**DEPARTAMENTO DE INFORMÁTICA**

# Compression in Machine Translation

## MULTIMODAL SYSTEMS

CHi — CENTER FOR INFORMATICS AND INFORMATION TECHNOLOGIES

## Jorge Costa

(PhD Student)

Advised by Gabriel P. Lopes and Luís Russo.

Research focus on Algorithms and Compressed Data Structures

## Objectives

Machine Translation (MT) tasks are typically space demanding. Tasks like parallel text alignment, extraction of translation pairs or concordancer applications, demand querying over the texts, which can have sizes of 1 (or more) Gigabyte per language. The space consuming nature of text indexes like suffix trees or suffix arrays (at least 4 times the text size), makes it difficult to index the texts in main memory, which slows down the applications.

Our aim is to solve this space consumption problem, by developing a framework, or several frameworks, based on data compression, for supporting the mentioned MT tasks in main memory, without losing efficiency on query time response.
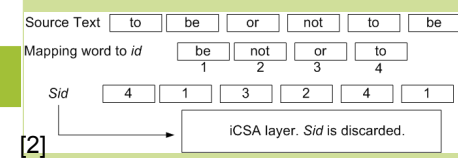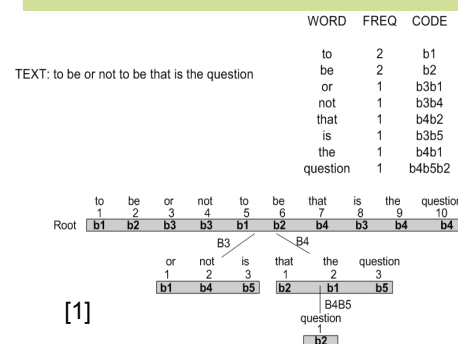


[1]

## Methodology

For development a compressed framework, we need several layers:
- **Compressed Text Indexing:** Word-Based Compressed Data Structures.
  - Byte-Oriented Wavelet Tree [1] – text words are codified in codewords and organized in a wavelet tree-like fashion.
  - Word-Based Compressed Suffix Array [2] – words replaced by integer ids (Sid) and added to an integer compressed suffix array (iCSA).
- **Alignment Layer:** represent the links between aligned segments [3], independently of the granularity, which can be at sentence or word level. This representation works as well for bilingual lexicon entries [4].
- **Document Layer:** represent the several documents indexed and their home directories, distinguishing them in the compressed index and enabling the filtering of the occurrences to certain documents.



[2]



[3]

## Expected Results

- Space consumption lower than the size of the texts indexed.
- Logarithmic query response, as the known text indexes, where a potential slowdown due to compression is compensated by avoiding secondary memory.
- Ability to index several texts at the same time in one single index, reducing the number of loading operations from disk to main memory.
- Ability to index and query not only aligned parallel corpora, but also bilingual lexicons and monolingual text collections for capturing language models.
- Support gapped alignments (Hierarchical Phrase Based Translation) and gapped bilingual lexicon entries, with variables ($) between string literals as in Figure [5].



[4]



[5]