

SZIRMAI MONIKA

BEVEZETÉS A KORPUSZNYELVÉSZETBE



SEGÉDKÖNYVEK  
A NYELVÉSZET TANULMÁNYOZÁSÁHOZ XLVI.

SZIRMAI MONIKA

BEVEZETÉS  
A KORPUSZNYELVÉSZETBE

*A korpusznyelvészeti alkalmazása az anyanyelv  
és az idegen nyelv tanulásában és tanításában*

TINTA KÖNYVKIADÓ  
BUDAPEST, 2005

SEGÉDKÖNYVEK  
A NYELVÉSZET TANULMÁNYOZÁSHOZ XLVI.

Sorozatszerkesztő:  
KISS GÁBOR

Lektor:  
ANDOR JÓZSEF

ISSN 1419-6603  
ISBN 963 7094 42 3

© Szirmai Monika, 2005

A kiadásért felel  
a TINTA Könyvkiadó igazgatója  
Felelős szerkesztő: Temesi Viola  
Műszaki szerkesztő: Bagu László

## TARTALOM

Ábrák jegyzéke .....	9
Táblázatok.....	11
Köszönetnyilvánítás.....	12
Kinek szól ez a könyv? .....	13
Bevezetés .....	15
1. Mi a korpusznyelvészet?.....	17
1.1. Bevezetés .....	17
1.2. Mi a korpusz? .....	17
1.3. A korpusz tervezése.....	23
1.3.1. A reprezentativitás .....	23
1.3.1.1. Mintavétel.....	23
1.3.1.2. A korpusz mérete .....	27
1.3.2. A korpuszok fajtái.....	32
1.3.2.1. A mintavétel módja szerint .....	32
1.3.2.2. A korpusz felhasználásának módja szerint .....	32
1.3.3. Jogi problémák .....	36
1.3.4. Átírás és annotáció .....	37
1.3.4.1. A beszéd átírása .....	37
1.3.4.2. A standard annotáció.....	38
1.3.4.3. Speciális annotációk .....	43
1.4. Összefoglalás .....	45
2. Számítástechnika és nyelvtudomány .....	47
2.1. Bevezetés .....	47
2.2. A számítástechnika fejlődése .....	47
2.3. A korpuszok fejlődése .....	50
2.3.1. A szellemi háttér .....	51
2.3.2. A korpusznyelvészet és a kapcsolódó tudományágak .....	53
2.3.2.1. A számítógépes nyelvészeti .....	53
2.3.2.2. A mesterséges intelligencia.....	53
2.3.2.3. A számítógépes nyelvészeti kutatási területe .....	54
2.3.3. A magyarországi számítógépes nyelvészetről .....	55
2.4. Folyóiratok .....	58
2.5. Összefoglalás .....	59

3. A korpuszokról.....	60
3.1. Bevezetés .....	60
3.2. Az elektronikus korpuszok előfutárai .....	60
3.2.1. A Szerb Nyelv Korpusza .....	60
3.2.2. A SEU Korpusz (Survey of English Usage Corpus) .....	61
3.3. A Brown Korpusz (1964).....	63
3.4. A LOB Korpusz .....	63
3.5. A COBUILD projekt.....	64
3.6. A Brit Nemzeti Korpusz – British National Corpus (BNC).....	66
3.7. Az Angol Nyelv Nemzetközi Korpusza (International Corpus of English – ICE) .....	67
3.8. A nem anyanyelvi angol korpuszok .....	69
3.8.1. A Longman Angol Nyelvtanulói Korpusz – Longman Corpus of Learners' English (LCLE).....	70
3.8.2. A Nemzetközi Angol Nyelvtanulói Korpusz (International Corpus of Learners' English (ICLE)) .....	70
3.8.3. A Hongkongi Műszaki és Természettudományi Egyetem Angol Tanulói Korpusza (Hong Kong University of Science and Technology [HKUST] Corpus of Learner English).....	72
3.8.4. Japán diákok angol nyelvű korpuszai.....	72
3.8.5. A Janus Pannonius Tudományegyetem Korpusza .....	72
3.8.6. Az Eötvös Loránd Tudományegyetem Korpusza .....	73
3.9. A korpuszok nyelvenként .....	73
3.9.1. További angol nyelvű korpuszok.....	74
3.9.1.1. A Brown Korpusz klónjai .....	74
3.9.1.2. Könyvkiadók korpuszai .....	75
3.9.1.3. Történeti nyelvészeti korpuszok .....	77
3.10. Az angol nyelvű korpuszok áttekintése.....	78
3.11. Magyar nyelvű korpuszok .....	80
3.11.1. A Magyar Nemzeti Szövegtár (MNSZ) .....	80
3.11.2. A Magyar Irodalmi és Köznyelv Nagyszótárának korpusza / Magyar Történeti Korpusz .....	81
3.11.3. Szeged Korpusz ( <a href="http://www.inf.u-szeged.hu/projectdirs/hlt/">http://www.inf.u-szeged.hu/projectdirs/hlt/</a> ) .....	86
3.11.4. Magyar dalszövegek .....	87
3.11.5. CHILDES Database: <a href="http://childepsy.cmu.edu/">http://childepsy.cmu.edu/</a> magyar nyelvű korpusza .....	87
3.11.6. A Hunglish Korpusz .....	87
3.11.7. A Magyar Webkorpusz .....	87
3.12. Egyéb nyelvek korpuszai .....	88
3.12.1. Német nyelvű korpuszok.....	89
3.12.1.1. A negr@ korpusz .....	89
3.12.1.2. A Tiger Korpusz.....	89
3.12.1.3. Freiburger Korpus .....	90
3.12.1.4. Dialogstrukturenkorpus .....	90
3.12.1.5. Pfeffer-Korpus .....	90

3.12.1.6. Telefonbeszélgetések (Brons-Albert, 1984) .....	90
3.12.2. Francia nyelvű korpuszok .....	91
3.12.2.1. PAROLE Francia Korpusz http://www.elda.org/catalogue/en/text/W0020.html.....	91
3.12.2.2. Francia Beszélt Nyelvi Korpusz.....	91
3.12.2.3. Kanadai Francia Korpusz .....	91
3.12.2.4. Le corpus VALIFLOUI (Variétés Linguistiques du Français en Louisiane) http://languages.louisiana.edu/French/Valifloui.html University of Louisiana at Lafayette .....	92
3.12.2.5. Le Corpus du Théâtre religieux français du Moyen Âge (Középkori Francia Vallásos Színház Korpusza).....	93
3.12.3. A Szerb Nyelv Korpusza .....	93
3.12.4. A horvát nyelv korpusza.....	95
3.12.5. Szlovén nyelvű korpuszok.....	97
3.12.5.1. Szlovén – FIDA .....	97
3.12.5.2. BESEDA.....	97
3.12.6. Cseh nyelvű korpuszok .....	98
3.12.7. Lengyel korpuszok .....	98
3.13. Összefoglalás .....	98
4. A szoftverekről.....	100
4.1. Bevezetés .....	100
4.2. A korpuszok készítésekor használt programok .....	100
4.3. A konkordancia programok.....	102
4.3.1. A kezdet kezdetén.....	110
4.3.2. Internetes felületen futó ingyenes programok .....	111
4.4. Konkordanciák készítése.....	112
4.4.1. Az MLCT .....	115
4.4.2. Simple Concordance Program SCP.....	120
4.4.3. ConcApp.....	126
4.4.4. AntConc.....	127
4.5. Összefoglalás .....	129
5. Korpusznyelvészeti módszerek az oktatásban .....	130
5.1. Bevezetés .....	130
5.2. Konferenciák és publikációk.....	130
5.3. Számítógéppel és nélküle .....	132
5.4. A késztermékek .....	133
5.4.1. Az egynyelvű tanulói szótárakról .....	133
5.4.1.1. Bevezetés .....	133
5.4.1.2. Anyanyelvi szótár – tanulói szótár.....	134
5.4.2. COBUILD kiadványok .....	139
5.4.2.1. Tankönyv .....	139
5.4.2.2. Segédanyagok .....	139

5.4.3. A Longman Grammar of Spoken and Written English (LGSWE) .....	142
5.4.4. Touchstone – új korpusz alapú tankönyv .....	143
5.5. Saját készítésű feladatok .....	144
5.5.1. Konkordanciák nyomtatásban .....	144
5.5.2. A „számok tükrében” .....	156
5.6. Számítógépes feladatok.....	157
5.7. Összefoglalás .....	162
A könyvben szereplő nyelvi korpuszok, szövegtárak és adatbázisok .....	164
Korpusznyelvészeti alapfogalmak .....	167
Bibliográfia .....	175
Tárgymutató.....	184
Névmutató.....	186
Korpuszok mutatója .....	189

## KÖSZÖNETNYILVÁNÍTÁS

E könyv soha nem jöhetett volna létre Cseresnyési László biztatása és segítsége nélkül. Szintén neki köszönhetem Kiss Gáborral, a Tinta Könyvkiadó igazgatójával való megismeredésem, melynek eredményeként könyvem itt került kiadásra.

Egy könyv megjelenése nem csak a szerző érdeme, még akkor sem, ha a címlapon csak az ő neve szerepel. Mivel ez az első magyar nyelven megjelenő írásom, úgy érzem, hogy ez a legjobb alkalom arra, hogy nyilvánosan köszönjem meg a segítséget mindenaknak, akik nélkül nem lettem volna azzá, aki vagyok. Középiskolás koromtól kezdve sok kiváló szakember irányított és segített nyelvi és nyelvészeti ismereteim megszerzésében, így lehetetlen lenne mindenkit felsorolni, annak ellenére, hogy az én emlékeztemben minden élénken élnek. Vannak azonban olyanok is, akik az életem alakulására is nagy hatást gyakoroltak.

Pályaválasztásomban döntő szerepet játszott két középiskolai tanárom: Albert Sándor és Novák György. Most is hálás vagyok akkori tanácsaikért. Debreceni egyetemi éveim alatt az angol és a francia tanszék minden tanára nagy hatással volt szellemi fejlődésem alakulására. A Birminghami Egyetemen a British Council ösztöndíjasaként eltöltött egy év alatt pedig olyan kutatókkal kerültem napi kapcsolatba, mint John Sinclair, Susan Hunston, Ramesh Krishnamurthy, Dave és Jane Willis. Szintén hálával tartozom Hollósy Bélának, aki doktori tanulmányaim alatt elfoglaltsága ellenére is mindenben és mindenkor segítségemre sietett. Doktori tanulmányaim alatt ismerkedtem meg Andor Józseffel, aki hihetetlen energiával és lelkesedéssel végzi mind pedagógusi, mind kutatói munkáját. E könyvet nem csak kész állapotában lektorálta, hanem azon túlmenően, az írás megkezdésétől kezdve készségesen válaszolt kérdéseimre, és segített tanácsaival. Fenyvesi Anna és Heitzmann Judit a kéziratra vonatkozó észrevételeikkel segítették munkámat.

A szakmai segítség mellett azonban elengedhetetlen a család és a barátok türelme, támogatása és biztatása. Köszönöm szüleimnek, hogy mindig mindenben mellettem álltak és segítettek. A bajai Kerényi és Rácz családnak, a hiroshima Szaszaki Tadaszunak és Murakami Midorinak külön köszönettel is tartozom. Számítógépes problémáim megoldásában Rácz Györgyre, Schmíz Istvánra, Oláh Gyulára és Varjú Tiborra bármikor számíthattam.

Ha e sok segítség és támogatás ellenére mégis előfordulnak hiányosságok vagy pontatlanságok a könyvben, akkor az csakis rajtam műlött.

Hiroshima, 2005 tavasza

*Szirmai Monika*

## **KINEK SZÓL EZ A KÖNYV?**

Könyvemnek az a célja, hogy a lehető legszélesebb közönséggel ismertesse meg a korpusznyelvészeti alapjait és alkalmazási lehetőségeit. A könyv írásakor igyekeztem egyaránt szem előtt tartani a nyelvszakos egyetemi hallgatókat, az általános és középiskolai tanárokat, a magyar nyelvet és irodalmat vagy idegen nyelvet tanító nyelvtanárokat, valamint – kortól függetlenül – a nyelvtanulókat. Arra törekedtem, hogy a nyelvészeti zsargon elkerülve vagy megmagyarázva, egyszerű, könnyen olvasható formába öntsem gondolataimat.

A korpusznyelvészeti számos területen alkalmazható, de számomra elsősorban az a tény fontos, hogy jelentősen elősegíthesi minden az anyanyelvvel kapcsolatos ismereteink bővítését és pontosítását, minden az idegen nyelvek elsajátítását. A nyelvtanároknak vagy az anyanyelvi beszélőknek segítséget nyújthat olyan kérdések megválaszolásában, amelyekre a diákok a nyelvtankönyvekben nem találnak választ. Munkám megjelentetését az is indokolja, hogy magyar nyelvű könyv még nem készült a korpusznyelvészetről, így csak idegen nyelven, elsősorban angolul olvasható ezzel foglalkozó szakirodalom.



## BEVEZETÉS

Néha az ember életét egy-egy találkozás örökre megváltoztathatja. Ilyen fordulópont volt az én életemben az 1993–94-es tanév, amikor felnőtt fejjel ismét diákként ütem a padban, és egy előadást hallgattam, vagyis inkább egy számítógépes program bemutatóját néztem. A program neve Contexts volt (ebben a könyvben is bemutatásra kerül), az előadó pedig a Birminghami Egyetem tanára, Tim Johns. Hogy miért éppen egy számítógépes program hatott rám ilyen elemi erővel, amikor számos elméleti előadáson már túl voltam minden különösebb lelkesedés nélkül? Bizonyára az előadó személyes varázsa is közrejátszott, de elsősorban a játék izgalma és a szellemi kihívás ragadott meg mint nyelvtanulót. Mint nyelvtanárnak, azonnal az alkalmazási lehetőségek sora futott végig a fejemen, hiszen ha nekem ennyire tetszett, talán a diákjaimnak is hasznos és szórakoztató lehet.

A nagy lelkesedés nem hozott olyan gyors előrelépést, mint azt gondolhatnánk, mert sajnos abban az időben sem számítógéppel, sem pedig informatikai ismeretekkel nem rendelkeztem. Szeretném megnyugtatni az olvasót, hogy ma már nincs is szükség olyan programozási ismeretekre, mint a Windows 95 operációs rendszer megjelenése előtt, és a számítógép kezelése is egyre könnyebb lett. E könyv írásakor azonban feltételezem, hogy az olvasó rendelkezik alapvető számítógépes ismeretekkel, azaz tudja, hogy hogyan kell megnyitni egy könyvtárt vagy fájlt, tud menteni, el tud indítani egy programot, tudja kezelní a menürendszeret stb.

A modern korpusznyelvészeti feltételezi a számítógép használatát, de nem azonos a számítógépes nyelvészettel annak ellenére, hogy számos átfedés található a kettő között. A kutatás céljai és az eredmények felhasználásának lehetőségei is különböznek. A programozási ismeretek a korpusznyelvészeti művelésekor is sokat segíthetnek, de a számítógépes nyelvészeti elengedhetetlenek.

A nyelvtanulás számomra olyan, mint a felfedezés. Akár írott szövegről, akár élőbeszédről legyen szó, ha magam „fedezem fel” a szabályt vagy veszek észre valamit, az sikeres eredményt nyújt, és sokáig megmarad az emlékezetemben. A mások által közölt vagy könyvben olvasott szabályokra azonban már nem emlékszem olyan könnyen és hamar el is felejtem őket. A diákokat is kutatóknak tekintem. Minél magasabb szinten szeretnének elsajátítani egy nyelvet, annál jobb és önállóbb kutatókká kell válniuk. A diákoknak három nagyon fontos dologra van szükségük: egyszerű megfigyelőképességre, másrészt képesnek kell lenniük helyes következtetések levonására a rendelkezésre álló adatok alapján, harmadrészt pedig képesnek kell lenniük intelligens találkagatásra. Egy egyszerű példával szeretném ezt bemutatni.

Tegyük fel, hogy valaki hirtelen olyan környezetbe kerül, amelynek nyelvét (mondjuk a japánt) nem beszéli. Mivel élenken figyel, egy idő után észreveszi, hogy minden étkezés előtt hallja az *itadakimasz* kifejezést, étkezés után pedig azt, hogy *gocsiszó-*

*szamadesta*. Ennek alapján levonhatja a következtetést, hogy ezt a két kifejezést akkor használják, amikor a magyarban a *Jó étvágyat!* és az *Egészségünkre!* fordulatokat. Az *itadakimász* kifejezés jelentését a helyzetből fakadóan a *Jó étvágyat!* kifejezéssel fogja azonosítani mindaddig, amíg új, más helyzetben nem találkozik ugyanezzel a kifejezéssel, ahol nyilvánvalóan más jelentésben használják.

Azok a nyelvtanulók lesznek sikeresek, akik az ilyen helyzeti megfigyelőképességeket rendszeresen használják. Sokan automatikusan teszik ezt, másoknak fel kell erre hívni a figyelmét, és még gyakorlásra szoruló tanulók is vannak. De ezek a képességek gyakorlással fejleszthetők. Tapasztalataim szerint erre különösen alkalmasak a korpusznyelvészeti szellemében készült feladatok.

Ez a könyv tehát az olvasók széles skálájának igényeit igyekezik kielégíteni – a szakembertől kezdve az érdeklődő diákokig. Öt fejezetből áll. Az első a korpusznyelvészeti alapjait mutatja be. Mivel a számítástechnika fejlődése meghatározó volt a korpusznyelvészeti szempontjából, a második fejezetben a számítástechnika és nyelvészeti kapcsolatáról esik majd szó. Ezt követi a legfontosabb korpuszok bemutatása a harmadik fejezetben. A negyedik fejezet a korpuszok használatához szükséges számítógépes programokat ismerteti. A legtöbb figyelmet e fejezetben az úgynevezett konkordancia-programoknak szenteljük: illusztrációk és részletes leírás segítségével szeretném lehetővé tenni, hogy még az angolul nem tudó érdeklődők is kezelhessék e programokat. Az ötödik fejezet a korpusznyelvészeti oktatásban való felhasználásának lehetőségeit mutatja be. Számos mintafeladatot és alkalmazási ötletet tartalmaz, melyeket bármely nyelv tanításakor vagy tanulásakor fel lehet használni. Eddigi kutatásaim és tanítási tapasztalataim, valamint a korpusznyelvészeti szakirodalom jelentős többsége az angol nyelvre vonatkozik, így annak ellenére, hogy igyekeztem minél több esetben magyar nyelvű példát is hozni, azok nagy része angol nyelvű. Remélem, hogy így is haszonnal forgatja majd e könyvet minden olvasója.