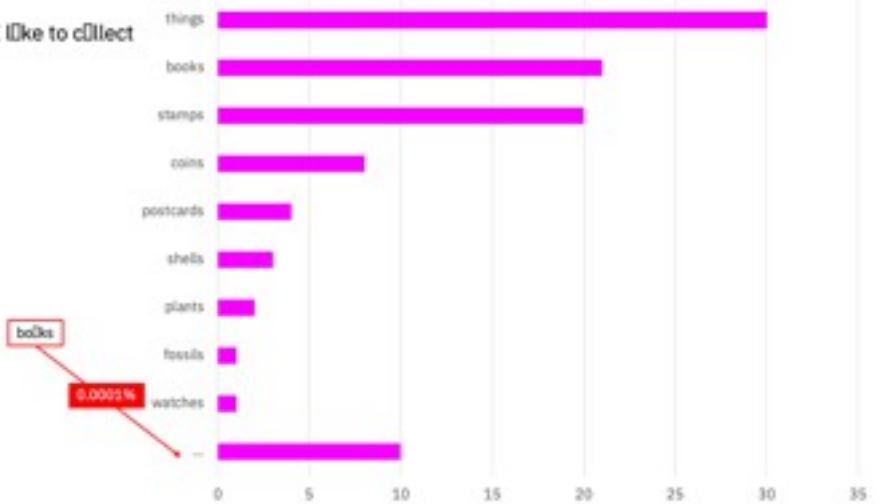


Idea 2: Perplexity

As a hobby, I like to collect

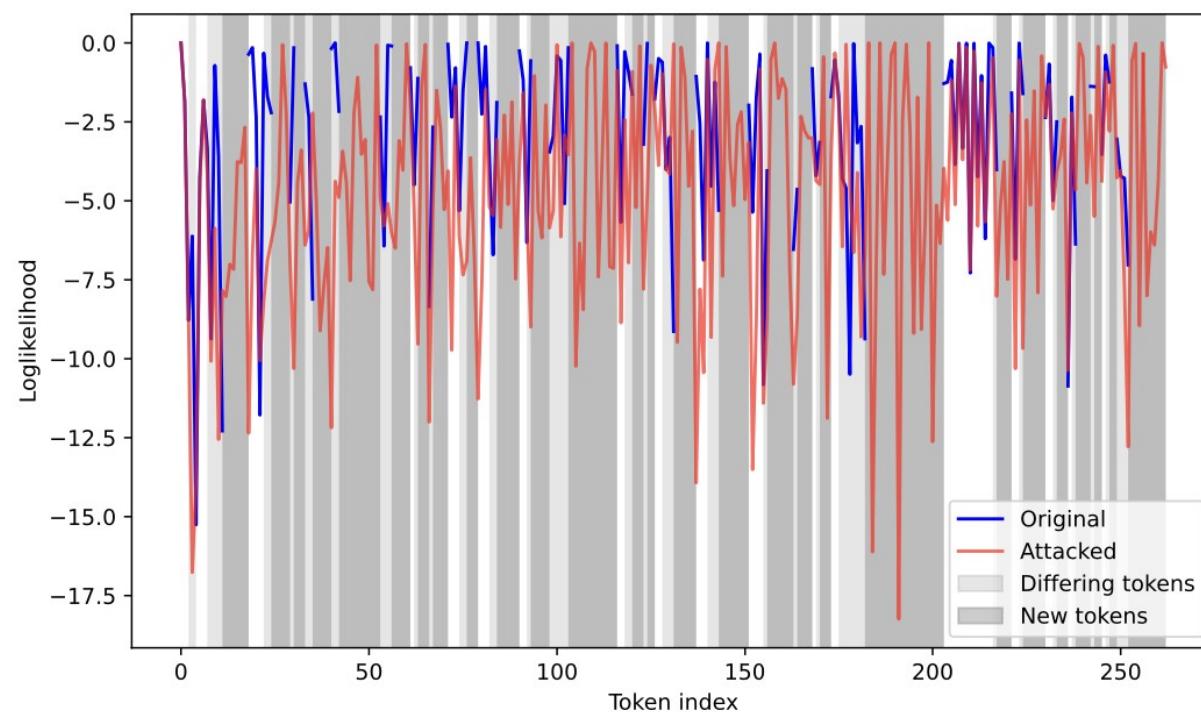


As a hobb0, I l0ke to c0llect

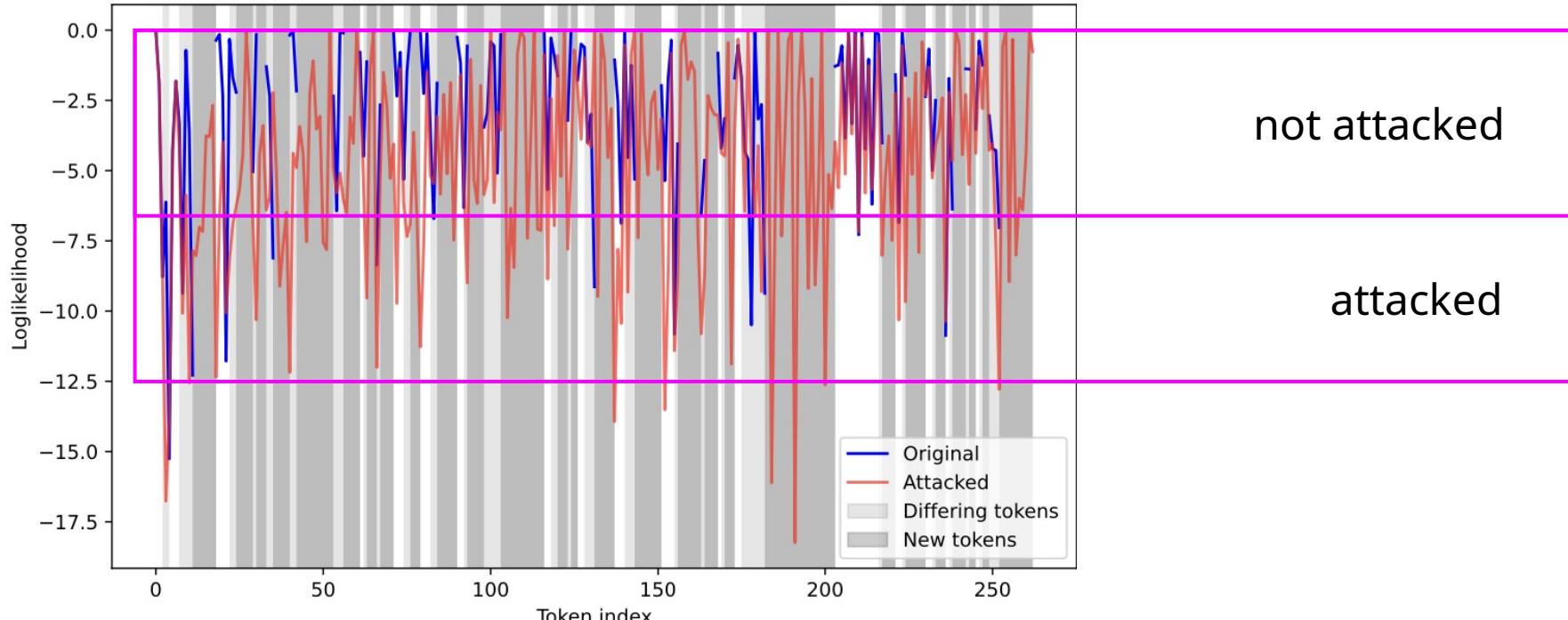


Dr. Capy Cosmos, a capybara unlike any other, astounded the scientific community with his groundbreaking research...

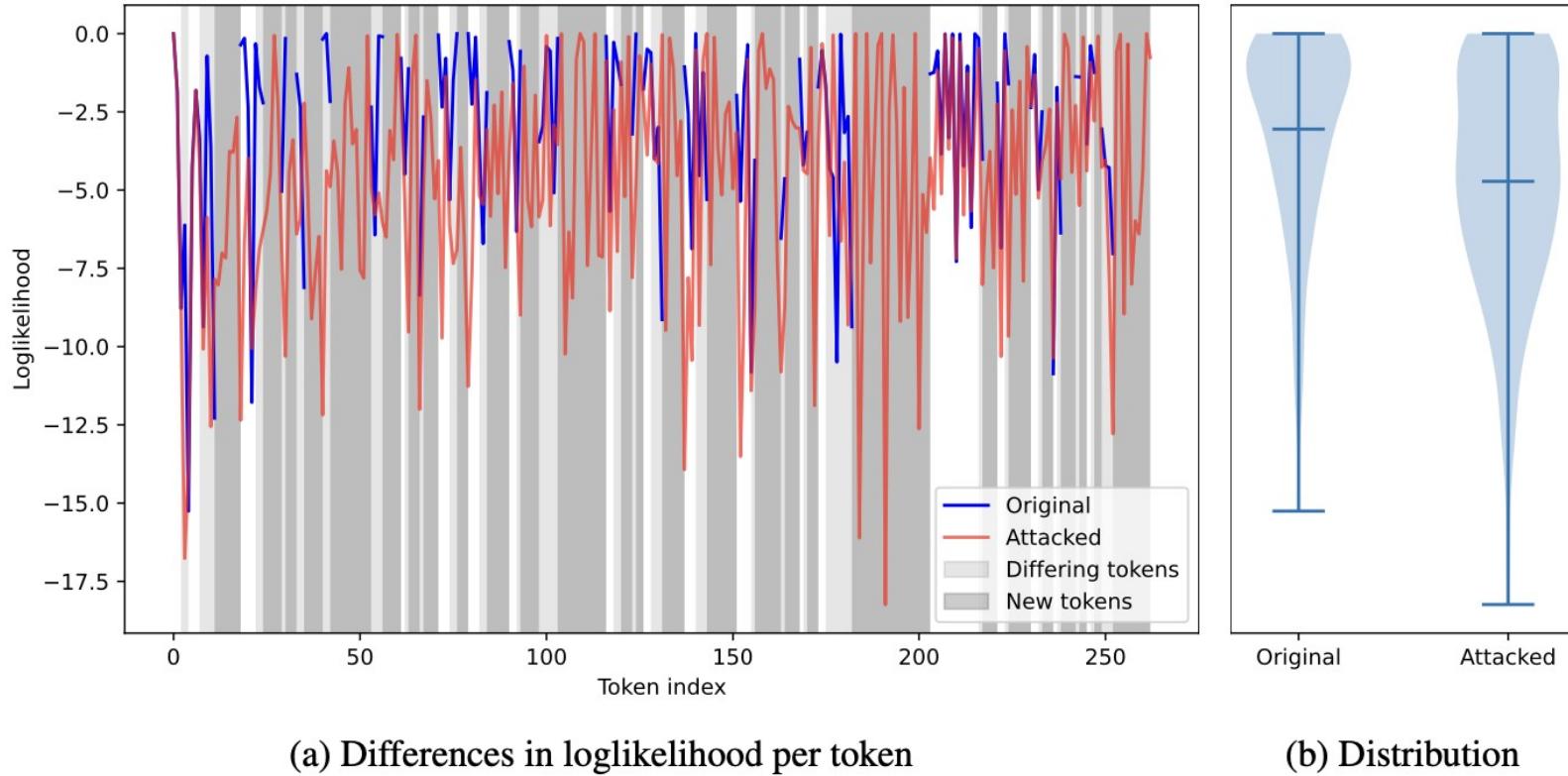
Dr. Capy Cosmos, a cap~~y~~bara unlike any other, astounded the scientific community with his groundbreaking research...



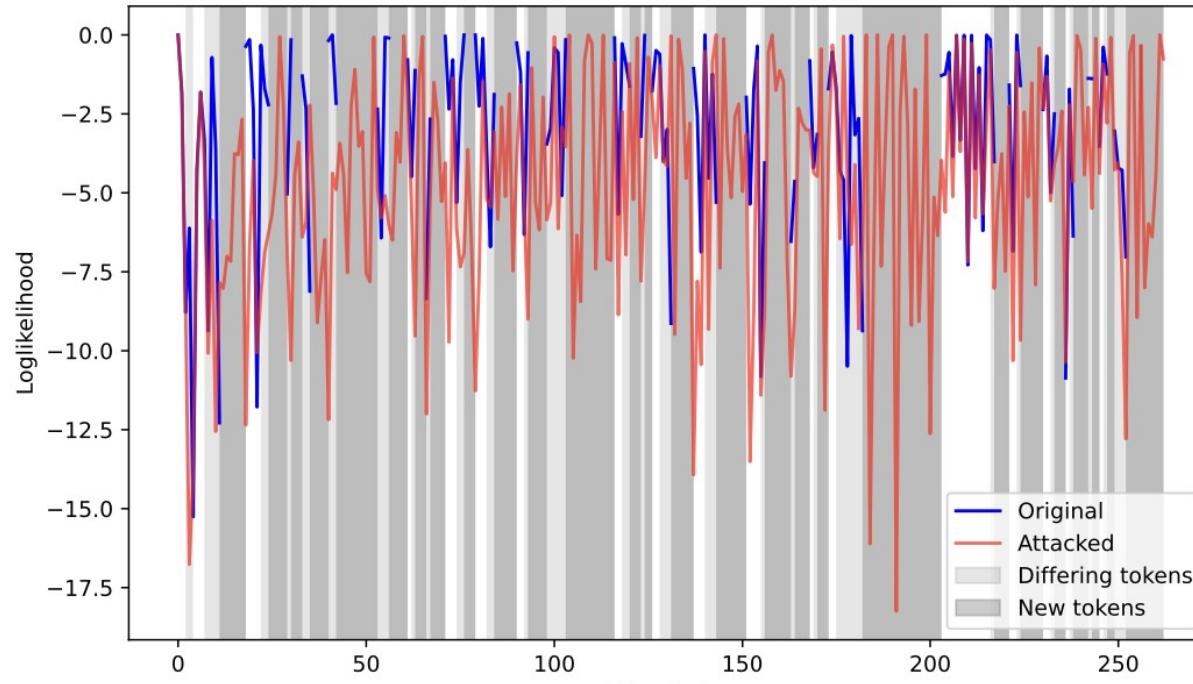
(a) Differences in loglikelihood per token



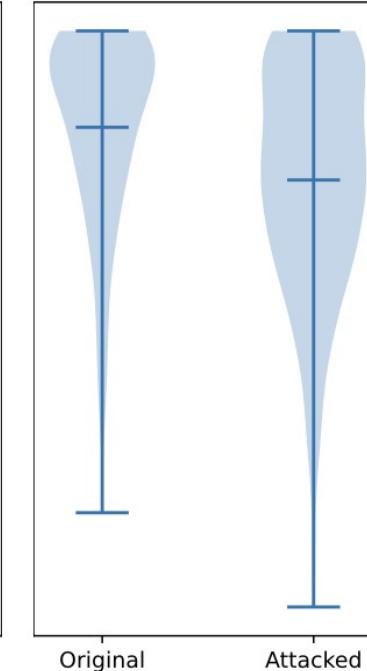
(a) Differences in loglikelihood per token



Perplexity values are not enough

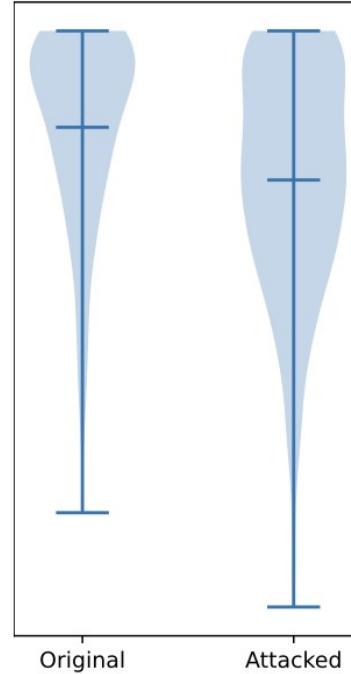


(a) Differences in loglikelihood per token



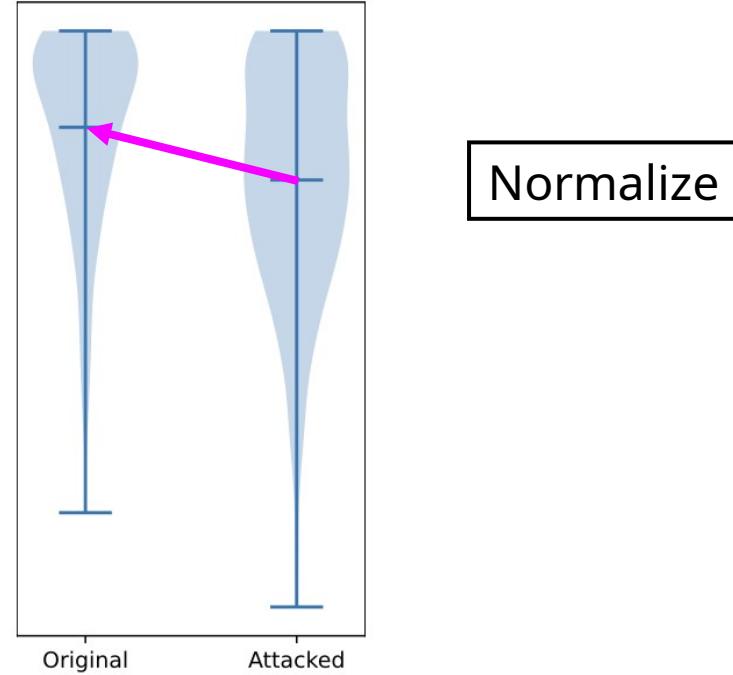
(b) Distribution

Perplexity values are not enough



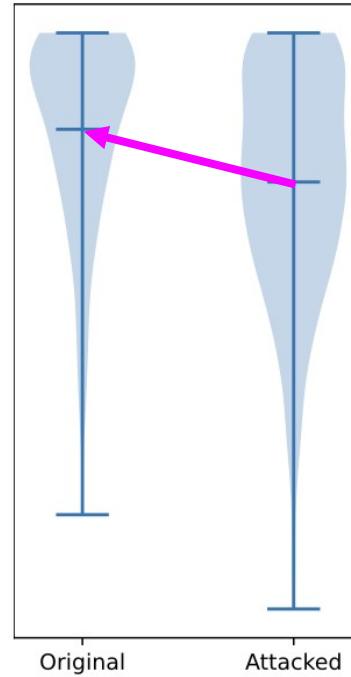
(b) Distribution

Perplexity values are not enough



(b) Distribution

Perplexity ratios are enough



Ratio: 1.3

(b) Distribution

US Constitution

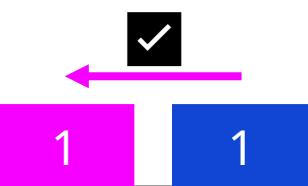
We the People of the
United States, in
Order to form a more
perfect Union,
establish Justice,
insure domestic
Tranquility...

1

Ratio: 1.3

US Constitution

We the People of the
United States, in
Order to form a more
perfect Union,
establish Justice,
insure domestic
Tranquility...



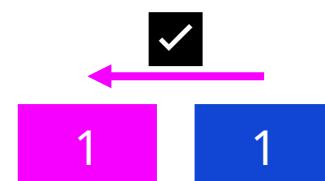
Ratio: 1.3

US Constitution

We the People of the United States, in Order to form a more perfect Union, establish Justice, insure domestic Tranquility...

US Constitution attack

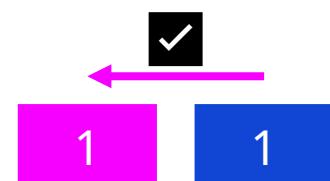
We the People of the United States, in Order to form a more perfect Union, establish Justice, insure domestic Tranquility...



Ratio: 1.3

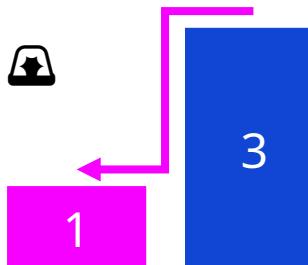
US Constitution

We the People of the United States, in Order to form a more perfect Union, establish Justice, insure domestic Tranquility...



US Constitution attack

We the People of the United States, in Order to form a more perfect Union, establish Justice, insure domestic Tranquility...

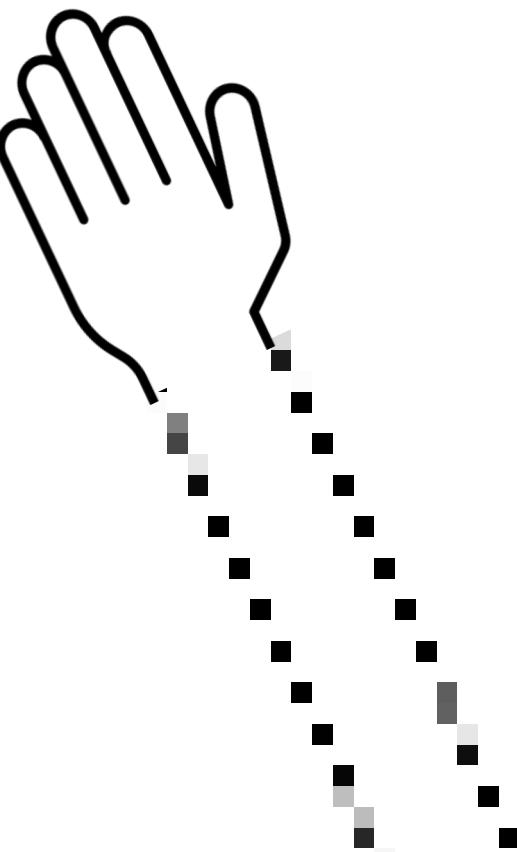


Ratio: 1.3

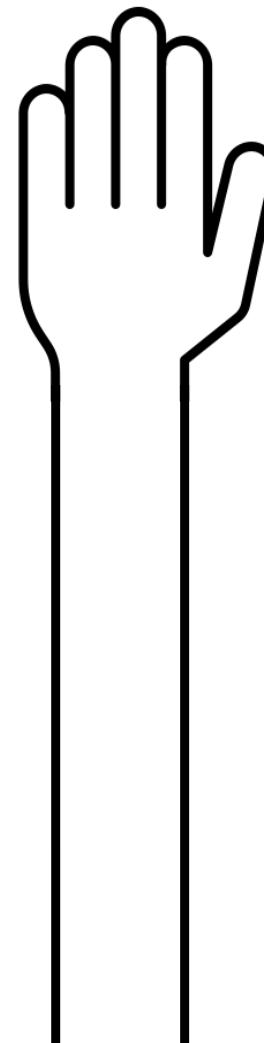
Demo
!

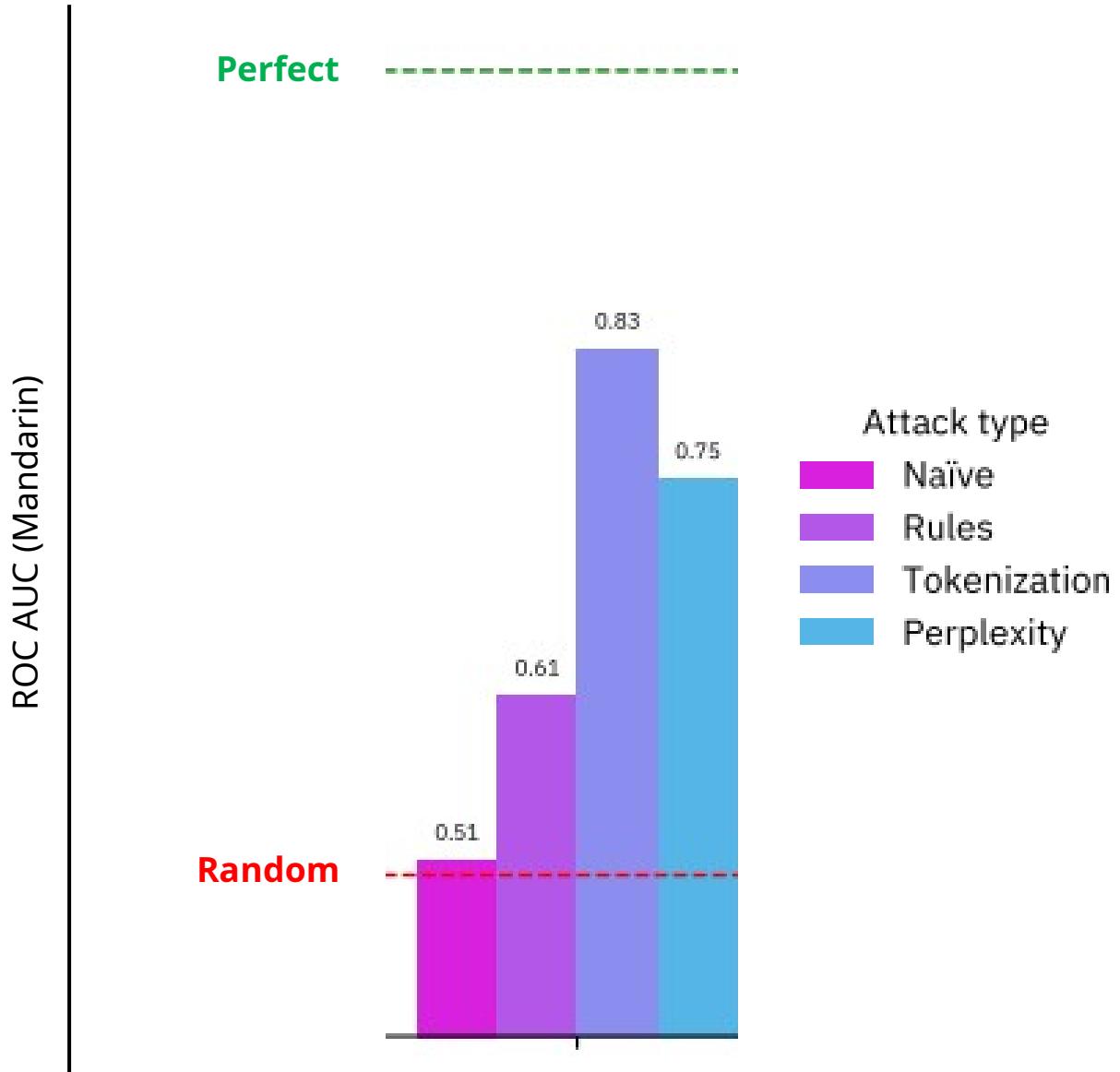
Tokenization or Perplexity?

Tokenization

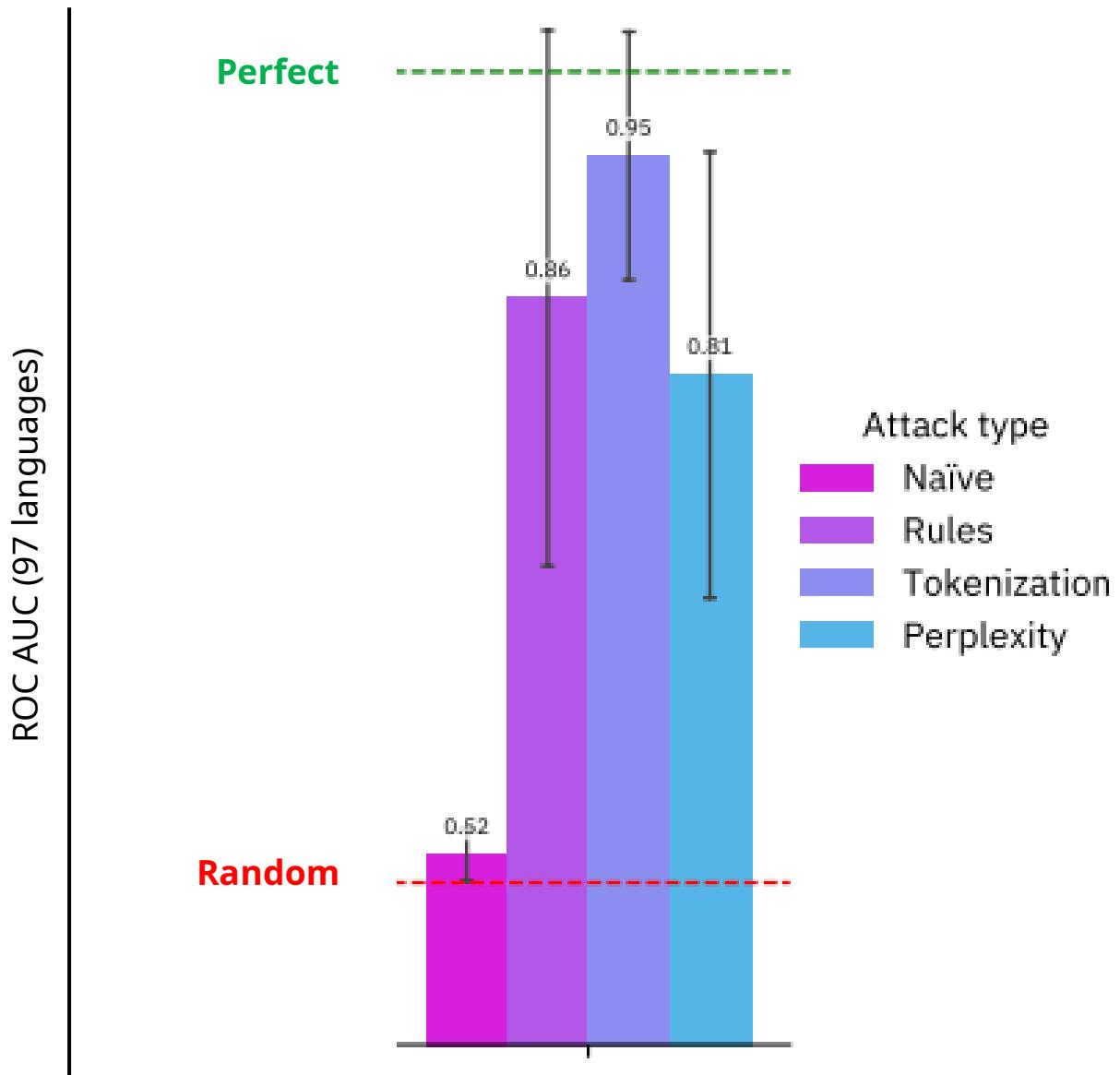


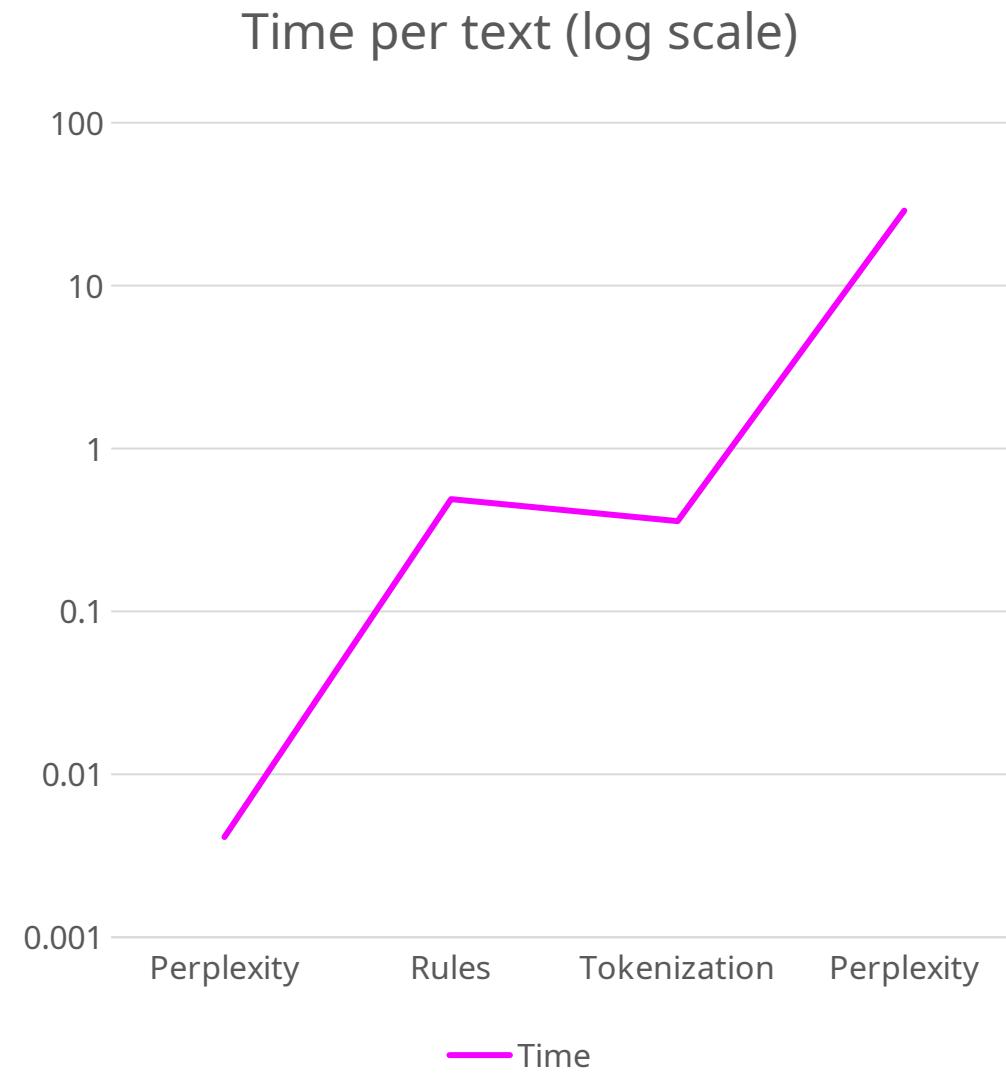
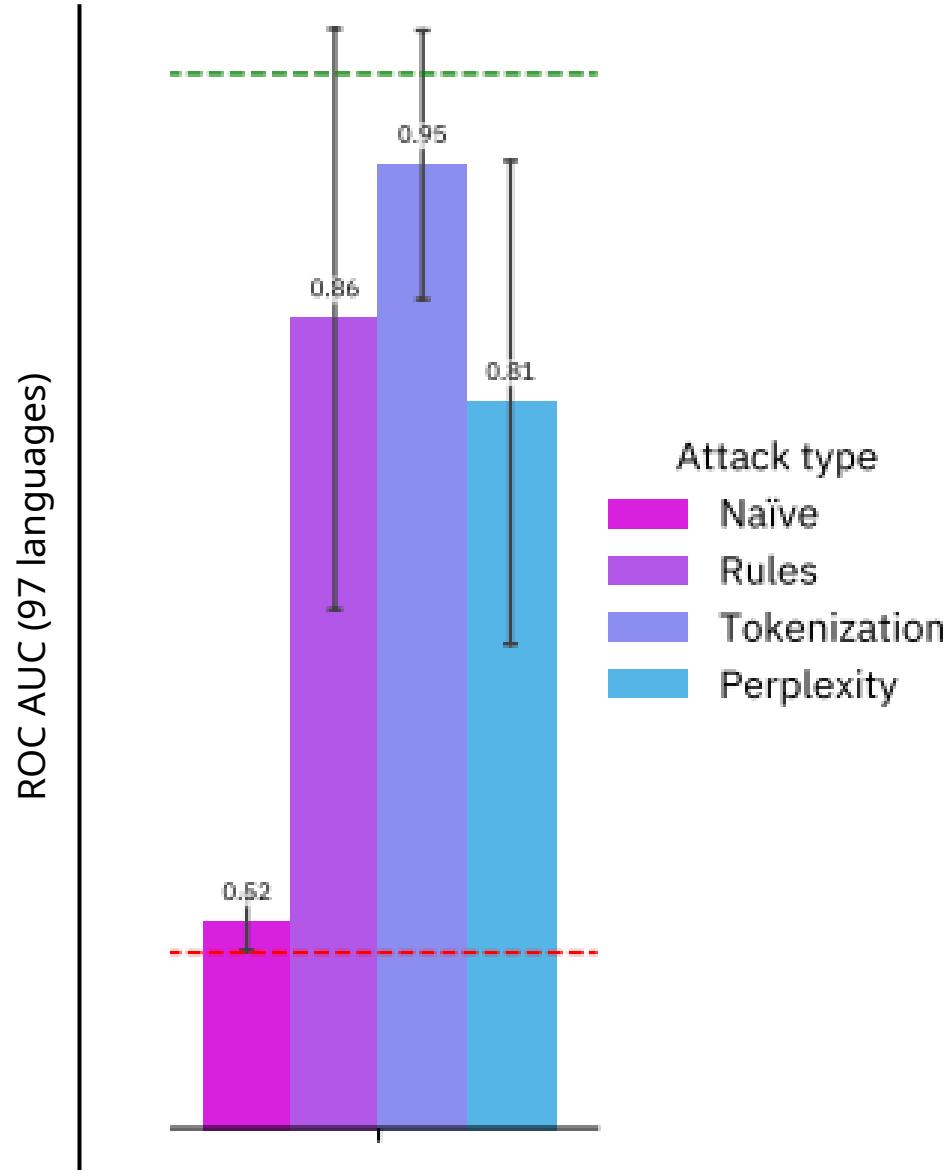
Perplexity

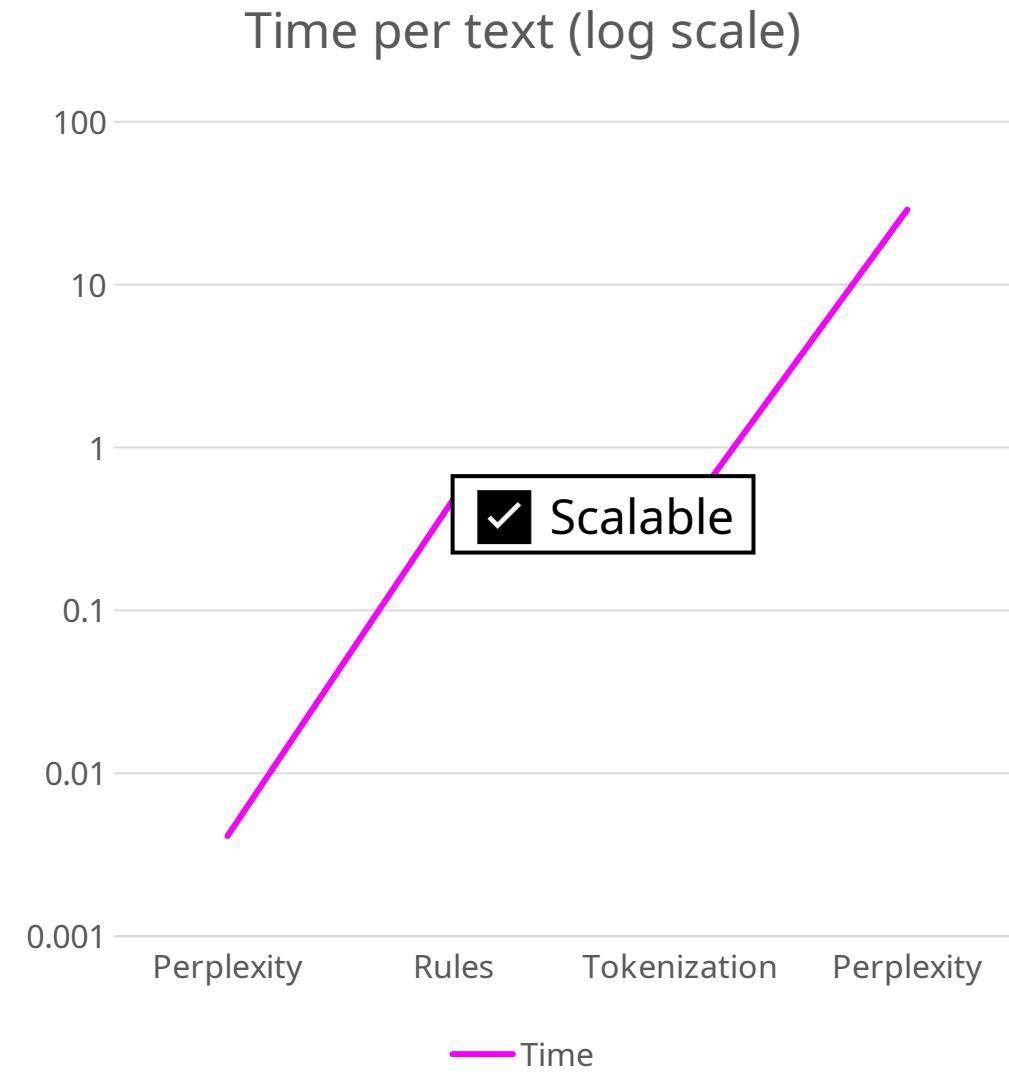
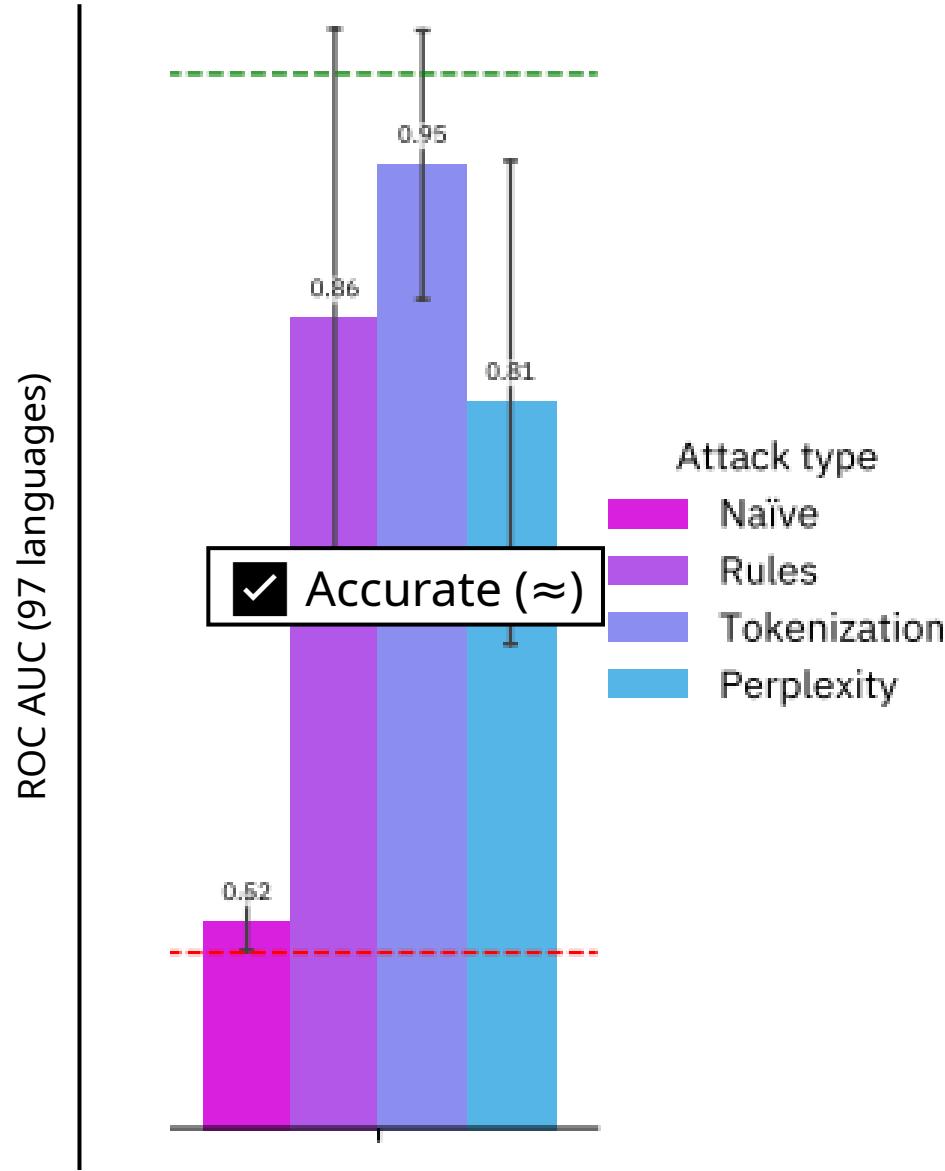


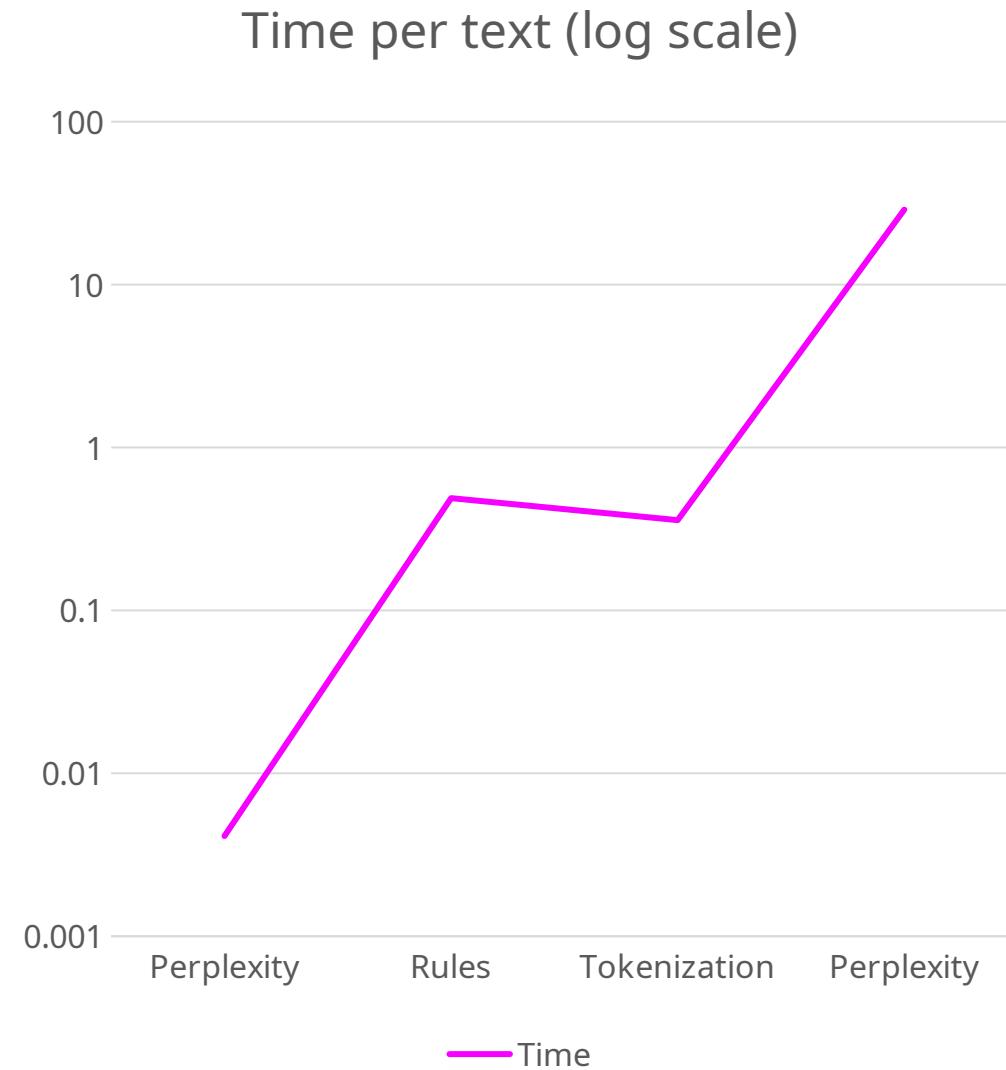
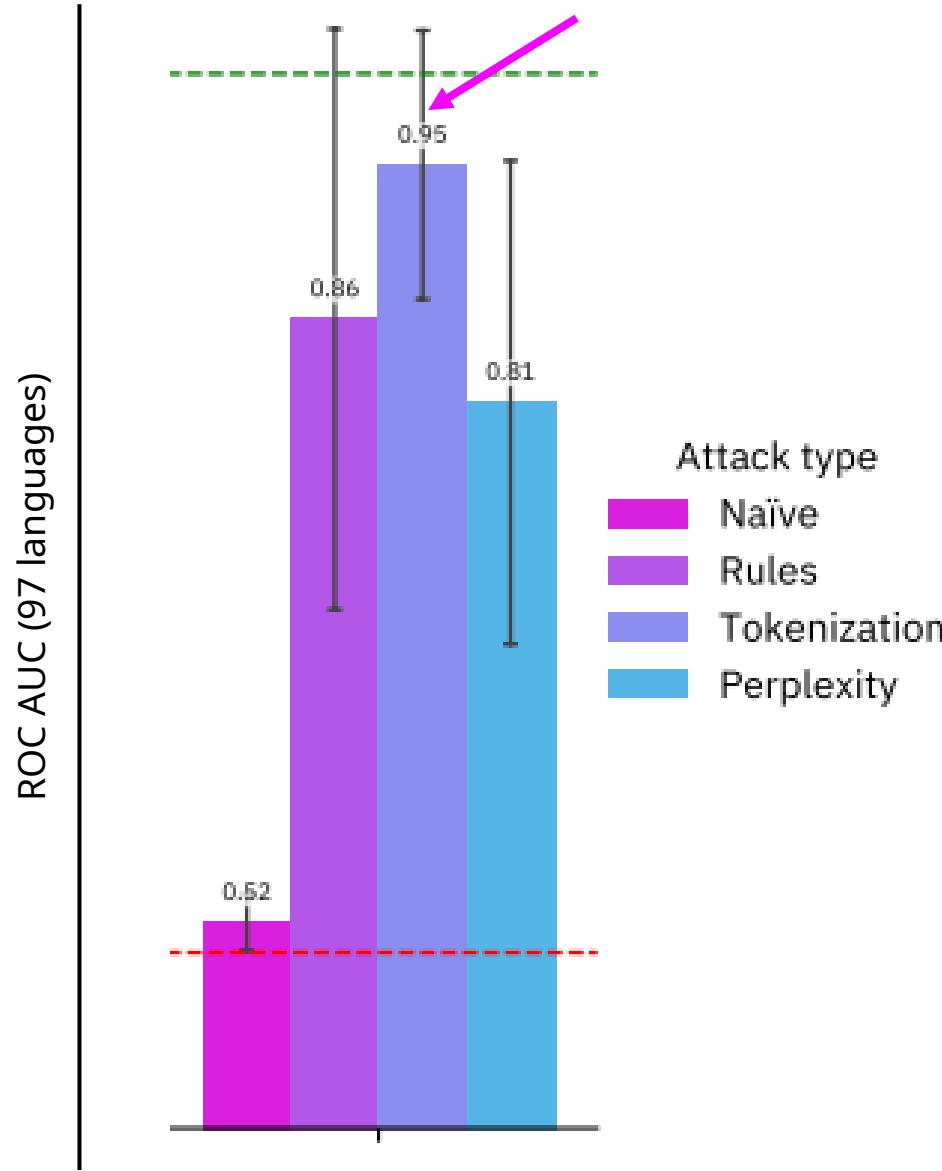


Does this extend to other
languages?







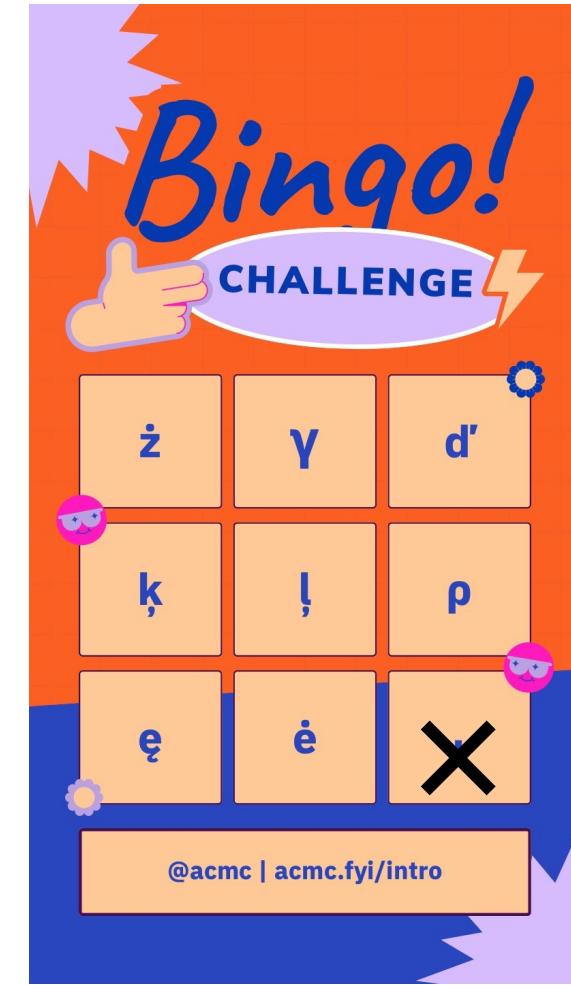


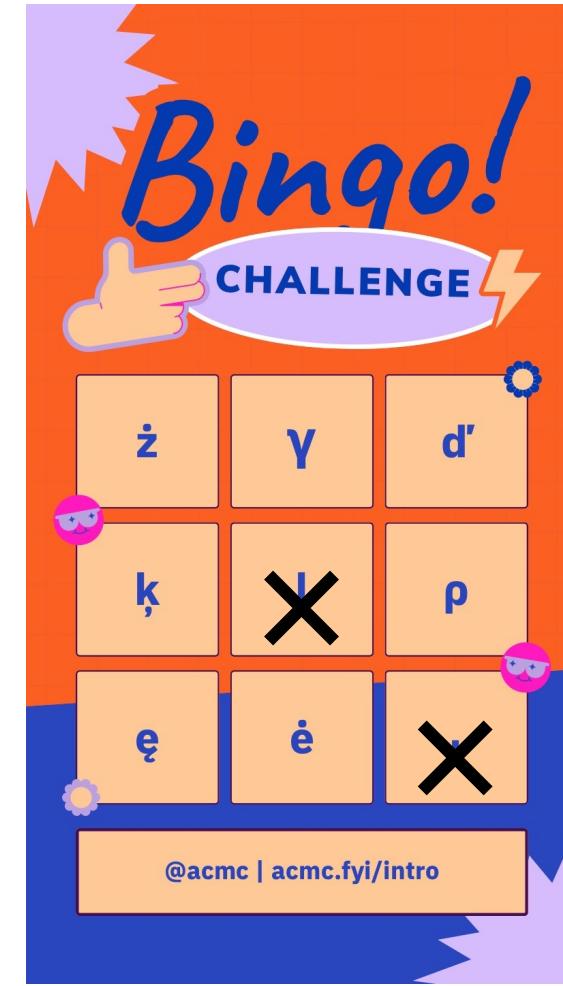
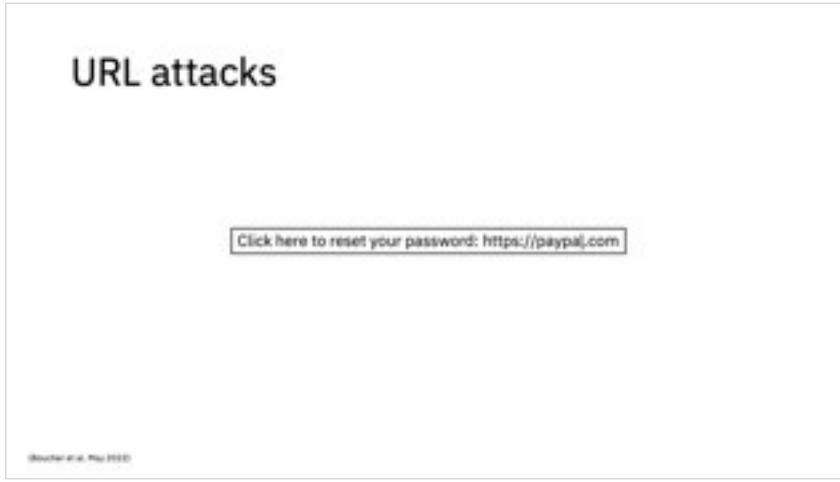
Work to do is on us – let's keep innovating!

Bingo!

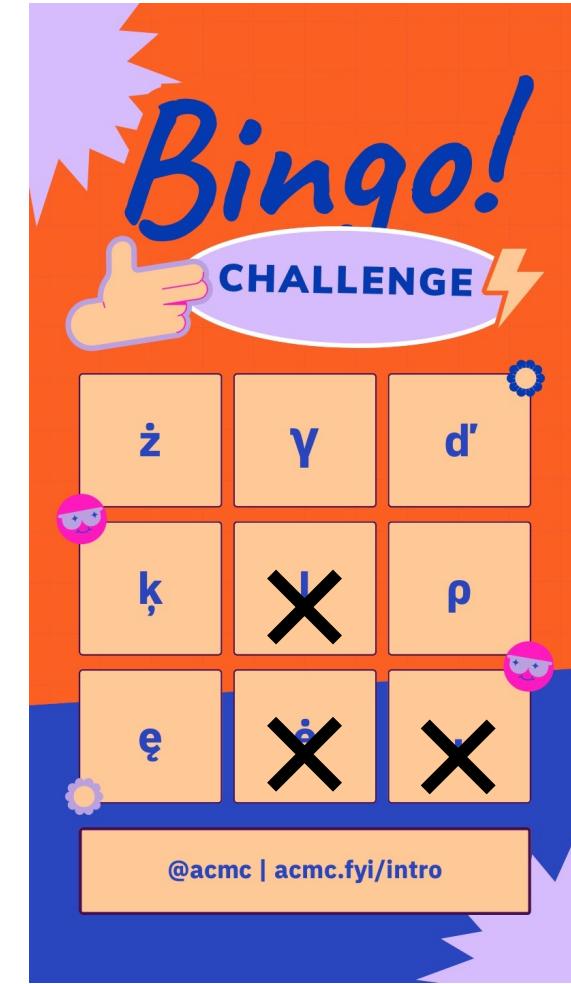
Steganography

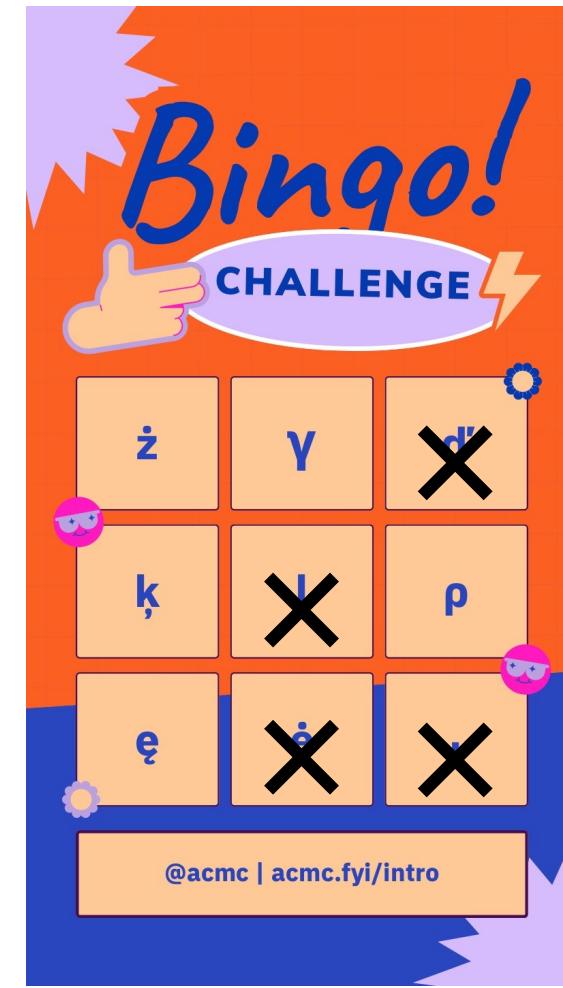
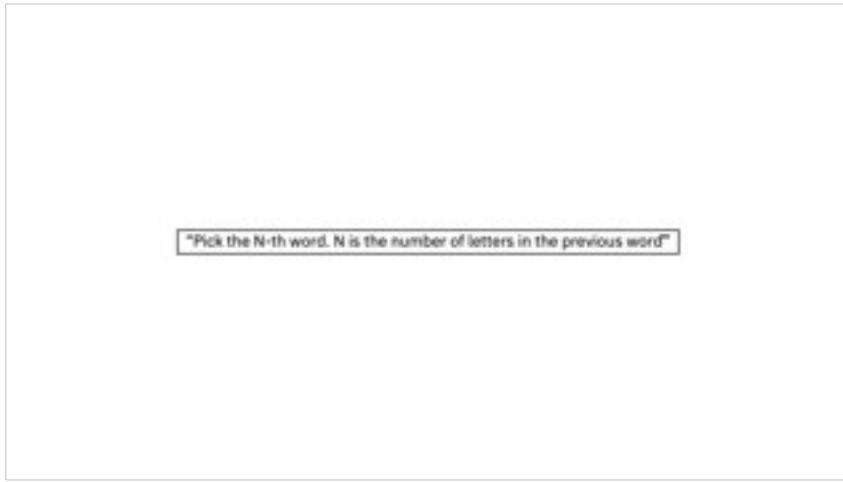
Hiding data in plain sight
to avoid detection

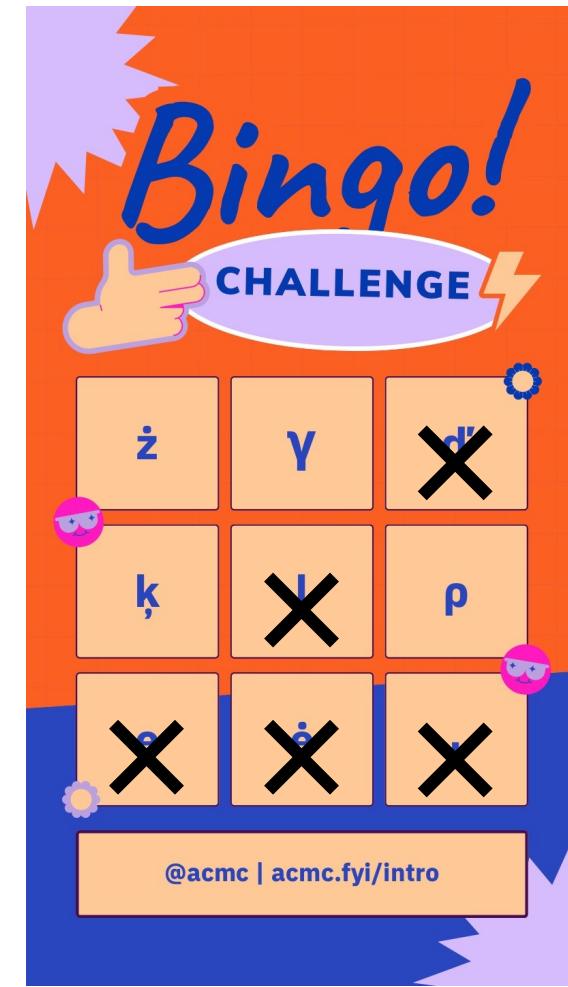




Perplexity-based detectors

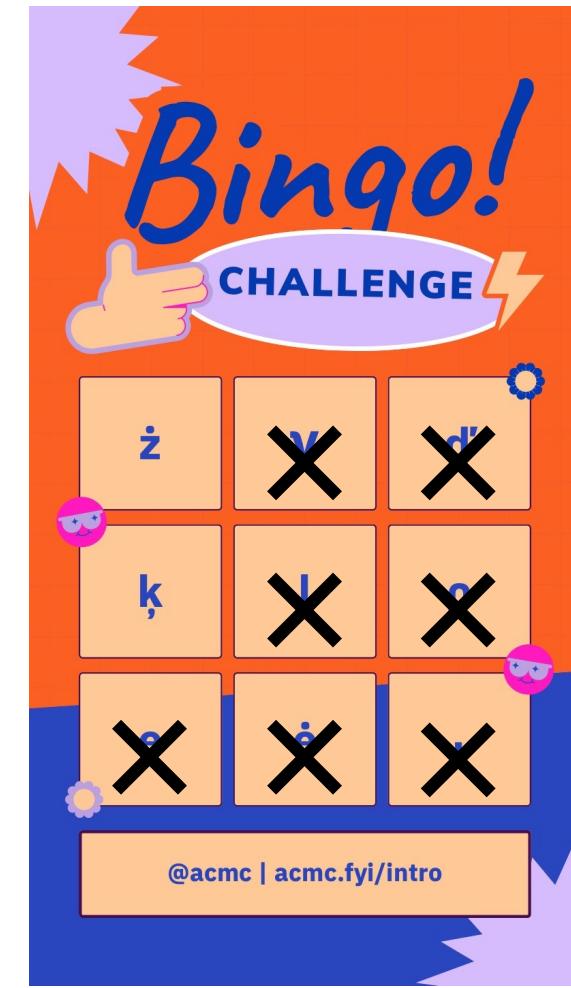
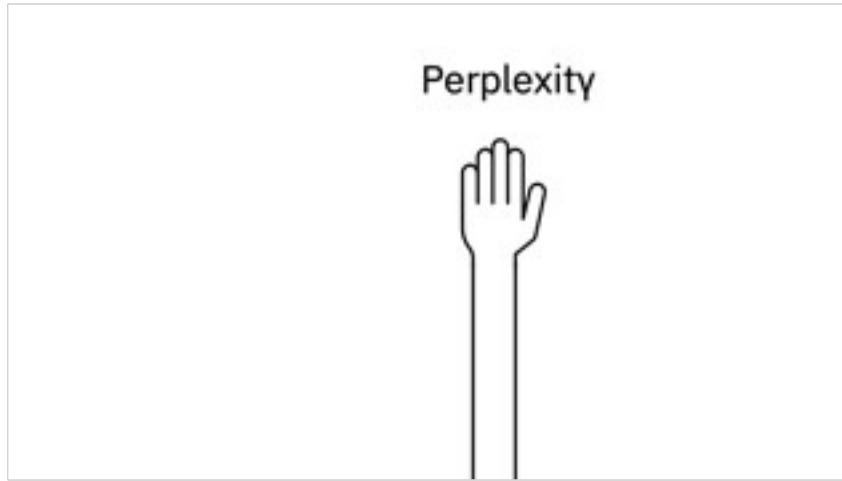






Idea 2: Perplexity





Recap

Homoglyphs: characters we can't discern

Homoglyphs: characters we can't discern

Detection of AI-generated texts: different techniques

Homoglyphs: characters we can't discern

Detection of AI-generated texts: different techniques

Homoglyph-based attacks: text becomes **unrecognizable**

Homoglyphs: characters we can't discern

Detection of AI-generated texts: different techniques

Homoglyph-based attacks: text becomes unrecognizable

Different mechanisms of action – all exploit “confusion”

Homoglyphs: characters we can't discern

Detection of AI-generated texts: different techniques

Homoglyph-based attacks: text becomes unrecognizable

Different mechanisms of action – all exploit “confusion”

Highly effective, renders all detectors **ineffective**

Homoglyphs: characters we can't discern

Detection of AI-generated texts: different techniques

Homoglyph-based attacks: text becomes unrecognizable

Different mechanisms of action – all exploit “confusion”

Highly effective, renders all detectors ineffective

Lower access barrier: increased **risk**

Homoglyphs: characters we can't discern

Detection of AI-generated texts: different techniques

Homoglyph-based attacks: text becomes unrecognizable

Different mechanisms of action – all exploit “confusion”

Highly effective, renders all detectors ineffective

Lower access barrier: increased risk

Some ideas we can try

Homoglyphs: characters we can't discern

Detection of AI-generated texts: different techniques

Homoglyph-based attacks: text becomes unrecognizable

Different mechanisms of action – all exploit “confusion”

Highly effective, renders all detectors ineffective

Lower access barrier: increased risk

Some safeguards we can try

Work to do is on us – let's keep **innovating!**

Work to do is on us – let's keep innovating!



*get the bibliography & connect!

<https://forms.gle/JzkJTLGNNiCAt4h99>

