# LING-L 445: Computation and Linguistic Analysis Practical 01B

Instructor: Francis M. Tyers

**Dante Razo**

February 11, 2019

Dante Razo

# Segmentation

**How should you segment sentences with semicolons? As a single sentence or as two sentences? Should it depend on context?**

A semicolon conjoins two independent clauses with a common idea. They can usually be rewritten as two separate sentences, but I feel it's more appropriate to keep them as one.

**Should sentence with ellipsis (...) be treated as a single sentence or as several sentences?**

In most cases, ellipsis denote a skip. In online conversation or texts, they can represent a pause for thought. In both cases, it would be inappropriate to

**If there is an exclamation after the first word in a sentence should it be a separate sentence? How about if there is a comma?**

**Can you think of some hard tasks for the segmenter?**

# Tokenization

**Why should we split punctuation from the token it goes with?**

**Should abbreviations with space in them be written as a single token or two tokens? How about numerals like 134 000?**

**If you have a case suffix following punctuation, how should it be tokenized ?**

**Should contractions and clitics be a single token or two (or more) tokens?**