# 1   Project description

## (a)   Description

The goal of this project is to be able to find the original language a text was written in, from the translated English version. We believe that even with completely bilingual translators, languages have a specific structure that will influence the translator in the way he creates the translation. We hope to recover this structure during our task.

## (b)   Methodology

From Project Gutenberg [], we built a dataset of books written in different languages, and translated into English. We then cut these books into slices of 500 sentences, and compute some features on these slices. After feature extraction, we feed these features into a classifier and test it on a completely different set of books.

# 2   Building the dataset

## (a)   Pitfalls

Since we built it ourselves, our dataset is not big enough for us to