



# Detecting Vague Clauses in Italian Privacy Policies using Transformers, LLMs, and Cross-Lingual Techniques

**Giulia Grundler, Mariaceleste Musicco,**

Andrea Galassi, Francesca Lagioia, Ruta Liepina,

Giorgio Resta, Sara Roccu, Giovanni Sartor, Paolo Torroni



ALMA MATER STUDIORUM  
UNIVERSITA DI BOLOGNA

# Context

Growing **information asymmetry** between **digital service providers** and **consumers** about the process of their **personal data**



**Privacy Policies**  
should inform instead overwhelm

# Context

Multilingualism ensures **clarity and accountability** towards consumers and Data Protection Authorities

Multilingual legal LLMs **challenges:**

- Nature of Legal Language
- Lack of annotated datasets for less resource languages



# Objective and Contribution

**Automatically detect insufficiently informative clauses  
in PPs across multiple European languages**



- **New corpus of Italian** privacy policies
- Experiments with BERT models and LMM for detecting insufficiently informative clauses **on Italian PP**
- Experiments with multilingual techniques to **transfer knowledge** between Italian and English

# Hierarchical annotation method



Detecting clauses regarding **the kind of personal data processed**

- **Sufficiently informative**

The required information are comprehensively specified and not vague

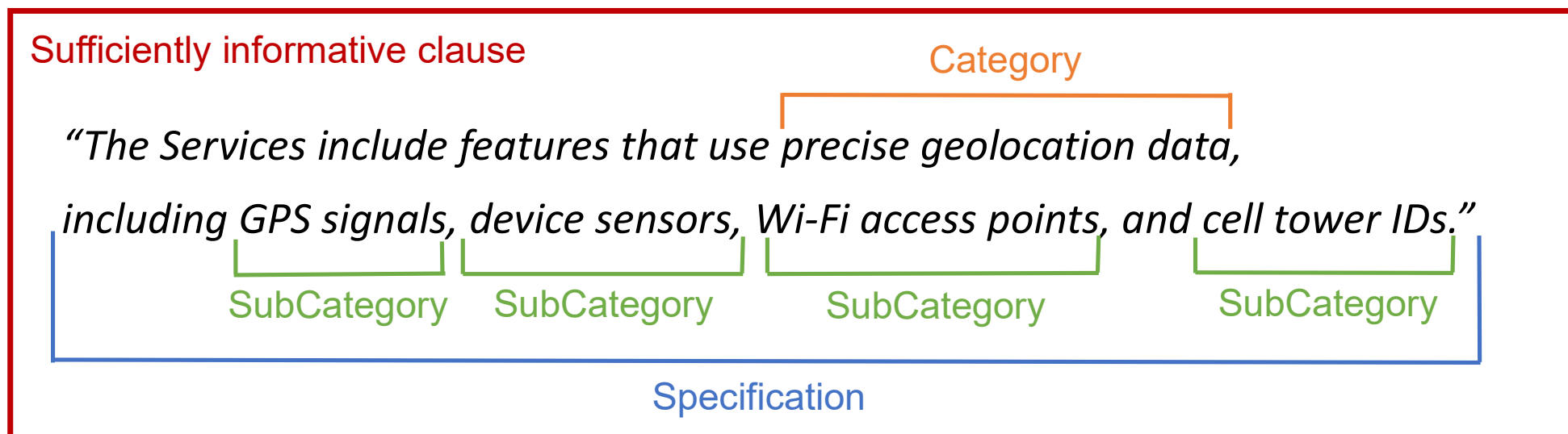
Example: "*We collect your payment information such as, name, surname, credit card number and expiration date*" [PayPal]

- **Insufficiently informative**

The required information are missing or they are present but vague and imprecise

Example: "*We collect **information about your interactions with the Airbnb Platform**, like the bookings you have made*" [AirBnB]

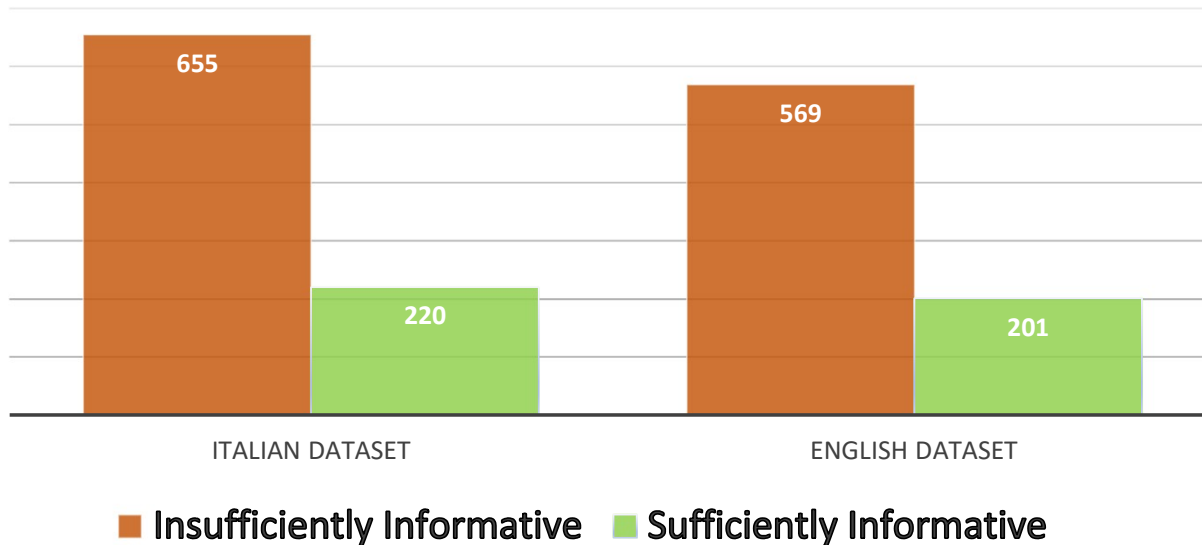
# Hierarchical annotation method



The assessment depends on a set of sub-elements and attributes:

- **Closed Terms:** concrete and exhaustive (e.g. Payment information)
- **Open Terms:** vague and ambiguous (e.g. Usage information)

# Dataset

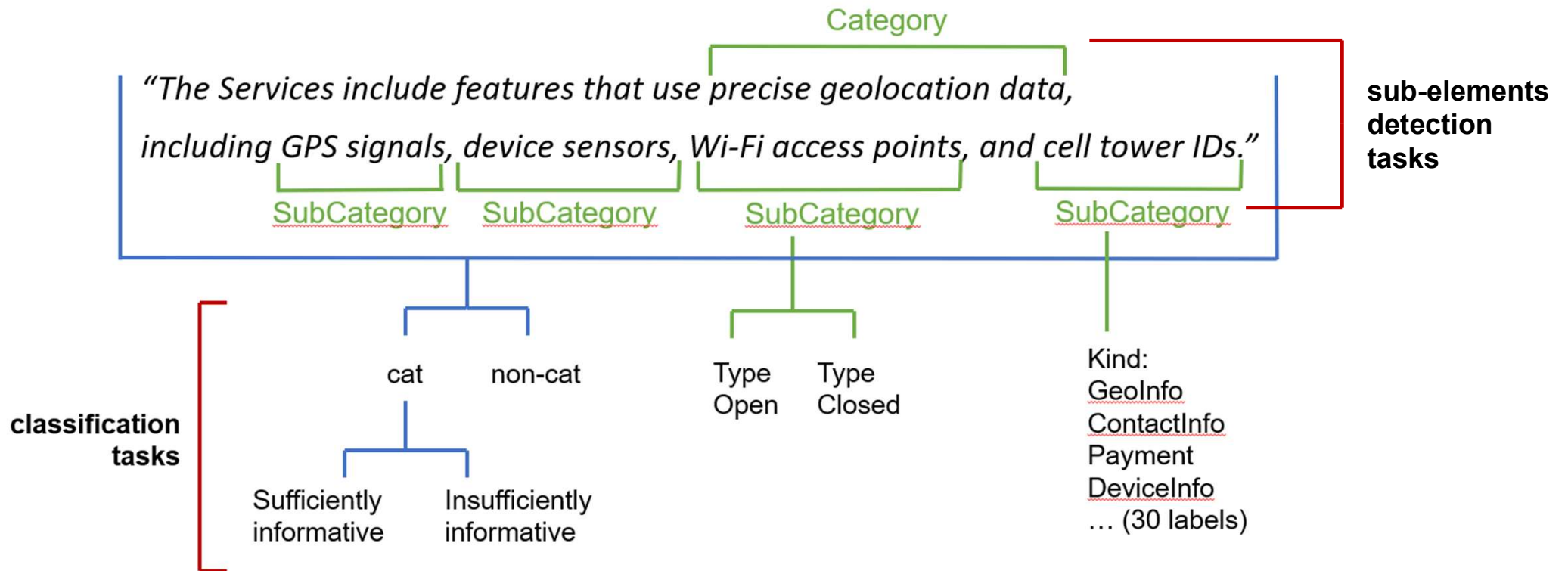


**NO PP WAS FULLY COMPLIANT**

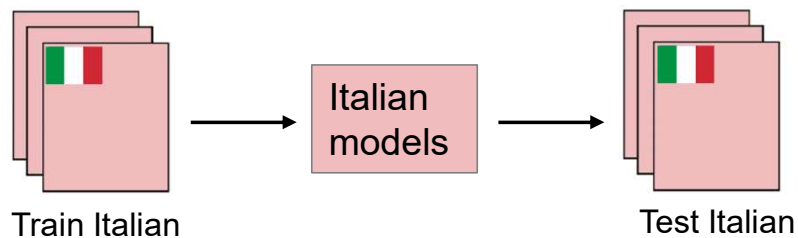
Element	Ita	Eng
documents	30	30
sentences	5862	6156
cat	875	770
<i>Sufficiently inf.</i>	220	201
<i>Insufficiently inf.</i>	655	569

- Italian PPs longer and more detailed
- **Selection criteria:** different market sectors, global relevance, users number
- **Manually annotated by legal expert**

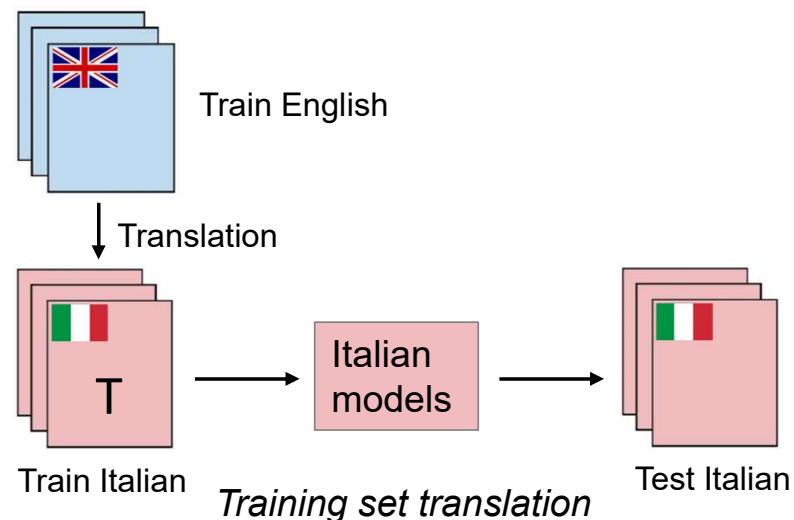
# Tasks



# Cross-lingual setting with BERT models



*Learning on Italian data*

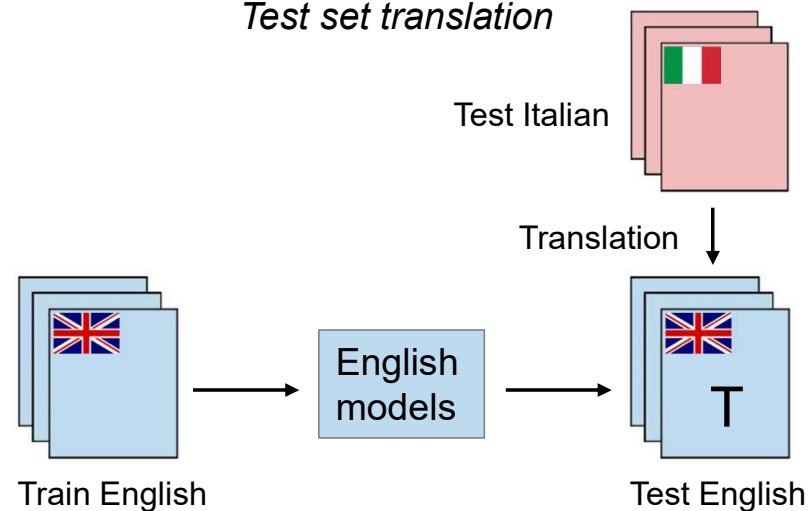


*Training set translation*

*Multilingual models*



*Test set translation*



# Cross-lingual results

(macro F1 score)

Method and Model	cat	Inform.	Type	Kind
<i>Learning on Italian data</i>				
Italian BERT	<b>0.86</b>	0.79	<b>0.80</b>	0.44
<i>Multilingual model</i>				
BERT multilingual	0.68	0.65	0.73	0.30
<i>Training set translation</i>				
Italian BERT	0.85	0.78	0.77	0.45
UmBERTo	0.85	<b>0.80</b>	0.73	0.32
<i>Test set translation</i>				
DeBERTa	0.85	<b>0.80</b>	0.79	0.43

- Comparable results in all settings, except for the multilingual model
- Top score for **INFORMATIVENESS** when training in English and translating the test set at inference time (most convenient setting in term of resources)

# LLM experimental setting

- **Zero-shot:** definitions from the annotation guidelines
- **Static few-shot:** 15 random examples
  - balancing the number of examples per class
- **Dynamic few-shot:** 15 most similar examples to the given query
  - k-NN algorithm
- **Models:** Gemini 2.0 Flash, Llama-3.1-8B-Instruct

# Classification results

(macro F1 score)

Method and Model	cat	Inform.	Type	Kind
<i>No learning</i>				
gemini zero-shot	0.85	0.64	0.67	0.56
llama zero-shot	0.77	0.64	0.69	0.44
<i>Learning on Italian data</i>				
gemini static few-shot	0.83	0.70	0.65	0.54
gemini dynamic knn few-shot	0.86	0.72	0.76	0.54
llama static few-shot	0.73	0.67	0.63	0.45
llama dynamic knn few-shot	0.77	0.71	0.73	0.51
Italian BERT	0.86	0.79	0.80	0.44
DeBERTa	0.85	0.80	0.79	0.43

- BERT > generative LLMs
  - Except for Kind

- Zero-shot > static few-shot
- Dynamic few-shot best among LLMs



- Example selection is crucial

# Detection results

## (F1 score)

Model	Category-Subcategory					Specification				
	<i>B-C</i>	<i>I-C</i>	<i>B-S</i>	<i>I-S</i>	<i>O</i>	Avg.	<i>B-Sp</i>	<i>I-Sp</i>	<i>O</i>	Avg.
<i>BIO tagging</i>										
gemini zero-shot	0.48	0.51	0.71	0.72	0.85	0.65	0.64	0.51	0.80	0.65
gemini few-shot	0.66	0.60	0.81	0.77	0.88	0.75	0.51	0.69	0.82	0.67
llama few-shot	0.56	0.51	0.67	0.62	0.81	0.63	0.29	0.33	0.62	0.41
Italian BERT	0.74	0.68	0.81	0.80	0.90	<b>0.79</b>	0.74	0.91	0.93	<b>0.86</b>
UmBERTo	0.67	0.64	0.84	0.80	0.91	<b>0.77</b>	0.69	0.94	0.96	<b>0.86</b>
<i>IO tagging</i>										
gemini zero-shot		0.56		0.76	0.85	0.72		0.52	0.80	0.66
gemini few-shot		0.63		0.81	0.88	0.77		0.70	0.82	0.76
llama few-shot		0.54		0.65	0.81	0.67		0.34	0.62	0.48
Italian BERT		0.70		0.83	0.90	<b>0.81</b>		0.89	0.92	<b>0.91</b>
UmBERTo		0.68		0.84	0.90	<b>0.81</b>		0.92	0.94	<b>0.93</b>

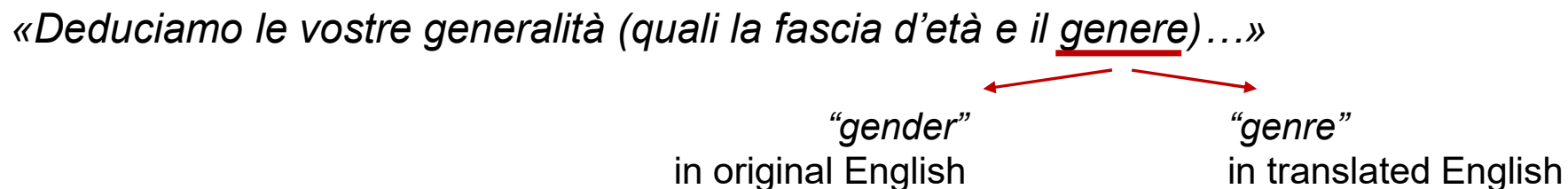
- BERT > generative LLMs
- Few-shot > zero-shot

# Error Analysis

- We analyze the results of the **test set translation** approach
- Lexicon issues – differences in the domain-specific terminology in the two languages



- Lack of context – when classifying the sub-elements



# Conclusions and Future works

- BERT-based models yield better results than LLM
- English resources can be used to train a system for Italian with comparable results
  - Potentially generalizable to other languages
- New directions for LLM example selection:
  - Expert selected instead of random
  - Change number of examples
  - Dynamic selection with similarity+variety

Working on AI-empowered tool to detect potentially unlawful clauses in PPs

**PRIMA**

An automated detector of potentially unfair clauses

Copy your text here

Submit

# THANK YOU!



Detecting Vague Clauses in Italian Privacy Policies using  
Transformers, LLMs, and Cross-Lingual Techniques

**28th European Conference on Artificial Intelligence**

# Hierarchical Annotation System

## Element and Sub-elements:

- **<Category>** is the genus of personal data processed [e.g. “your location”]

This element may include two optional sub-elements:

- **<Specification>** is the term that introduces a list of sub-categories [e.g., “such as”]
- **<SubCategory>** is the data denoting a species of the given genus [e.g., “GPS information”]

## Mandatory and Optional Attributes:

**<Type>**: The precision of terms describing the data;

**<Kind>**: the terms describing the category of data and so identifying its content;

**<Ref>**: The link between the subelement and Category it refers to.

# Hierarchical Annotation System

---

<b>Sufficiently Informative (Level="1")</b>		<b>Insufficiently Informative (Level="2")</b>	
Rule 1	Category Type = "Closed" No Specification No Subcategories	Rule 1	Category Type = "Open" No Specification No Subcategories
Rule 2	Category Type = "Closed" Specification Type = "Closed"/"Open" SubCategory Type = "Closed"/"Open"	Rule 2	Category Type = "Open" Specification Type = "Open" Subcategory Type = "Closed"/"Open"
Rule 3	Category Type = "Open" Specification Type = "Closed" SubCategory Type = "Closed"	Rule 3	Category Type = "Open" Specification Type = "Closed" Subcategory Type = "Open"

---

Method and Model	cat				Informativeness				Type				Kind	
	cat	non-cat	Avg.	$\sigma$	suff.	insuff.	Avg.	$\sigma$	open	closed	Avg.	$\sigma$	Avg.	$\sigma$
<i>No learning</i>														
gemini zero-shot	0.75	0.95	0.85	-	0.56	0.72	0.64	-	0.63	0.70	0.67	-	<b>0.56</b>	-
llama zero-shot	0.61	0.92	0.77	-	0.51	0.77	0.64	-	0.72	0.66	0.69	-	0.44	-
<i>Learning on Italian data</i>														
gemini static few-shot	0.72	0.94	0.83	-	0.61	0.79	0.70	-	0.59	0.70	0.65	-	0.54	-
gemini dynamic knn few-shot	<b>0.76</b>	0.95	<b>0.86</b>	-	0.64	0.79	0.72	-	0.76	0.77	0.76	-	0.54	-
llama static few-shot	0.57	0.89	0.73	-	0.56	0.78	0.67	-	0.57	0.69	0.63	-	0.45	-
llama dynamic knn few-shot	0.64	0.91	0.77	-	0.62	0.79	0.71	-	0.71	0.74	0.73	-	0.51	-
ITALIAN-LEGAL-BERT	0.72	0.95	0.83	0.008	0.62	0.86	0.74	0.027	0.79	0.76	0.78	0.005	0.37	0.020
Italian BERT	<b>0.76</b>	<b>0.96</b>	<b>0.86</b>	0.011	0.69	0.88	0.79	0.016	<b>0.82</b>	<b>0.78</b>	<b>0.80</b>	0.014	0.44	0.019
UmBERTo	<b>0.76</b>	<b>0.96</b>	<b>0.86</b>	0.008	0.67	<b>0.89</b>	0.78	0.043	0.79	0.74	0.76	0.013	0.34	0.036
BERT multilingual	0.75	<b>0.96</b>	0.85	0.010	0.68	0.86	0.77	0.009	0.81	0.77	0.79	0.016	0.41	0.021
<i>Multilingual model</i>														
BERT multilingual	0.41	0.94	0.68	0.015	0.45	0.85	0.65	0.109	0.75	0.72	0.73	0.011	0.30	0.024
<i>Training set translation</i>														
ITALIAN-LEGAL-BERT	0.71	0.95	0.83	0.004	0.68	0.87	0.78	0.030	0.76	0.75	0.76	0.008	0.38	0.007
Italian BERT	0.73	<b>0.96</b>	0.85	0.005	0.68	0.88	0.78	0.056	0.76	<b>0.78</b>	0.77	0.020	0.45	0.026
UmBERTo	0.73	<b>0.96</b>	0.85	0.003	<b>0.72</b>	0.88	<b>0.80</b>	0.022	0.73	0.73	0.73	0.024	0.32	0.049
BERT multilingual	0.75	<b>0.96</b>	0.85	0.017	0.67	0.88	0.78	0.026	0.78	0.76	0.77	0.019	0.43	0.001
<i>Test set translation</i>														
LEGAL-BERT	0.74	<b>0.96</b>	0.85	0.016	0.67	0.87	0.77	0.015	0.80	<b>0.78</b>	0.79	0.006	0.45	0.015
DistilRoBERTa	0.73	<b>0.96</b>	0.85	0.008	0.71	0.88	<b>0.80</b>	0.013	0.79	0.75	0.77	0.012	0.45	0.003
DeBERTa	0.74	<b>0.96</b>	0.85	0.017	<b>0.72</b>	<b>0.89</b>	<b>0.80</b>	0.005	0.80	<b>0.78</b>	0.79	0.003	0.43	0.009
<hr/>														
Method and Model	cat				Informativeness				Type				Kind	
	cat	non-cat	Avg.	$\sigma$	suff.	insuff.	Avg.	$\sigma$	open	closed	Avg.	$\sigma$	Avg.	$\sigma$
<i>Learning on English data</i>														
LEGAL-BERT	0.75	<b>0.96</b>	<b>0.86</b>	0.004	<b>0.77</b>	<b>0.92</b>	<b>0.84</b>	0.029	<b>0.81</b>	<b>0.79</b>	<b>0.80</b>	0.009	<b>0.46</b>	0.017
<i>Training set translation</i>														
LEGAL-BERT	<b>0.76</b>	<b>0.96</b>	<b>0.86</b>	0.014	0.69	0.91	0.80	0.025	0.80	0.76	0.78	0.016	0.36	0.017
<i>Test set translation</i>														
Italian BERT	0.74	<b>0.96</b>	0.85	0.012	0.76	<b>0.92</b>	<b>0.84</b>	0.022	0.80	0.75	0.77	0.031	0.35	0.013