# **Course:** Natural language processing

#### 2024/25, Spring semester

Lecturer: Prof. Dr. Marko Robnik-Šikonja Assistants: Assoc. Prof. Dr. Slavko Žitnik, Assist. Aleš Žagar

*Course objectives*: Learn about the theory and main practical approaches in natural language processing and understanding. Use modern large language models and statistical techniques for language processing.

### **Student's obligations:**

- five web quizzes
- assignments
- written exam

## Grading

The practical work encompasses work with natural language processing tools and large language models. It is graded through assignments, which have to be finished on time. The assignments are done in groups of three students. The topics of the assignments are set at the start of the semester. The results of the assignments shall be described in a paper and publicly presented in front of the class.

The exam is in the form of a written test. The preconditions for the written exam are successfully passing quizzes and assignments. During the written exam, students are allowed to use one double-sided A4 sheet of paper. In case of any doubts about the score of assignments, quizzes, or written exams, the oral exam is obligatory. The final grade is the sum of assignment scores and the written exam. The contribution of each part to the final grade and the conditions to pass the exam are listed below:

Work	% of total	subject to
Five quizzes	0%	≥ 50% altogether
Three assignments	50%	≥ 50%
Written exam	50%	≥ 50%

The grades are valid in the current year. The students not passing the exam in the current year have to redo the quizzes and assignments next year.

## **Syllabus**

The syllabus is based on a selection of modern natural learning techniques and their practical use. The lectures introduce the main tasks and techniques and explain their operation and theoretical background. The knowledge gained during practical sessions, and seminars is applied to practical tasks using open-source tools. Students investigate and solve assignments based on real-world research and industrial problems, mainly dealing with English and Slovene.

- 1. Introduction to natural language processing: motivation, language understanding, ambiguity, traditional, statistical, and neural approaches.
- 2. Text preprocessing: normalization, string similarity, lemmatization, regular expressions and grammars.
- 3. Text similarity: measures, clustering approaches, language networks, and graphs.
- 4. Language resources: corpora, datasets, dictionaries, thesauri, semantic databases, WordNet.
- 5. Text representation: text similarity, sparse and dense embeddings, cosine distances, language models.
- 6. Deep neural networks for text: recurrent neural networks, CNNs for text, transformers.
- 7. Neural embeddings: word2vec, fastText, ELMo, BERT, cross-lingual embeddings, sentence, and document embeddings.
- 8. Large pretrained language models: BERT, GPT, and T5 families, multimodal models.
- 9. Large generative language models, prompt engineering, and retrieval augmented generation.
- 10. Linguistic tasks: part-of-speech tagging, dependency parsing, named entity recognition, word sense disambiguation.
- 11. Affective computing: sentiment, emotions, etc.
- 12. LLM-based tasks: text summarization, question answering, machine translation.
- 13. Semantic tasks: natural language inference, commonsense reasoning, paraphrasing, reasoning.
- 14. Advanced topics: long context in LLMs, agents with LLMs, alternatives to transformers.

#### Literature (all freely available):

1. Jurafsky, David and Martin, James H. Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics and Speech Recognition, 3rd edition draft, 2025.

This is the primary course literature; available on <u>authors' webpages</u>

 Simon J.D. Prince: Understanding Deep Learning. MIT Press, 2023 <u>https://udlbook.github.io/udlbook/</u> This recent book intuitively explains many important concepts and techniques in deep learning.