

# Multi-Lingual Natural Language Processing

## Homeworks Report



**SAPIENZA**  
UNIVERSITÀ DI ROMA

**Dario Loi — 1940849**

M.Sc. in AI & Robotics,  
Sapienza, University of Rome.

**A. Y. 2023–2024**

# Tasks

For the first homework, I was assigned **two** tasks:

1. WiC-ITA: Detect whether two italian words in two **different** sentences are used with the **same** meaning.
2. ITAmoji: Predict which emoji was used in a given italian tweet.

We will spend a few slides on each task.

## WiC-ITA

The first task involved the parsing of JSONL files containing the Word-in-Context-ITA dataset, an example of which is shown below:

[Listing](#): Sample 23 from WiC-ITA train.jsonl.

```
1 {
2   "id": "lira.noun.15",
3   "lemma": "lira",
4   "sentence1": "In caso di inosservanza degli obblighi stabiliti dal comma 1 , si applica la
5     sanzione amministrativa pecuniaria da lire dieci milioni a lire cento milioni .",
6   "sentence2": "Per le finalità di cui all' Art. 1 , comma 2 , della legge regionale 26 aprile 1
7     995 , n. 31 recante \" Norme in materia di musei degli Enti locali e di interesse locale
8     \" è autorizzata per l' esercizio finanziario 2000 la spesa di lire 200.000.000 .",
9   "start1": 115,
10  "end1": 119,
11  "start2": 230,
12  "end2": 234,
13  "label": 1
14 }
```

## Desired Format

The samples were rewritten for use by an LLM, as shown below:

**Listing:** Same sample with required changes.

```
1  {
2    "id": "lira.noun.15",
3    "lemma": "lira",
4    "sentence1": "In caso di inosservanza degli obblighi stabiliti dal comma 1 , si applica la
                    sanzione amministrativa pecuniaria da lire dieci milioni a lire cento milioni .",
5    "sentence2": "Per le finalità di cui all' Art. 1 , comma 2 , della legge regionale 26 aprile
                    1995 , n. 31 recante \" Norme in materia di musei degli Enti locali e di interesse
                    locale \" è autorizzata per l' esercizio finanziario 2000 la spesa di lire 200.000.000
                    .",
6    "start1": 115,
7    "end1": 119,
8    "start2": 230,
9    "end2": 234,
10   "choices": [
11     "DIVERSO",
12     "UGUALE"
13   ],
14   "label": 1
15 }
```

## ITAmoji

The second task also involved manipulation of JSONL files, a typical sample of ITAmoji is shown below:

**Listing:** Sample 27 from ITAmoji\_i\_2018\_TRAINdataset\_v1.ANON.list.

```
1 {  
2   "uid": "447352763",  
3   "text_no_emoji": "... il rumore del mare \ufe0f #28Settembre <URL>",  
4   "created_at": "Thu Sep 28 15:32:06 +0000 2017",  
5   "label": "red_heart",  
6   "tid": "913426094002458626"  
7 }
```

For this task, we also have to add **distractors**, that is, plausible alternatives to the correct label.

## Desired Format... Again

The output from our distractor generation process on the previous sample is as follows:

**Listing:** Sample 27 with generated distractors.

```
1  {
2    "id": "ITA-emoji-train-00000027",
3    "sentence": "... il rumore del mare \ufe0f #28Settembre <URL>",
4    "choices": [
5      "red_heart",
6      "two_hearts",
7      "rose",
8      "kiss_mark"
9    ],
10   "label": 0
11 }
```

## Distractor Clusters

To generate distractors, we used a hand-crafted list of **semantic clusters**, which we define as a set of emoji that are semantically related. For example, the cluster love contains the emojis red\_heart, two\_hearts, blue\_heart, rose, and kiss\_mark. When augmenting a sample, we select a cluster in which the correct label is present, and then randomly select three other labels from the same cluster. This way, we ensure that the distractors are capable to confuse the model.

## Distractor Sampling

To obtain a set of distractors from an ITAmoji sample, we follow this process:

- Randomly sample a cluster from the list of clusters that contain the correct emoji.
- Randomly sample three other emojis from the same cluster.
- Return the new sample with the correct label and the three distractors.



# Distractor Sampling

To obtain a set of distractors from an ITAmoji sample, we follow this process:

- Randomly sample a cluster from the list of clusters that contain the correct emoji.
- Randomly sample three other emojis from the same cluster.
- Return the new sample with the correct label and the three distractors.

# Distractor Sampling

To obtain a set of distractors from an ITAmoji sample, we follow this process:

- Randomly sample a cluster from the list of clusters that contain the correct emoji.
- Randomly sample three other emojis from the same cluster.
- Return the new sample with the correct label and the three distractors.

# Distractor Sampling

To obtain a set of distractors from an ITAmoji sample, we follow this process:

- Randomly sample a cluster from the list of clusters that contain the correct emoji.
- Randomly sample three other emojis from the same cluster.
- Return the new sample with the correct label and the three distractors.