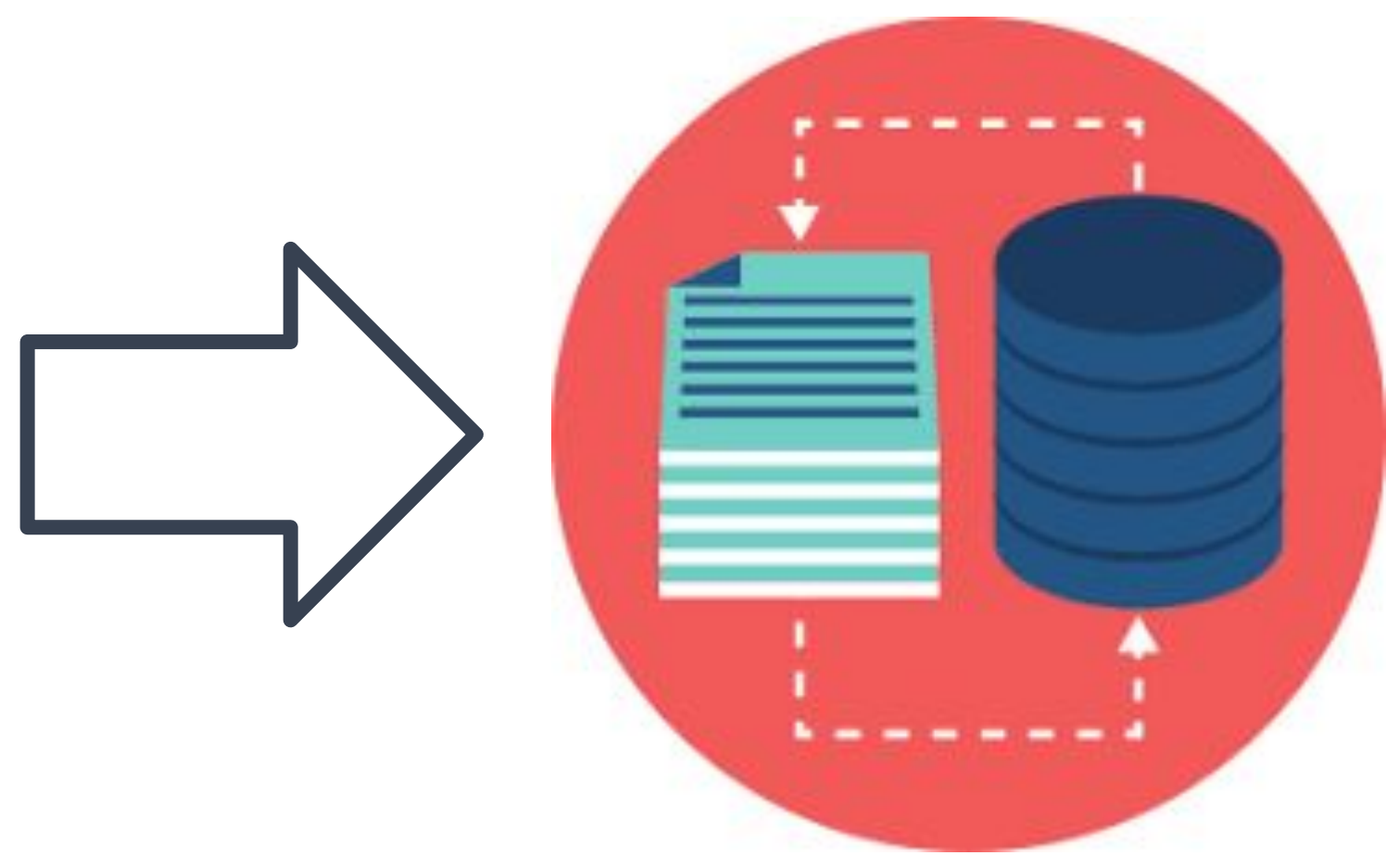


# Using Arguments to Assess Models' Critical Thinking Skills

Xabier de Zuazo, Carlos Domínguez, and Jon Manzanal  
NLP Applications II / University of the Basque Country, Donostia, Basque Country  
xzuazo002@ikasle.ehu.eus, cdominguez019@ikasle.ehu.eus, jmanzanal001@ikasle.ehu.eus

## Objective

We propose different methods to classify, detect, and generate arguments using transformers. The project's main proposal is to evaluate the critical thinking skills of varied natural language processing models through a classification task on fallacy prediction and compare them. As a secondary task, we also analyze the models' capacity to generate counter-arguments by themselves, and detect fallacies within texts.



## Data

We have used two different datasets, which come from two papers:

- **Riposte! corpus (Reisert et al., 2019)**: Over 18,000 counter-arguments. It contains claims, premises, and counter-arguments for 4 different types of fallacies. Used in Tasks 1 and 3.
- **Argument Mining CL2017 corpus (Habernal and Gurevych, 2017)**: 340 documents annotated with sequences of arguments. Used in task 2.

## Tasks

### Task 1: Counter-Argument Multi-Label Natural Language Inference

It categorizes the fallacy types found on a given claim and the premise.

The Riposte dataset annotations have been converted to a multi-label corpus.

### Task 2: Sequence Tagging Arguments

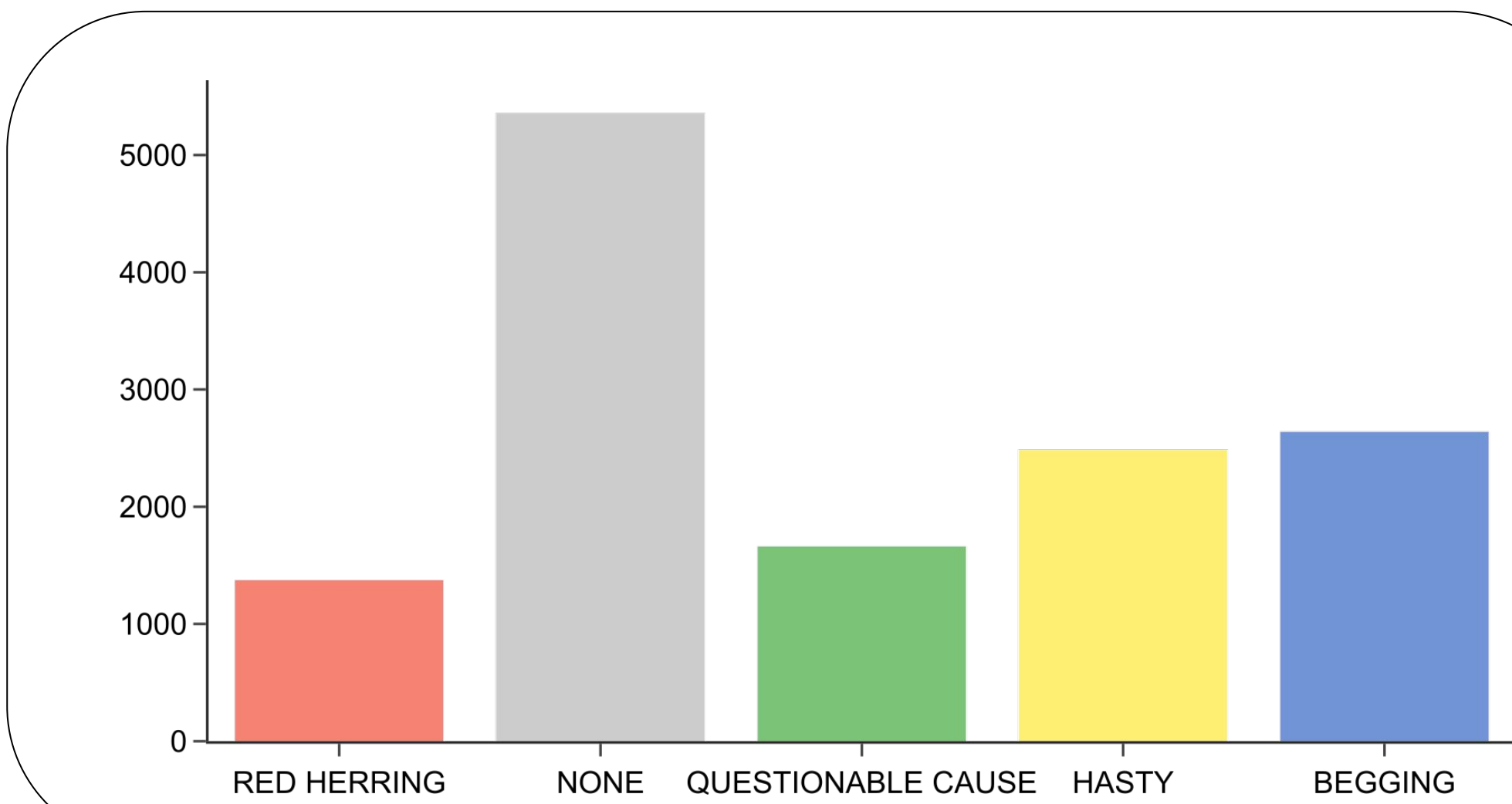
It detects arguments within the text by assigning to each token a tag. The available tags are O (other), B-ARG (beginning argument), and I-ARG (inside argument).

### Task 3: Counter-Argument Natural Language Generation

It gives one possible counter-argument using as information the claim, the fallacious premise, and the fallacy type.

## Fallacy Types

- **BEGGING** : The argument assumes the truth of the conclusion instead of supporting it.
- **HASTY** : Over-generalization, making a claim based on evidence that is just too small.
- **QUESTIONABLE CUASE** : The cause is incorrectly identified.
- **RED HERRING** : Something that misleads or distracts from a relevant or important question.



## Results

System	Accuracy ↑	Precision ↑	Recall ↑	F1-score ↑	BLEU ↑
Glove <sup>1</sup>	0.111	<b>0.674</b>	0.294	0.409	-
ELMo <sup>1</sup>	0.104	0.667	0.290	0.405	-
Flair <sup>1</sup>	0.111	<b>0.674</b>	0.294	0.409	-
bert-base-cased <sup>1</sup>	0.629	0.642	0.629	0.633	-
bert-base-uncased <sup>1</sup>	0.638	0.641	0.638	0.639	-
roberta-base <sup>1</sup>	0.605	0.606	0.605	0.604	-
distilroberta-base <sup>1</sup>	<b>0.664</b>	0.665	<b>0.664</b>	<b>0.665</b>	-
xlm-roberta-base <sup>1</sup>	0.605	0.606	0.605	0.604	-
Glove <sup>2</sup>	0.056	0.153	0.082	0.107	-
ELMo <sup>2</sup>	0.600	0.719	0.783	0.750	-
Flair <sup>2</sup>	0.517	0.664	0.701	0.682	-
bert-base-cased <sup>2</sup>	0.914	0.634	0.761	0.692	-
bert-base-uncased <sup>2</sup>	0.924	0.709	0.786	0.745	-
roberta-base <sup>2</sup>	<b>0.938</b>	0.774	0.800	0.787	-
distilroberta-base <sup>2</sup>	0.926	0.701	0.786	0.741	-
xlm-roberta-base <sup>2</sup>	<b>0.938</b>	<b>0.778</b>	<b>0.837</b>	<b>0.806</b>	-
TS <sup>3</sup>	-	-	-	-	<b>0.654</b>

## Example Task 2 (ST)

International tourism is now more common than ever before. The last decade has seen an increasing number of tourists traveling to visit natural wonder sights, ancient heritages and different cultures around the world. Firstly, international tourism promotes many aspects of the destination country's economy in order to serve various demands of tourists. These demands trigger related business in the surrounding settings, which in turn creates many jobs for local people and improve infrastructure and living standards. Therefore tourism has clearly improved lives in the tourist country. Without this support and profit from tourism, many traditional cultures would disappear due to their low-income work. International tourism has both triggered economic development and maintained the cultural and environmental values of the tourist countries. In addition, the authorities should adequately support these sustainable developments.

## Example Task 1 (NLI)

**Claim:** Is Google a harmful monopoly?

**Premise:** Monopolies are not good for the market.

**Types of fallacies found:** Hasty, Questionable Cause, and Begging.

## Example Task 3 (GEN)

**Source:** Women should not delay motherhood. Many young women don't realize the difficulties of IVF.

**BEGGING**

**Counter-argument:** If many young women don't realize the difficulties of IVF are assumed to be true, then women should not delay motherhood is already assumed to be true.

## Cascade Models

