# Vector Semantics

*Natural Language Understanding Lab*

Evgeny A. Stepanov,
Mahed Mousavi, Gabriel Roccabruna

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

# Objectives

- Understanding:
  - different methods of representing words as vectors
  - vectors and similarity between vectors
  - evaluation of word embeddings
- Learning how to:
  - train word embeddings with `gensim`
  - use pre-trained word embeddings for similarity computation

# Outline

## Outline

# Recommended Reading

- Dan Jurafsky and James H. Martin. Speech and Language Processing (3rd ed. draft)
  - Chapter 6: Vector Semantics and Embeddings