# Combining Textual and Visual Features to Identify Anomalous User-generated Content

**Lucia Noce, Ignazio Gallo, Alessandro Zamberletti**
Università dell'Insubria, Varese, Italy
*http://artelab.dicom.uninsubria.it/*

## 1. Problem

Nowadays many websites allow users to post their own content.

- It exposes the platform to a wide number of possible threats e.g. fake, low quality or malicious user-generated content may damage the credibility of the website.
- Our GOAL: automatically identifying fake or low quality user-generated content.

We apply anomaly detection to a novel task: discovering fake or suspicious descriptions of commercial products.

## 2. Solution

- To infer the key features describing the behavioral traits of expert users.
- To combine textual descriptors with visual information extracted from the media resources associated with each product description.
- The joint use of textual and visual features helps in obtaining a robust detection model, that automatically reports whenever a newly generated description is anomaluos.

## 3. Method

### Textual Features:

**Description Length.**
Length of the description in terms of number of text characters.
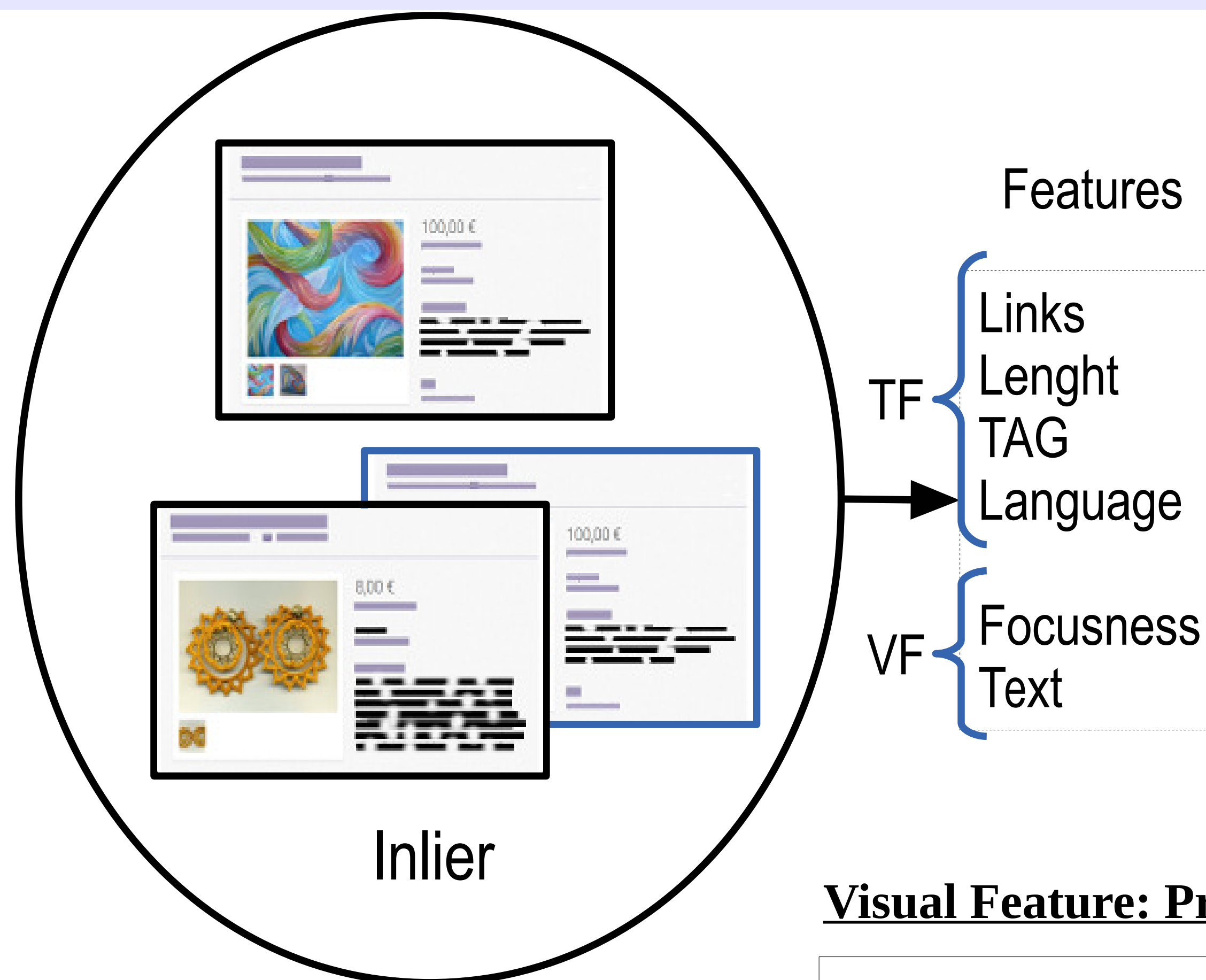**Description Language.**
Descriptions should always be written using the platform main idiom.
**Description Keywords.**
Index of consistency of the textual information provided by users.
**Presence of Hyperlinks.**
Presence of email addresses and website URLs.



Inlier

Features

TF { Links, Lenght, TAG, Language }

VF { Focusness, Text }

One-Class SVM → Anomaly / OK

### Proposed Method:

The one-class Support Vector Machine model is trained using both textual and visual features (TF and VF respectively) extracted from high quality genuine commercial product descriptions (inliers).

### Visual Feature: Focusness

Variance of Laplacian, is the operator that determines whether an image is focused or unfocused.

**Focused.**     **Unfocused.**

### Visual Feature: Presence of Visual Hyperlinks.

- Tesseract detects whether an image contains text.
- Regex based rules determine whether it represents a hyperlink/email address or not.
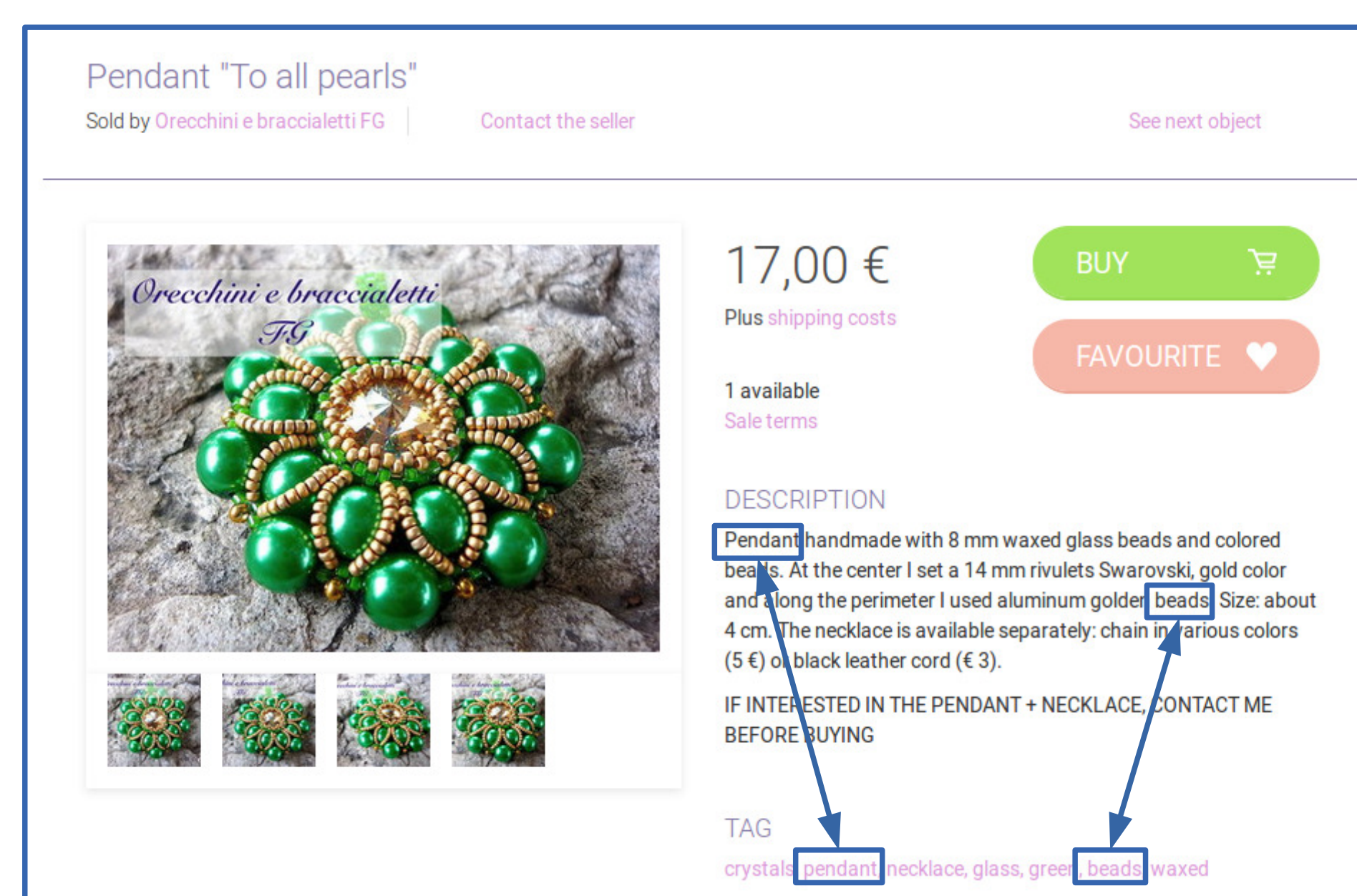
**Allowed Text.**     **Not Allowed Text.**
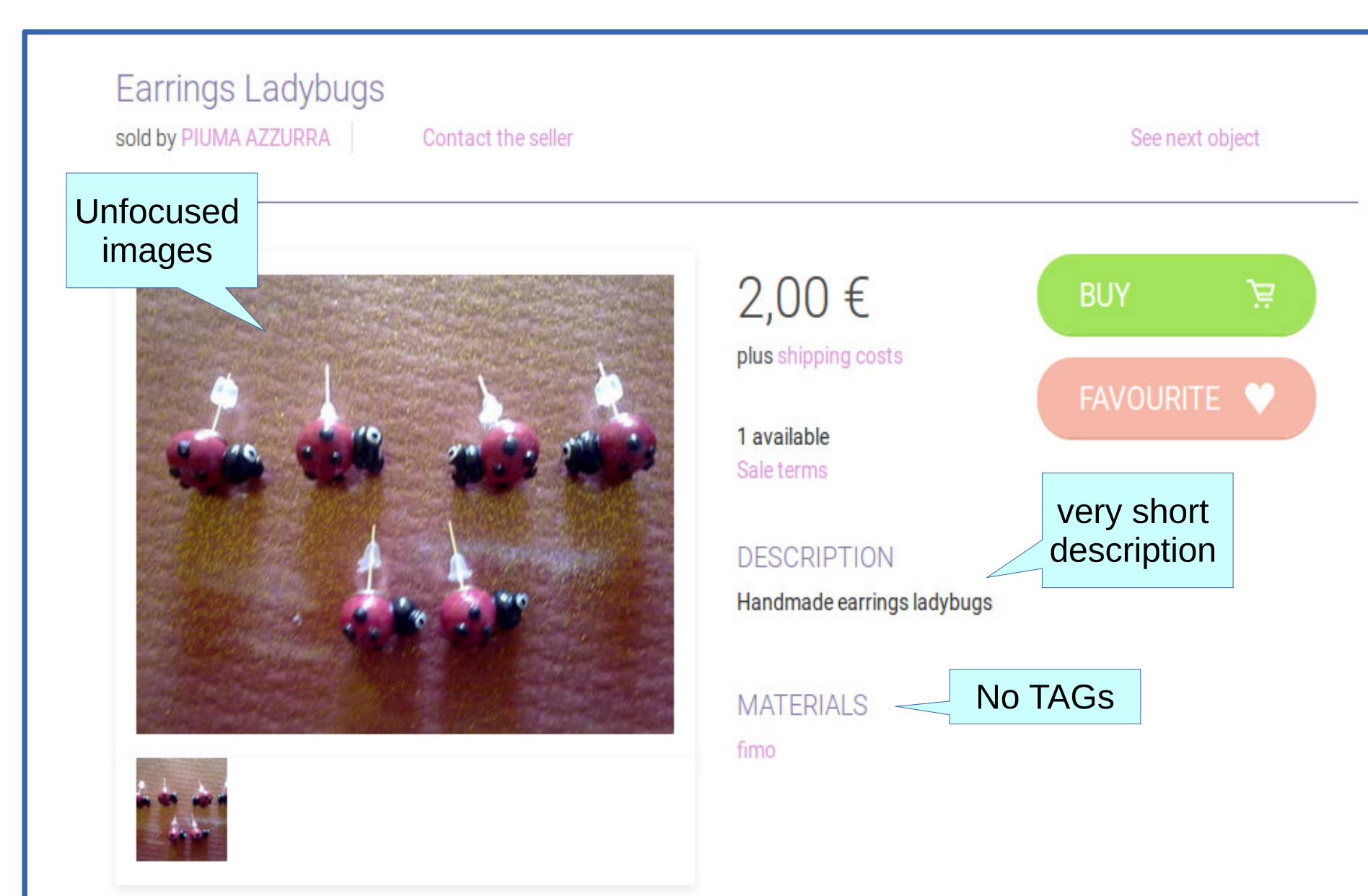
## 4. Results

### Train Dataset

We gather a dataset composed of 41635 documents from a website specialized on selling handmade products. Each insertion is user-generated and composed of a textual description, a set of keywords or tags and one or more images representing the item. In our experiments, 80% of the high quality genuine documents are used for training, while the remaining 20% are used for testing.

### Anomalous Dataset

The set of possible anomalous descriptions has been articially built by randomly injecting dierent kinds of anomalies into high quality genuine user-generated documents, trying to cover a large set of possible anomalous behaviors. We articially create a total of 1000 negative documents.

### Overall accuracies

In our experiments, we divide the 8327 positive test documents into 8 equal splits, and we build 8 different test sets by adding each split to the set of anomalous documents. Each test set is composed of 1040 positive and 1000 anomalous documents. The goodness of the proposed system at detecting anomalous commercial product descriptions was measured.

|  | Anomalies | Correct |  |  |
|---|---|---|---|---|
| Identified as anomalous | 858 | 94 | Accuracy | 0.88 |
| Identified as correct | 142 | 946 | Precision | 0.90 |
|  |  |  | Recall | 0.85 |
|  |  |  | F-measure | 0.88 |

### Detection rate

To evaluate the capability of the proposed method at detecting different types of anomalies. The detection rate has been calculated for each type of anomalous content considered.

| Anomalies | Detection Rate |
|---|---|
| Anomalous length | 39% |
| Anomalous language | 79% |
| Missing/wrong *tags* | 43% |
| Presence of hyperlinks | 99% |
| Unfocused images | 97% |
| Images containing hyperlinks | 98% |
| ≥ 2 random anomalies | 97% |
| Average | 78% |