# UNIVERSITY OF LEEDS

This is a repository copy of *Computer-aided error annotation: a new tool for annotating Arabic error*.

White Rose Research Online URL for this paper:
https://eprints.whiterose.ac.uk/83911/

**Conference or Workshop Item:**

Alfaifi, AYG and Atwell, ES (2015) Computer-aided error annotation: a new tool for annotating Arabic error. In: 8th Saudi Students Conference, 31 Jan - 01 Feb 2015, Queen Elizabeth II Conference Centre, London.

eprints@whiterose.ac.uk
https://eprints.whiterose.ac.uk/

# Computer-Aided Error Annotation
## A New Tool for Annotating Arabic Error

Abdullah Alfaifi and Eric Atwell
University of Leeds {scayga, e.s.atwell}@leeds.ac.uk

## WHY THIS TOOL?

Existing tools for annotating errors in learner corpora are developed for languages other than Arabic.
Thus, this poster introduces a new tool for computer-aided error annotation in Arabic learner corpora.

The new tool includes the Arabic Error Taxonomy (Alfaifi et al., 2013) which was designed specifically for Arabic

**Arabic Error Taxonomy**
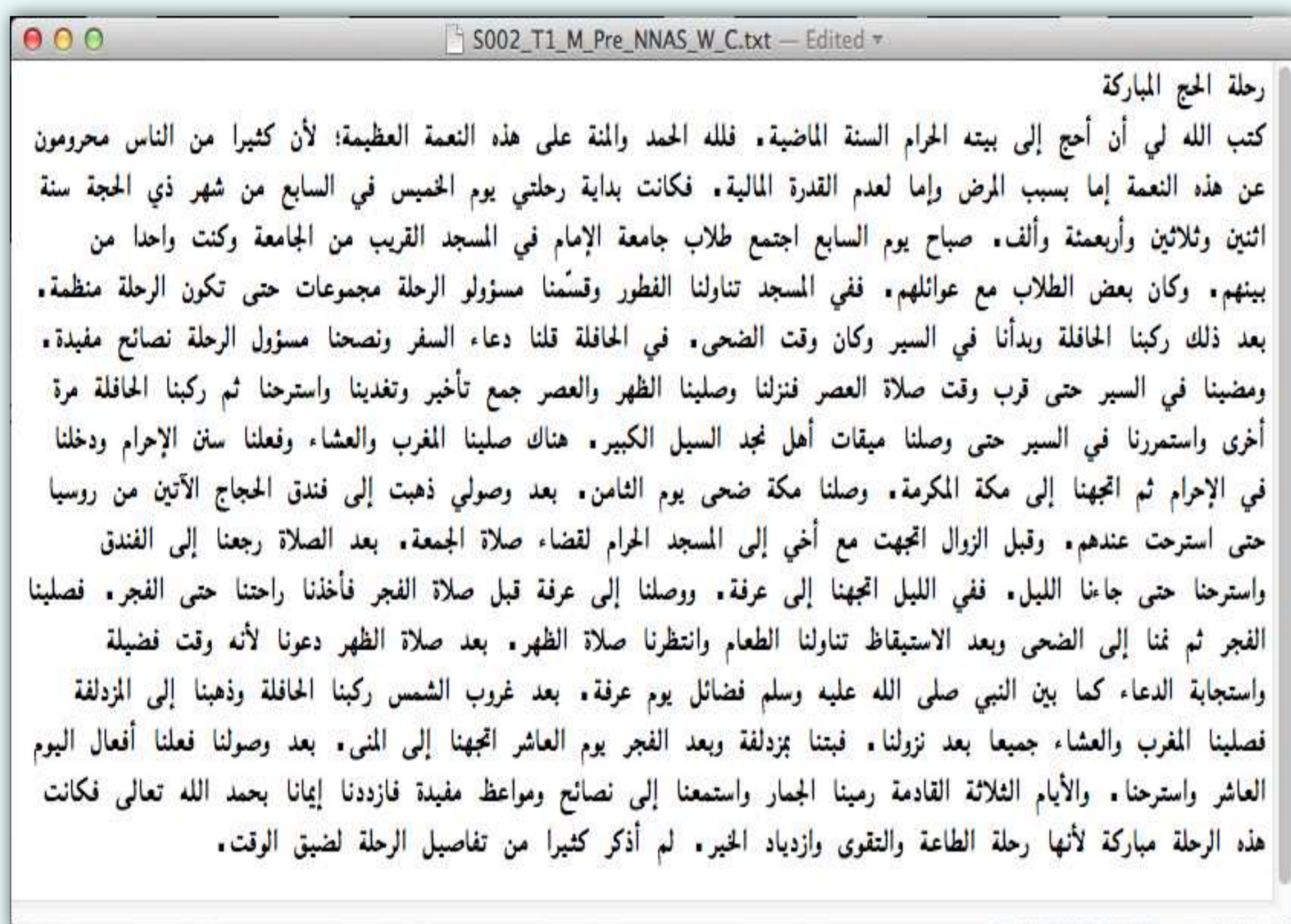


## ARABIC LEARNER CORPUS

The error annotation will be applied to the Arabic Learner Corpus (ALC), it adds more value to the corpus data.

The corpus includes 282,732 words, produced by 942 learners of Arabic.

**ALC ARABIC LEARNER CORPUS**
المدونة اللغوية لمتعلمي اللغة العربية

## CURRENT DATA

The Arabic Learner Corpus currently includes 1585 texts in a raw format.



**Example of raw text**

## Annotation

## TARGET DATA
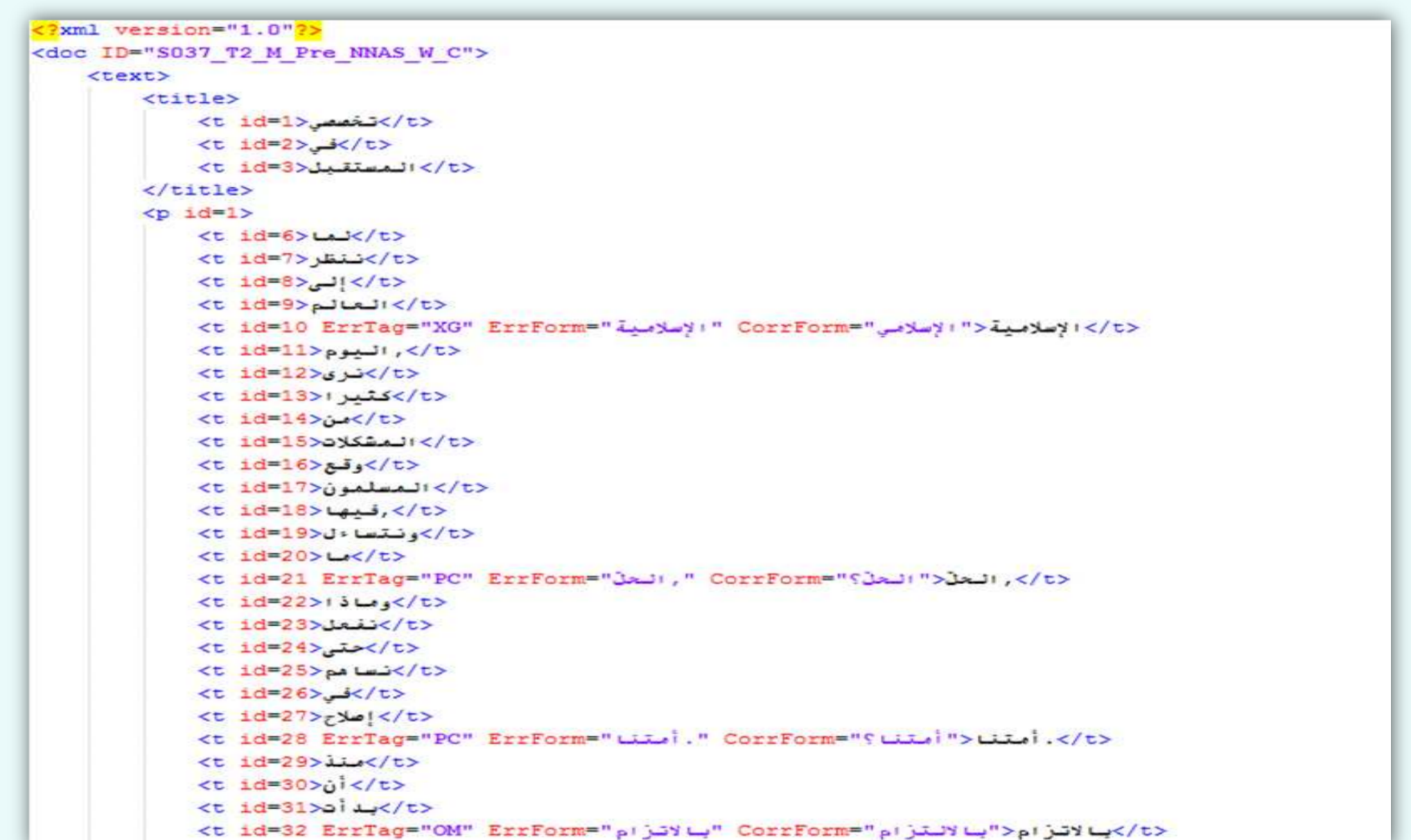
The ALC will include further information describing the language errors

| Error Type | Error Form | Correct Form |

<t id=10 ErrTag="XG" ErrForm="الإسلامية"  CorrForm="الإسلامي"الإسلامية</t>



**Example of annotated text**

## THE ANNOTATION TOOL
## Main functions



### 1. Text Tokenisation

The *Text Tokenisation* function helps in segmenting the text into separate word-based tokens, this enables the annotator to attach tags to those tokens include errors.
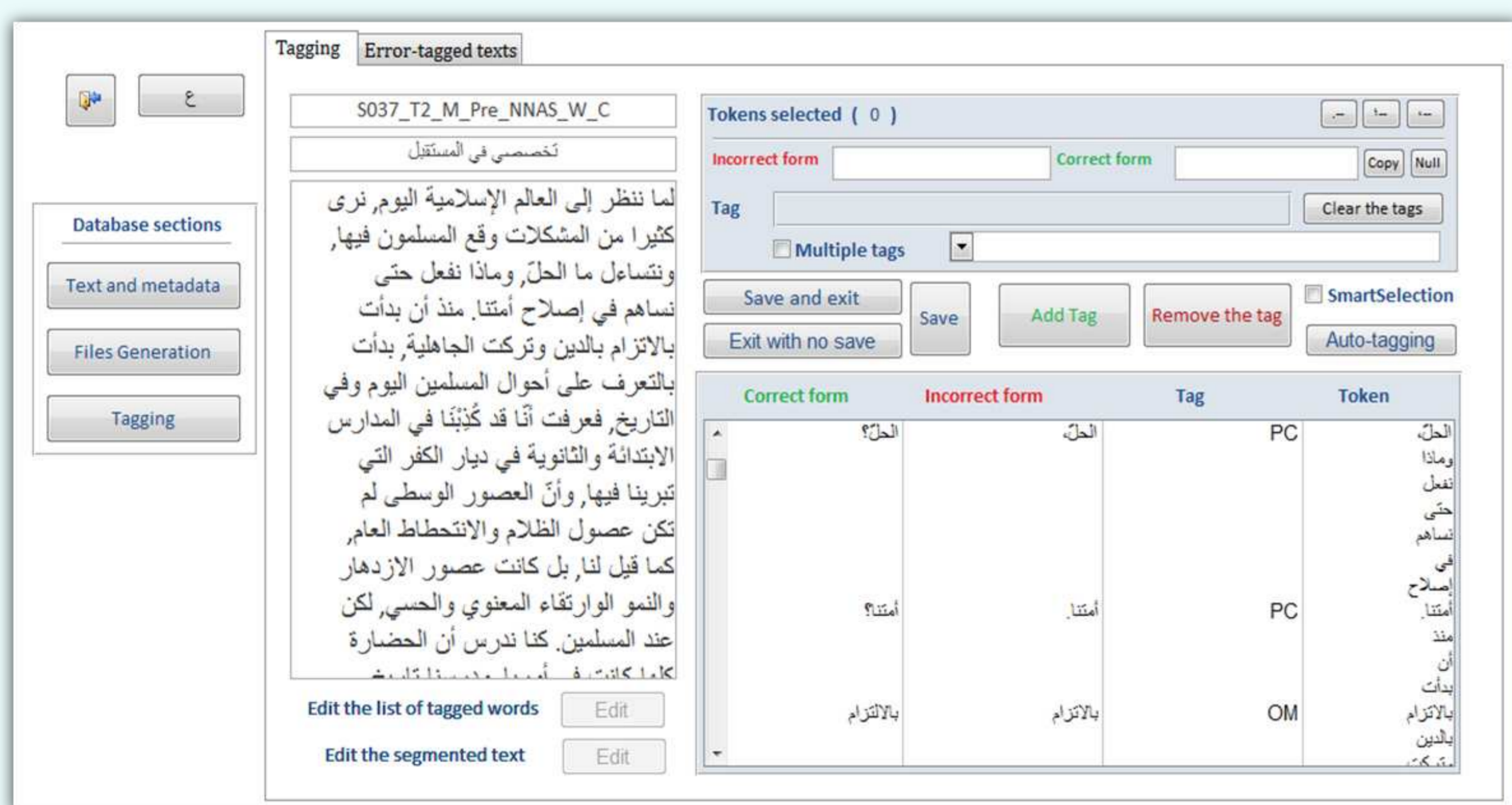
### 2. Smart-Selection

When an error exists more than once, the *Smart-Selection* function enables the annotator to find them all and tag them in a single step with no need to repeat the annotation process with each error.

### 3. Auto-Tagging

The *Auto-Tagging* feature, which is similar to translation memories, recognises the tokens that have been manually annotated and stores them in a database. So, tokens with the same error can be detected and annotated automatically next time.

For further information about ALC please visit: www.arabiclearnercorpus.com

UNIVERSITY OF LEEDS