# Statistical Language Modeling with NLTK

*Natural Language Understanding Lab*

Evgeny A. Stepanov,
Mahed Mousavi, Gabriel Roccabruna

SISL, DISI, UniTN & VUI, Inc.
evgeny.stepanov@unitn.it

# Objectives

- Understanding:
  - relation between lexicon and language model
  - smoothing techniques
  - OOV effects and remedies
- Learning how to:
  - prepare data for ngram modeling
  - count ngrams in a corpus
  - train an ngram model with NLTK
  - use ngram model to
    - compute probability of a sequence
    - generate a sequence
  - evaluate ngram model

# Outline

# Outline

# Recommended Reading

- Dan Jurafsky and James H. Martin. Speech and Language Processing (3rd ed. draft)
  - Chapter 3: N-gram Language Models