

Tantárgy neve/kódja:	BNY-DK-200 Bevezetés a nyelvtechnológiába
Tárgyfelelős:	A doktori iskola vezetője
Tantárgy oktatója:	Ligeti-Nagy Noémi, PhD
Kreditszám:	8 kredit
Óratípus:	szeminárium
A tantárgy céljának rövid ismertetése:	<p>A kurzus gyakorlati áttekintést nyújt az empirikus nyelvészeti elemzésben és a természetesnyelv-feldolgozásban alkalmazott számítógépes módszerekről. Bemutatja a számítógépes nyelvészet, a nyelvtechnológia és az NLP központi fogalmait és technikáit, beleértve a reguláris kifejezéseket, a korpuszépítést, a nyelvi annotációt, a benchmarkokat, a számítógépes szemantikát, a gépi fordítást, a neurális nyelvi modelleket, a nagy nyelvi modelleket és a promptprogramozást. A kurzus az elméleti megbeszéléseket gyakorlati munkával ötvözi. A hallgatók kijelölt fejezeteket és tanulmányokat olvasnak, heti egyéni feladatokat készítenek, valamint csoportprojektben vesznek részt, amely során egy kis annotált korpuszt hoznak létre, és azt nyelvmodell-finomhangolási kísérletekhez használják. A kurzus végére a hallgatók megértik a számítógépes nyelvészet fő módszertani alapelveit, és képesek lesznek egy kis léptékű nyelvtechnológiai projekt megtervezésére, dokumentálására, annotálására, értékelésére és bemutatására.</p>
Elsajátítandó elméleti ismeretanyag:	<p>A hallgatók elméleti ismereteket szereznek a számítógépes nyelvészet és a természetesnyelv-feldolgozás főbb területeiről, beleértve:</p> <ul style="list-style-type: none"> <li>• a nyelvészet, a számítógépes nyelvészet, a nyelvtechnológia és az NLP kapcsolatát;</li> <li>• a reguláris kifejezéseket és szerepüket a szövegfeldolgozásban;</li> <li>• a korpuszokat, a korpusztervezést, az adatgyűjtést és a benchmarkokat;</li> <li>• az annotációelméletet, annotációs irányelveket és az annotátorok közötti egyezést;</li> <li>• a számítógépes szemantika alapfogalmait, beleértve a szójelentést, szemantikai hasonlóságot, szóértelmezési egyértelműsítést és szemantikai erőforrásokat;</li> <li>• a felügyelt és felügyelet nélküli tanulást a nyelvészeti elemzésben;</li> <li>• a gépi fordítást és a neurális megközelítéseket a nyelvfeldolgozásban;</li> <li>• a nyelvi modellezést, előtanítást, finomhangolást és a nagy nyelvi modelleket;</li> <li>• a prompt engineeringet és a generatív MI-eszközök gyakorlati használatát;</li> <li>• az NLP és az empirikus nyelvtechnológia értékelési módszertani alapelveit.</li> </ul>

<p>Elsajátítandó gyakorlati ismeretanyag:</p>	<p>A hallgatók gyakorlati készségeket szereznek az alábbi területeken:</p> <ul style="list-style-type: none"> <li>• reguláris kifejezések használata szövegfeldolgozási feladatokhoz;</li> <li>• nyelvi adatok gyűjtése, rendszerezése és dokumentálása;</li> <li>• annotációs kategóriák és annotációs irányelvek tervezése;</li> <li>• pilot annotáció végrehajtása és az irányelvek fejlesztése az eredmények alapján;</li> <li>• annotációs konzisztencia és adatkészlet-minőség értékelése;</li> <li>• kis korpusz előkészítése számítógépes nyelvészeti feladathoz;</li> <li>• benchmark adatkészletek és értékelési módszerek értelmezése;</li> <li>• nyelvmodell-finomhangolási vagy kapcsolódó nyelvtechnológiai munkafolyamatok kipróbálása;</li> <li>• kiválasztott NLP-eszközök és online források használata;</li> <li>• rövid kutatási tanulmány írása, amely bemutatja egy nyelvtechnológiai projekt szakirodalmát, módszertanát, adatait és eredményeit;</li> <li>• projektelőrehaladás és végső eredmények szóbeli bemutatása.</li> </ul>
<p>A 2-4 legfontosabb kötelező irodalom felsorolása bibliográfiai adatokkal (szerző, cím, kiadás adatai, (esetleg oldalak), ISBN)</p>	<p>Jurafsky, Daniel, és James H. Martin. <i>Speech and Language Processing: An Introduction to Natural Language Processing, Computational Linguistics, and Speech Recognition</i>. 3. kiadás, tervezet. Kijelölt fejezetek: 2. fejezet 2.1–2.2 szakasz; 3. fejezet 3.1–3.2.1 szakasz; 4. fejezet, Bevezetés; 6. fejezet 1–10. oldal. Online elérhető.</p> <p>Jurafsky, Daniel, és James H. Martin. <i>Speech and Language Processing</i>. Korábbi kiadás, 2014-es verzió. 1. fejezet.</p> <p>Wang, Alex, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy és 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 és 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>A 2-4 legfontosabb ajánlott irodalom felsorolása bibliográfiai adatokkal (szerző, cím, kiadás adatai, (esetleg oldalak), ISBN)</p>	<p>Wissler, Lars, Mohammed Almashraee, Dagmar Monett és 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 és 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>. Rövid kurzus. Ajánlott a promptprogramozásról és a ChatGPT gyakorlati használatáról szóló órához.</p>

Elmélet-gyakorlat aránya:	Elméleti óra óraszám: 0%	Gyakorlati óra óraszám: 100%
Az alkalmazott oktatási módszerek:	<p>A kurzus szemináriumi formában zajlik, amely rövid előadásokat, irányított vitát, olvasmányalapú felkészülést, gyakorlati feladatokat és projektalapú tanulást ötvöz. Az alkalmazott módszerek:</p> <ul style="list-style-type: none"> <li>• bevezető előadások az alapfogalmakról és módszerekről;</li> <li>• kötelező olvasmányok és kiválasztott kutatási tanulmányok megbeszélése;</li> <li>• órai gyakorlati feladatok, beleértve reguláris kifejezés-feladatokat és kis léptékű NLP-tevékenységeket;</li> <li>• heti egyéni feladatok;</li> <li>• csoportmunka korpuszépítésen, annotáción és kísérletezésen;</li> <li>• pilot annotáció és annotációs irányelvek megvitatása;</li> <li>• hallgatói prezentációk benchmark adatkészletekről és projektelőrehaladásról;</li> <li>• vendégelőadás a gépi fordításról;</li> <li>• ChatGPT és prompt engineering irányított kipróbálása;</li> <li>• a projekt eredményeinek záró szóbeli bemutatása.</li> </ul> <p>A hallgatóknak notebookot/laptopot kell hozniuk az órákra, mivel egyes alkalmak gyakorlati számítógépes munkát tartalmaznak.</p>	
Az értékelés módja:	<p>Az értékelés folyamatos félévi munkán és egy záró csoportprojektben alapul. A hallgatók rendszeres egyéni feladatokat készítenek a szemeszter során, körülbelül heti egy feladatot. A záróprojekt egy kis korpusz létrehozásából áll egy kiválasztott számítógépes nyelvészeti feladathoz, az annotációs folyamat dokumentálásából, az adatkészlet és/vagy a modell eredményeinek értékeléséből, valamint a munka bemutatásából. A végső beadandó egy 4–6 oldalas tanulmány, amely összefoglalja a releváns szakirodalmat, a projekt módszertanát, az adat- és annotációs folyamatot, valamint az elért eredményeket. A hallgatók a szemeszter során projektelőrehaladási beszámolókat is tartanak, és a kurzus végén záró prezentációt adnak elő.</p>	
Az értékelés kritériuma:	<p>Osztályzás háromfokú skálán:</p> <ul style="list-style-type: none"> <li>- 80-100% = Kiválóan megfelelt (5)</li> <li>- 60–79% = Megfelelt (3)</li> <li>- 0–59% = Nem felelt meg (1)</li> </ul>	
Miként járul hozzá a tantárgy a doktori iskola képzési programjában tanulási eredményként megjelölt MKKR 8. szintű tanulási eredmények eléréséhez.	<p>Tudás: A tantárgy doktori szintű tudást nyújt a számítógépes nyelvészet és az empirikus nyelvtechnológia elméleti alapjairól és aktuális módszereiről. A hallgatók megtanulják, hogyan operacionalizálhatók a nyelvészeti kérdések számítási feladatokká, hogyan épülnek fel a korpuszok és benchmarkok, valamint hogyan értékelik a nyelvi modelleket és NLP-rendszereket. A kurzus megismerteti a hallgatókat a neurális nyelvi modellezés, finomhangolás és nagy nyelvi modellek legújabb fejleményeivel is, lehetővé téve számukra, hogy saját kutatásukat a kortárs számítógépes megközelítések kontextusában helyezték el.</p> <p>Képességek: A kurzus fejleszti a hallgatók képességét egy kis léptékű empirikus nyelvtechnológiai projekt megtervezésére és</p>	

	<p>megvalósítására. A hallgatók megtanulják a nyelvi adatok gyűjtését és strukturálását, annotációs kategóriák kialakítását, annotációs irányelvek írását, pilot annotáció végzését, annotációs minőség értékelését, benchmark adatkészletek értelmezését, valamint módszertani döntések tudományos formában történő bemutatását. A zárótanulmány és prezentáció erősíti a kutatási terv, adatok, módszerek és eredmények világos tudományos kommunikációját.</p> <p>Attitűdök: A tantárgy kritikus és reflektív hozzáállást alakít ki a számítógépes módszerekkel kapcsolatban a nyelvészeti kutatásban. A hallgatókat arra ösztönzi, hogy értékeljék a korpuszok, annotációk, benchmarkok, nyelvi modellek és generatív MI-rendszerek megbízhatóságát, korlátait és torzításait. A kurzus támogatja az interdiszciplináris módszerek iránti nyitottságot is, ötvözve a nyelvészeti szakértelmet a számítógépes gondolkodással, empirikus értékeléssel és együttműködésen alapuló kutatási gyakorlattal.</p> <p>Felelősség és autonómia: A kurzus támogatja a doktori szintű autonómiát azáltal, hogy a hallgatóknak önálló döntéseket kell hozniuk az adatválasztásról, annotációtervezésről, módszertani döntésekről és az eredmények értelmezéséről egy felügyelt projektkeretben. A hallgatók felelősséget vállalnak munkájuk minőségéért, dokumentálásáért és reprodukálhatóságáért, miközben társaikkal együttműködnek egy közös korpuszépítési és annotációs feladaton. A heti előrehaladási jelentések és a zárótanulmány révén elsajátítják az átlátható, felelősségteljes és kutatóorientált projektmunka gyakorlatát.</p>
--	---