

# Beyond the Parameters:

## Measuring Actual Privacy in Obfuscated Texts

IIR 2024

14th Italian Information Retrieval Workshop

September 5-6, 2024 - Udine, Friuli Venezia Giulia, Italy

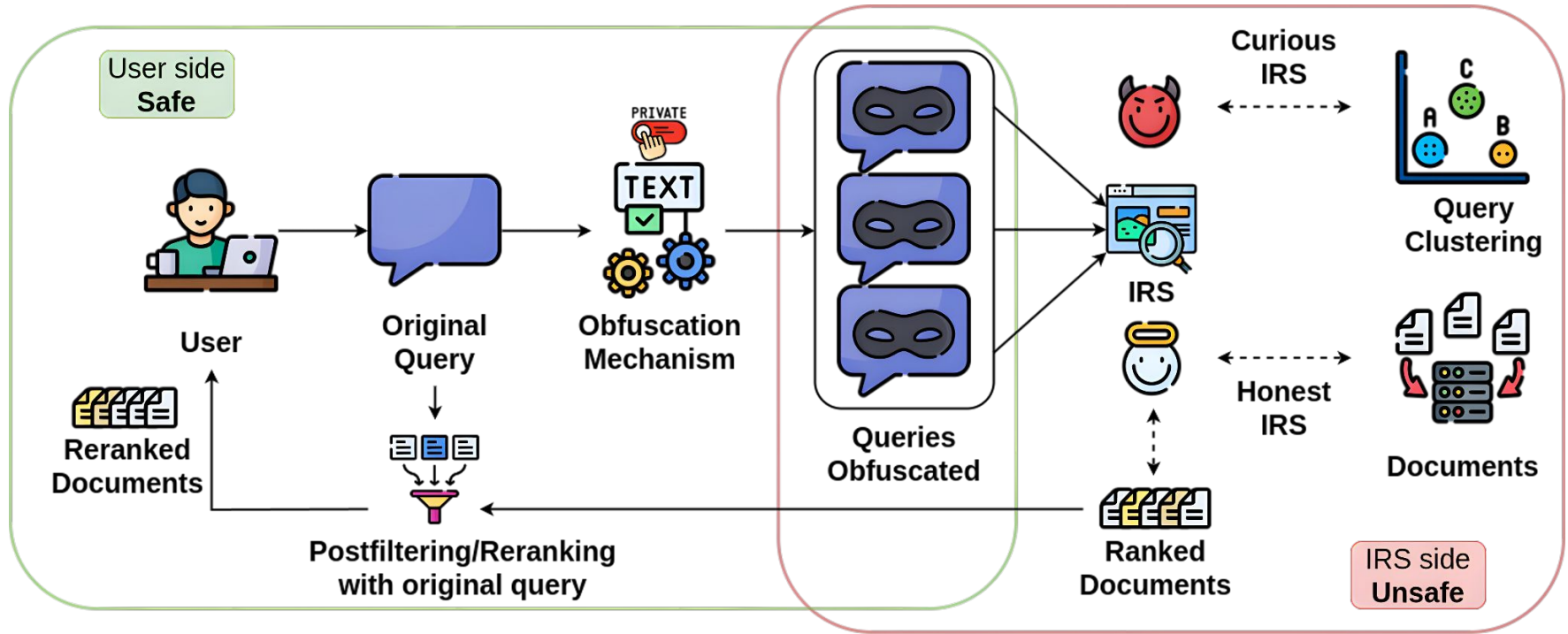
**Francesco L. De Faveri**

Guglielmo Faggioli

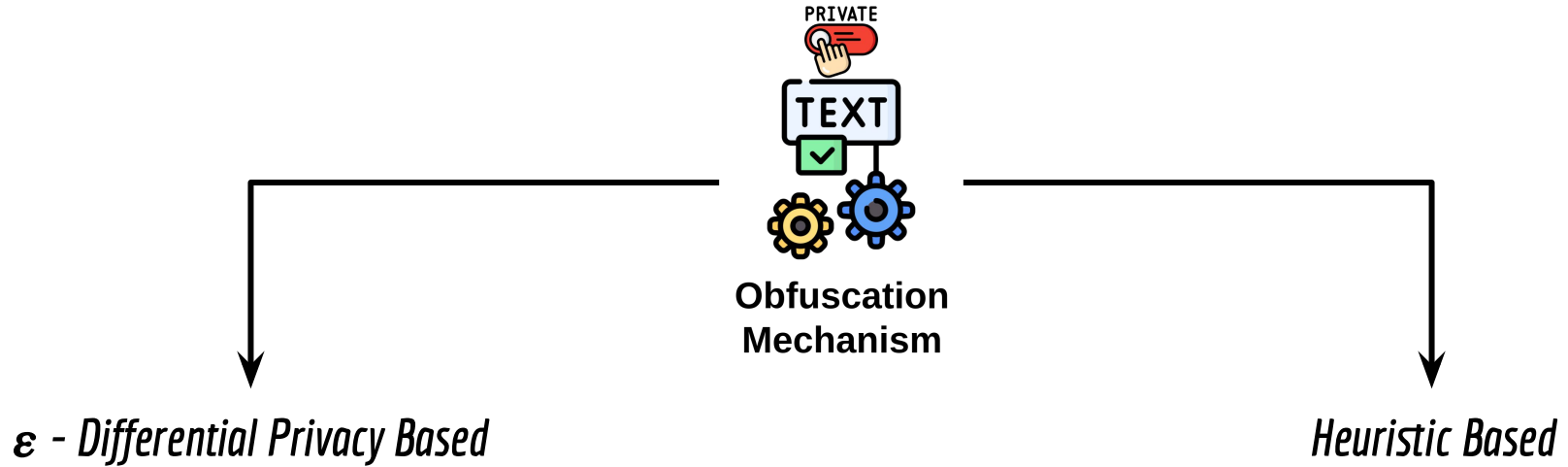
Nicola Ferro



# IR & Privacy



# Generating the obfuscated queries



E.g. Using the CMP mechanism with  $\epsilon = 12.5$ :

"do goldfish grow"  $\rightarrow$  "do xlvi grow", "host frangieh expands", "do goldfish grow"

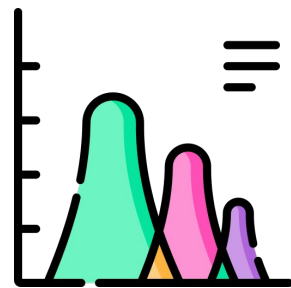
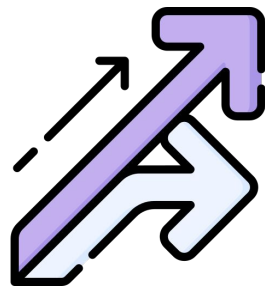
# $\epsilon$ - Differential Privacy obfuscation

Considering any pair neighbouring datasets\*,  $D$  and  $D'$ , a **privacy budget**  $\epsilon \in \mathbb{R}^+$ , a mechanism is  $\epsilon$  - Differentially Private if it holds:

$$\Pr [\mathcal{M}(D) \in S] \leq e^\epsilon \cdot \Pr [\mathcal{M}(D') \in S] \quad \forall S \subset \text{Im}(\mathcal{M})$$

Applied to **textual data**:

- Embedding Perturbation: CMP (2020), Mahalanobis (2020), Vickrey (2021)
- Sampling: CusText (2023), SanText (2021), TEM (2023)



\*Datasets that differs for at most one record.

# Heuristic based obfuscation

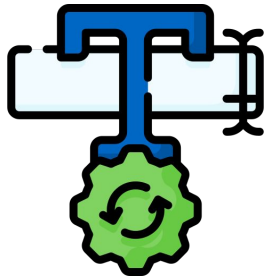
Arampatzis et al. (2013)

Obfuscate original query terms with **synonyms**, **hypernyms**, and **holonyms** from Wordnet.

E.g.:

“Cat” → “Feline”

“Cancer” → “Disease”

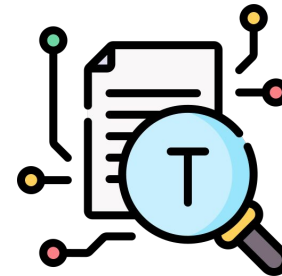


Fröbe et al. (2021)

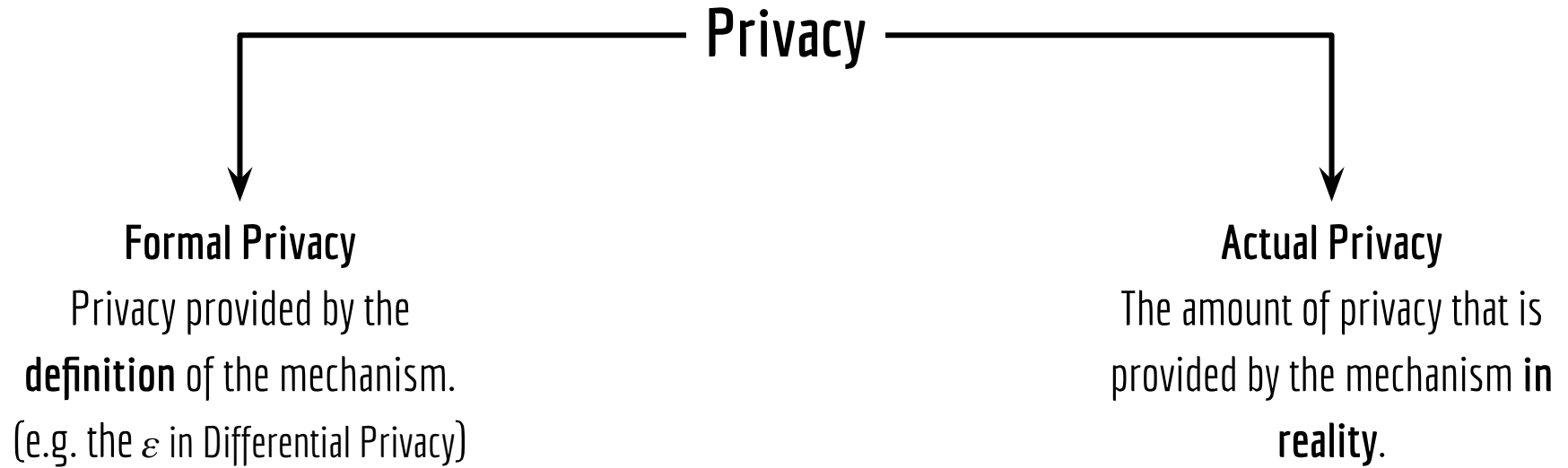
Generates **keyword** queries using a small local corpus to *obfuscate the original information need of the user*.

E.g.:

“A user wants health advice while hiding a potential disease.” submits to the IRS queries like “**lower heart rate**”, “**forearm pain**”, “**symptoms heart attack**”.



# What is the privacy?



# How to measure privacy?

Privacy

## Formal Privacy

Privacy provided by the **definition** of the mechanism.  
(e.g. the  $\epsilon$  in Differential Privacy)

## Actual Privacy

The amount of privacy that is provided by the mechanism **in reality**.



# Strategies to measure actual Privacy

Technique	Measure	PROs	CONs
<i>Translation Based</i>	<i>BLEU</i>	<ul style="list-style-type: none"> <li>- Effective short n-grams comparison</li> <li>- Computationally efficient</li> </ul>	<ul style="list-style-type: none"> <li>- No semantic similarity</li> <li>- Long distance dependences</li> </ul>
	<i>ROUGE</i>	<ul style="list-style-type: none"> <li>- n-gram comparison</li> <li>- Computationally efficient</li> </ul>	<ul style="list-style-type: none"> <li>- No semantic similarity</li> <li>- Long distance dependences</li> </ul>
	<i>METEOR</i>	<ul style="list-style-type: none"> <li>- Synonyms and stemming</li> <li>- Paraphrasing</li> </ul>	<ul style="list-style-type: none"> <li>- Heuristic-based</li> <li>- Computational expensive</li> </ul>
<i>Lexical Based</i>	<i>Jaccard Index</i>	<ul style="list-style-type: none"> <li>- Lexical similarity</li> <li>- Computationally efficient</li> </ul>	<ul style="list-style-type: none"> <li>- No semantic similarity</li> <li>- No synonyms or paraphrasing</li> </ul>
	$N_w$ and $S_w$	<ul style="list-style-type: none"> <li>- Lexical similarity</li> <li>- Measure of mechanism failure</li> </ul>	<ul style="list-style-type: none"> <li>- No semantic similarity</li> <li>- Easy to deceive</li> </ul>
<i>Contextual Based</i>	<i>BERTScore</i>	<ul style="list-style-type: none"> <li>- Semantic similarity</li> <li>- Long distance dependences</li> </ul>	<ul style="list-style-type: none"> <li>- Pre-trained model dependent</li> <li>- Computational expensive</li> </ul>
	<i>Transformers Sentence Embeddings Similarity</i>	<ul style="list-style-type: none"> <li>- Semantic similarity</li> <li>- Long distance dependences</li> </ul>	<ul style="list-style-type: none"> <li>- Pre-trained model dependent</li> <li>- Computational expensive</li> </ul>



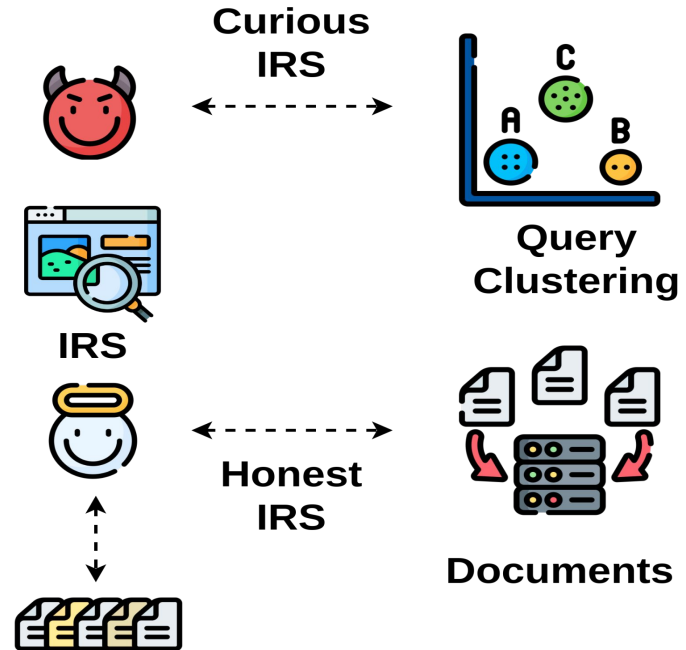
# Adversarial Risk of Breaking Privacy

How can we model the risk?

Risk = Probability of the attack x Vulnerability x Impact



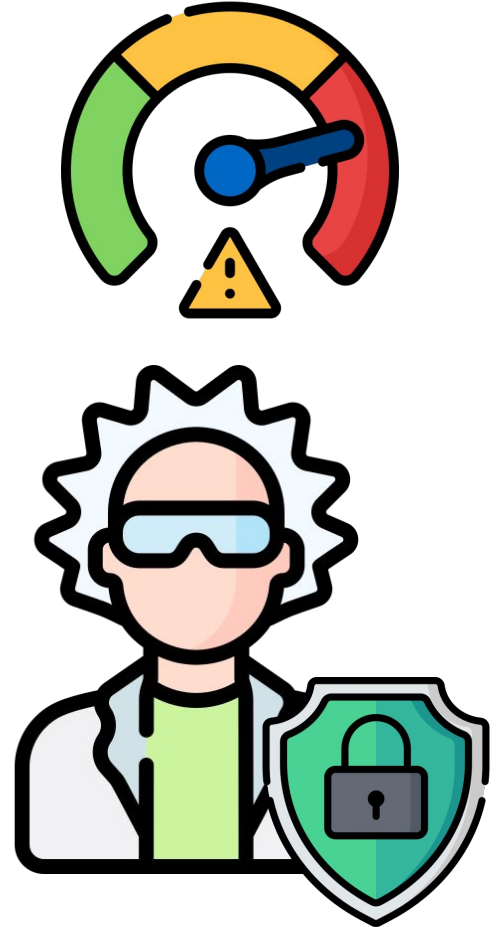
Can the privacy metrics be **good indicators** for such risk?

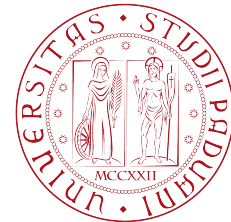


IRS side  
Unsafe

# Future directions

- Are the current measure enough to evaluate privacy?
- Do we need “human assessment” to evaluate privacy (e.g. Privacy relevance judgments) ?
- How can we understand if the measure is a good proxy for the probability of success for a class of attacks?



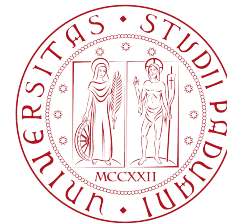


# Beyond the Parameters:

## Measuring Actual Privacy in Obfuscated Texts

Thanks for the attention!  
Question Time

Francesco L. De Faveri<sup>1</sup>  
Guglielmo Faggioli<sup>1</sup>  
Nicola Ferro<sup>1</sup>



# Backup Slides

## Beyond the Parameters:

Measuring Actual Privacy in  
Obfuscated Texts

Francesco L. De Faveri<sup>1</sup>  
Guglielmo Faggioli<sup>1</sup>  
Nicola Ferro<sup>1</sup>

<sup>1</sup>: Department of Information Engineering, University of Padua, Padova, Italy

# Backup 1 - Differential Privacy Mechanisms

Strategy	Mechanism	Params	Description
<i>Embedding</i>	CMP	-	The noise is sampled from an $n$ - dimensional Laplace distribution of scale $\frac{1}{\epsilon}$ .
	Mhl	$\lambda$	The noise is sampled from an $n$ - dimensional Normal distribution defined by the $\lambda$ regularized Mahalanobis norm of the term embedding.
	Vickrey	$t, \lambda$	The noise is sampled as defined by the parent mechanism (CMP or Mhl) and the obfuscation term is set based on a free parameter $t$ .
<i>Sampling</i>	CusText	$K$	Sampling of a new term is bounded to $K$ possible terms picked using the scores computed using the distances among word embeddings.
	SanText	-	Sampling of a new term is computed with a score based on the distances among embeddings, with terms closer to the obfuscation having a higher probability.
	TEM	$\beta$	Noise sampled from an $n$ - dimensional Gumbel distribution is added to the scores, and the final obfuscation term is sampled accordingly to the maximum noisy score.