

# Detailed Comparison Module in CoReMo 1.9 Plagiarism Detector

Diego A. Rodríguez-Torrejón<sup>1,2</sup>

José Manuel Martín-Ramos<sup>1</sup>



**Universidad  
de Huelva**

<sup>1</sup> Universidad de Huelva  
jmmartin@dti.uhu.es



<sup>2</sup> I.E.S. José Caballero  
dartsystems@gmail.com

**ELÍAS**  
evaluating information access systems

**EUROPEAN  
SCIENCE  
FOUNDATION**  
SETTING SCIENCE AGENDAS FOR EUROPE

Our attendance is sponsored and granted by ELIAS (ESF Research Networking Programme).



Universidad  
de Huelva

## Detailed Comparison Module in CoReMo 1.9 Plagiarism Detector



CoReMo is a classic PAN competitor since PAN2010.

Its **detection speed** was a highlighted goal ever .

Detailed Comparison module is new however, but also improved by **Surrounding Context N-grams (SCnG)**, an extended concept of former CTnG (case folding, stopwords/short removal, stemming and internal sort) by **including new special skip n-grams** to the classic consecutive. It gets almost 3 times more n-grams than words in the text, having a discriminative magnifier effect: “*The quick brown fox jumps over the lazy dog*”

*brown\_fox\_quick , brown\_jump\_quick , fox\_jump\_quick*



Universidad  
de Huelva

## Detailed Comparison Module in CoReMo 1.9 Plagiarism Detector



The highlighted runtime speed is due to:

- **C/C++ 64 bits** programming (single core however)
- **GNU Linux 64bits OS** and **ext4 file system** platform
- Internal sort of n-grams by ***Bubblesort*** algorithm
- n-grams into a document ordered by ***Quicksort***
- **Modified *Mergesort*** algorithm to compare both docs
- **Local translations** by dictionary mapping
- **Taking the advantage** of suspicious document modelling when repeated in **consecutive comparisons**





Universidad  
de Huelva

## Detailed Comparison Module in CoReMo 1.9 Plagiarism Detector



Detection is reached when minimum length and distance between n-grams conditions are got in both suspicious (counting n-grams) and source (counting chars) sections.

$$\text{maxNgramDist} = 2 \cdot \text{chunkLength}$$

$$\text{maxCharDist} = \text{chunkLength} \cdot \text{wordLengthAverage}$$

$$\text{minNgramLength} = (\text{monotony} - 1.5) \cdot \text{chunkLength}$$

$$\text{minCharLenght} = \text{minNgramLength} \cdot \text{wordLengthAverage}$$

$$\text{chunkLength} = 4 \text{ trigrams} \rightarrow 6 \text{ words (monolingual)}$$

$$\text{Monotony} = 2 \text{ chunks}$$

*Joining distance for direct detections: 80.000 chars*