

Subject name/code:	<b>Introduction to language technology / BNY-DK-155A</b>
Subject coordinator:	Head of the Doctoral School
Lecturer(s) of the subject:	Noémi Ligeti-Nagy, PhD
Credits:	8 credits
Lesson type:	lecture
Brief subject description:	<p>The course provides a practical overview of computational methods used in empirical linguistic analysis and natural language processing. It introduces central concepts and techniques in computational linguistics, language technology, and NLP, including regular expressions, corpus creation, linguistic annotation, benchmarking, computational semantics, machine translation, neural language models, large language models, and prompt programming. The course combines theoretical discussion with hands-on work. Students read selected chapters and papers, complete individual weekly assignments, and work on a group project in which they build a small annotated corpus and use it for experimentation with language-model fine-tuning. By the end of the course, students will understand the main methodological principles of computational linguistics and will be able to design, document, annotate, evaluate, and present a small-scale language technology project.</p>
Theoretical knowledge to be acquired:	<p>Students will acquire theoretical knowledge of the main areas of computational linguistics and natural language processing, including:</p> <ul style="list-style-type: none"> <li>• the relationship between linguistics, computational linguistics, language technology, and NLP;</li> <li>• regular expressions and their role in text processing;</li> <li>• corpora, corpus design, data collection, and benchmarking;</li> <li>• annotation theory, annotation guidelines, and inter-annotator agreement;</li> <li>• basic concepts of computational semantics, including word meaning, semantic similarity, word sense disambiguation, and semantic resources;</li> <li>• supervised and unsupervised learning in linguistic analysis;</li> <li>• machine translation and neural approaches to language processing;</li> <li>• language modelling, pretraining, fine-tuning, and large language models;</li> <li>• prompt engineering and the practical use of generative AI tools;</li> <li>• methodological principles of evaluation in NLP and empirical language technology.</li> </ul>

<p>Practical knowledge to be acquired:</p>	<p>Students will acquire practical skills in:</p> <ul style="list-style-type: none"> <li>• using regular expressions for text processing tasks;</li> <li>• collecting, organizing, and documenting linguistic data;</li> <li>• designing annotation categories and annotation guidelines;</li> <li>• carrying out pilot annotation and improving guidelines based on the results;</li> <li>• evaluating annotation consistency and dataset quality;</li> <li>• preparing a small corpus for a computational linguistics task;</li> <li>• interpreting benchmark datasets and evaluation methods;</li> <li>• experimenting with language-model fine-tuning or related language technology workflows;</li> <li>• using selected NLP tools and online resources;</li> <li>• writing a short research paper presenting the literature, methodology, data, and results of a language technology project;</li> <li>• presenting project progress and final results orally.</li> </ul>	
<p>List of the most important required literature (2–4 pieces) with bibliographical details (author, title, edition or specific pages, ISBN)</p>	<p>Jurafsky, Daniel, and James H. Martin. <i>Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition</i>. 3rd ed. draft. Selected chapters: Chapter 2, Sections 2.1–2.2; Chapter 3, Sections 3.1–3.2.1; Chapter 4, Introduction; Chapter 6, pp. 1–10. Available online.</p> <p>Jurafsky, Daniel, and James H. Martin. <i>Speech and Language Processing</i>. Earlier edition, 2014 version. Chapter 1.</p> <p>Wang, Alex, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2018. “GLUE: A Multi-Task Benchmark and Analysis Platform for Natural Language Understanding.” <i>Proceedings of the 2018 EMNLP Workshop BlackboxNLP</i>, 353–355.</p> <p>Wang, Alex, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019. “SuperGLUE: A Stickier Benchmark for General-Purpose Language Understanding Systems.” <i>Advances in Neural Information Processing Systems</i> 32.</p>	
<p>List of the most important recommended literature (2–4 pieces) with bibliographical details (author, title, edition or specific pages, ISBN)</p>	<p>Wissler, Lars, Mohammed Almashraee, Dagmar Monett, and Adrian Paschke. 2014. “The Gold Standard in Corpus Annotation.” <i>Proceedings of the 5th IEEE Germany Student Conference</i>.</p> <p>Yeomans, Michael, Alejandro Kantor, and Dustin Tingley. 2018. “The Politeness Package: Detecting Politeness in Natural Language.” <i>The R Journal</i> 10(2): 489–502.</p> <p>Arafa-Hilal, Marwa. 2023. “The Use of Politeness Strategies in Academic Conversations as Represented in a Corpus Linguistics MOOC.” <i>Journal of Pragmatic Research</i> 5(1): 85–106.</p> <p>DeepLearning.AI. <i>ChatGPT Prompt Engineering for Developers</i>. Short course. Recommended for the class on prompt programming and applied work with ChatGPT.</p>	
<p>Theory to practice ratio: 10 % practice – 90 % theory</p>	<p>Theory lessons: 100%</p>	<p>Practice lessons: 0%</p>

<p>Applied teaching methods:</p>	<p>The course is taught as a seminar combining short lectures, guided discussion, reading-based preparation, practical exercises, and project-based learning. Teaching methods include:</p> <ul style="list-style-type: none"> <li>• introductory lectures on core concepts and methods;</li> <li>• discussion of required readings and selected research papers;</li> <li>• in-class practical exercises, including regular-expression tasks and small-scale NLP activities;</li> <li>• weekly individual assignments;</li> <li>• group work on corpus creation, annotation, and experimentation;</li> <li>• pilot annotation and discussion of annotation guidelines;</li> <li>• student presentations on benchmark datasets and project progress;</li> <li>• guest lecture on machine translation;</li> <li>• guided exploration of ChatGPT and prompt engineering;</li> <li>• final oral presentation of project results.</li> </ul> <p>Students are expected to bring a notebook/laptop to class, as some sessions include practical computational work.</p>
<p>Method of assessment:</p>	<p>Assessment is based on continuous coursework and a final group project. Students complete regular individual assignments during the semester, approximately one assignment per week. The final project consists of building a small corpus for a selected computational linguistics task, documenting the annotation process, evaluating the dataset and/or model results, and presenting the work.</p> <p>The final submission is a 4–6-page paper summarizing the relevant literature, the project methodology, the data and annotation process, and the results obtained. Students also present their project progress during the semester and give a final presentation at the end of the course.</p>
<p>Assessment criteria:</p>	<ul style="list-style-type: none"> <li>- 80–100% = Excellent (5)</li> <li>- 60–79% = Satisfactory (3)</li> <li>- 0–59% = Unsatisfactory (1)</li> </ul>
<p>How the subject contributes to the achievement of the learning outcomes at level 8 of the HQF (MKKR), as identified as learning outcomes in the doctoral school's training programme:</p>	<p><b>Knowledge:</b> The subject contributes to doctoral-level knowledge by introducing students to the theoretical foundations and current methods of computational linguistics and empirical language technology. Students learn how linguistic questions can be operationalized as computational tasks, how corpora and benchmarks are constructed, and how language models and NLP systems are evaluated. The course also familiarizes students with recent developments in neural language modelling, fine-tuning, and large language models, enabling them to situate their own research in relation to contemporary computational approaches.</p> <p><b>Skills:</b> The course develops students' ability to design and implement a small-scale empirical language technology project. Students learn to collect and structure linguistic data, formulate annotation categories, write annotation guidelines, conduct pilot annotation, assess annotation quality, interpret benchmark datasets, and present methodological decisions in academic form. The final paper and presentation strengthen students' ability to communicate research design, data, methods, and results clearly to a scholarly audience.</p>

	<p><b>Attitudes:</b> The subject promotes a critical and reflective attitude toward computational methods in linguistic research. Students are encouraged to evaluate the reliability, limitations, and biases of corpora, annotations, benchmarks, language models, and generative AI systems. The course also supports openness to interdisciplinary methods, combining linguistic expertise with computational thinking, empirical evaluation, and collaborative research practice.</p> <p><b>Responsibility and autonomy:</b> The course supports doctoral-level autonomy by requiring students to make independent decisions about data selection, annotation design, methodological choices, and interpretation of results within a supervised project framework. Students take responsibility for the quality, documentation, and reproducibility of their work, while also collaborating with peers on a shared corpus-building and annotation task. Through weekly progress reports and the final paper, students practice accountable, transparent, and research-oriented project work.</p>
--	---