

The new Gigaword Version of the Hungarian National Corpus

Csaba Oravecz
oravecz@nytud.hu

Language Technology Research Group,
Research Institute for Linguistics,
Hungarian Academy of Sciences

21 November 2014



- Origins
- Motivation
- Objectives
- Design considerations
- Preparation
- Outcome



The Hungarian National Corpus (HNC)

- developed between 1998 and 2001
- representative sample of the language use of the second half of the 90s → empirical evidence for status of language and data for theoretical analysis and language technology
- first major annotated Hungarian corpus, available freely through a search interface
- 187 million words, covering language variants from beyond the border of Hungary → Hungarian Minority Language Corpus
- more than 7000 registered users, dozens of research papers based on HNC data



15 years after ...

- requirements against language resources have changed significantly
 - dominance of data oriented methods and applications in NLP
 - the more data the better results
 - better language processing tools
 - higher quality and finer level of analysis and annotation
 - preservation of representativity
 - subsequent sampling from language use needed
- HNC has become outdated



Increase ...

- *quality*. Use new technology for development and analysis.
- *size*. Extend the corpus to 1 Gw.
- *coverage and representativity*. Take new samples of language use and include further variants (transcribed spoken language data in particular).

HGC: Develop an up-to-date language resource that will service the research community as well as the interested public.



- Representativity?
Illusionary goal on this scale → **balancedness**
- Collection of data by web crawling or acquiring large amounts of newswire text?
Difficult to produce a solid, balanced resource with sufficient metadata
→ controlled, targeted resource collection, appropriate for each type of source. Easily processable source format is preferred
→ no pdf, no OCR.



- Clean IPR issues?

Difficult, sometimes even impossible, to collect appropriate licenses

→ different availability options are offered for various sections of the HGC.



Text collection

- (try to) clean up IPR issues
- collect data with extensive metadata

Preprocessing, normalization

- identify textual content and basic document structure
- filter out (near-)duplicates and non-Hungarian sections



Analysis and annotation

- detailed morphosyntactic analysis and disambiguation with updated processing toolchain (information on stem, each morph and compounding)
- NP chunking, Named Entity recognition (in prep.)
- XML annotation compatible with international standards



Annotation format

az	az	DET	D__D	compound=n;;hyphenated=n;;stem=az::DET BC BC az az	B
<i>token</i>	<i>stem</i>	<i>MSD code</i>	<i>corpus tag</i>	compound=n;;hyphenated=n;;stem=angol::A;; morphemes=ZERO::NOM	
angol	angol	A.NOM	AS_A	BCCBC BCCBC angol angol	I
			<i>morpheme level encoding</i>	compound=n;;hyphenated=n;;stem=nyelv::N;; morphemes=ZERO::NOM ü::_UKEP	
nyelvü	nyelvü	A.NOM	AS_A	CNCCF CNCCF Nelvü Nelvü	I
				compound=n;;hyphenated=n;;stem=szöveg::N;; morphemes=ZERO::NOM	
szöveg	szöveg	N.NOM	NS3NN	CFCNC CFCNC Söveg Söveg	I
			<i>IOB encoding (of an NP)</i>	compound=n;;hyphenated=n;;stem=van::IGE;; stemvar=vol::IGE;;morphemes=t::Me3	
volt	van	V.Me3	VS3PI	CBCC CBC volt van O	
			<i>CV skeleton of token and stem</i>	compound=y;;hyphenated=n;;stem=ad::VERB irány::N;;morphemes=ZERO::NOM ó::_OKEP	
irányadó	irányadó	A.NOM	AS_A	NCBCBCB NCBCBCB iráNadó iráNadó O	
			<i>phonemic transcription</i>	NA NA NA NA NA	
.	.	SPUNCT	SPUNCT		



Multiple derivational paths

in possibilities:

lehetőség [N] + ek [PL] + ben [INE]

lehető [A] + ség [_PROP] + ek [PL] + ben [INE]

*lesz [V] = le + hető [_HATO] + ség [_PROP] + ek [PL] + ben [INE]

=> lehetőség [N] [PL] [INE] / NP3N2

→ select the analysis with the highest number of morphemes

Compound overgeneration

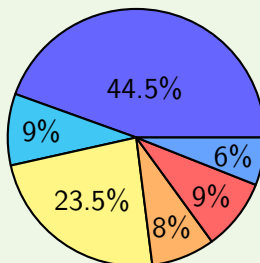
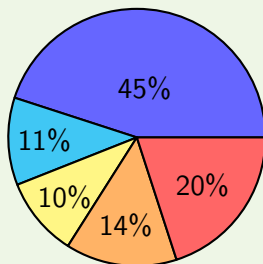
lázadó (“rebel”) = láz (“fever”) + adó (“tax”)

→ filter rules with regular expressions



HNC: 187 m.

HGC (+HNC): 1437 m.



- Journalism
- Official
- Personal
- Science
- Literature
- Spoken

Composition (m. words)

Register	HNC	HGC	Source
Journalism	84,5	639,2	Daily/weekly newspapers
Literature	38,2	130,7	Digital Literary Academy
(Popular) science	25,5	110,9	Hungarian Electronic Library
Personal	18,6	338,6	Social media
Official	20,9	135,4	Documents from public admin.
(Transcribed) spoken	–	83,0	Radio programs
	187,7	1,437,8	



Corpus query engine

- robust, able to handle several gigawords
- quick response (depending on complexity of queries)

"Intelligent" corpus

- complex searches based on every piece of information in the annotation
 - morpho(phono)logical phenomena
 - multiword expressions: collocations, verbal arguments
- display settings: context, metadata
- distributional analysis, built in post-processing (multilevel frequency lists, subsequent searches on previous results)



Searching for a specific verb form

User: **manina** Corpus: **HGC** in **HGC**

Concordance Word List [?](#)

Expert options:
Context
Text Types [?](#)

Corpus: **HGC**

Query Type: **Detailed search**

Detailed search:

word part of speech: **verb ...**

verb prefix: **ANY** **'-ha'** **yes** conjugation: **'-lak'** tense/mood: **declarative** number: **singular** person: **first**

ANY
subjective
objective
'-lak'

Lexis
Sketch Engine (ver. 0.1)
Interface (lanj)



Concordance list for *-hAt/ak* verbs

Concordance Word List



- Save
- View options
- Kwic/Sentence
- Sort
- Left | Right
- Node
- References
- Shuffle
- Sample
- Filter
- Frequency
- Node tags
- Node forms
- Doc IDs
- Collocations
- ConcDesc



Corpus: **HGC**
Hits: **753**

Page of 38 [Next](#) | [Last](#)

doc#2	<l-btw_p_resz_begin.resz-> <p> A. Milyen ized van ?	Megsókolhatlak	? Szeretsz ? Van -e öled ? Nőstényszagod
doc#2	szerelmet így tudta kifejezni : " nagyon	szerethetlek	, ha már a sálam is integer utánad "
doc#4	Gyere, gyere csak , foglaj helyet ! Mivel	hívhatlak	tól , hogy szabadon beszélhetsz velem .
doc#6	képre ? K : Á ... Á : S ha lekerülök ,	kinálhatlak	meg ? Azért ne siess ! A királylány jól
doc#8	Jó is volna , oly hosszú idő után egyszer	megkereshetlek	-e , mint az idősebb katonát , mutasd meg
doc#9	talán , talán lehet még , hiszen hiszem -	kaphatlak	meg , és te máris menni akarsz . Hisz az
doc#15	talán , talán lehet még , hiszen hiszem -	lehetlek	, s én se vagyok még emlék , és élj-se
doc#18	vessen , nem táltoslo , de táltosan szabad .	lehetlek	, s én se vagyok még emlék , és élj-se
doc#18	mosolygott . Túl udvariasan . Várt . <p><p>	Elveszithetlek	? Akár a száj a szót ? Szívdobogást a szív
doc#20	<p><p> - Meghívlak . Olyan a helyzet , hogy	Meghívhatlak	egy sörre ? - kérdezte Jócó . <p><p> - Megyek
doc#20	ide . A vörös bor talán , Ebéd után ...	meghívhatlak	. Megtaláltam az asszonyt ... Szerencsém
doc#25	, kiről beszélsz . Na megvársz ? <p><p>	Megkínálhatlak	borral is ? <p><p> - Ha nem édes . <p><p>
doc#41	befogadja . <p><p> - Jó - mondta Andrej Bodor ,	Megvárhatlak	- mondta Sivár . <p><p> Nopritz Andrea bement
doc#45	, kiről beszélsz . Na megvársz ? <p><p>	biztosíthatlak	, mindezt pont így tudatni fogom . <p> <head>
doc#46	mindennapi szökeledést irigylek ? Vagy hogy	Megvárhatlak	- mondta Sivár . <p><p> Nopritz Andrea bement
doc#47	után . Féltalomban és dagdoga hazámnak is	főtálhatlak	benneteket ? Véretek mocská ; szentkép
doc#47	Baranyl Antal vagyok , ülj le , kérlek ,	mondhatlak	téged . Visszatért lecséld itt rácsüroznak
doc#50	<title> <head> <p> Hagyjuk abba . Azt hittem ,	megkínálhatlak	egy kis disznósájtait , a vagonid igazgatója
doc#50	alakja születik meg az ő verseiben , hogy :	rászedhetlek	, hogy ellazítb magadat legalább felőrára
doc#50		lehetlek	. Azaz : én te lehetek majd , mert az szeretnék

Page of 38 [Next](#) | [Last](#)



Searching for phonological phenomena

User: **mama** Corpus: **HGC** Search in HGC

Concordance Word List

Expert options:
Context
Text Types

Corpus: **HGC**

Query Type: **Detailed search**

Detailed search: **lemma** {aff} {vow} {app} part of speech: **ANY**

vowels: ☐ all ☐ short ☐ long ☐ back ☐ neutral ☐ front Clear Add

consonants: ☐ all ☐ voiced ☐ voiceless ☐ labial ☐ alveolar ☐ palatal ☐ velar ☐ explosive ☐ spirant ☐ affricate ☐ nasal ☒ approximant ☐ trill ☐ fricative Clear Add

Make Concordance Clear All

Lexical Computing Ltd. 2005
Sketch Engine (v4r.open-2.5.9.1-open-2.9.1.1)
Interface language: [English](#) [magyar](#)



Frequency list generated from search result

User: manna Corpus: HGC

Concordance
Word List
[?](#)

Save
View
concordance
Sample
Filter
Frequency
Node tags
Node forms
Doc IDs
Collocations
ConcDesc
[?](#)

Frequency list

Frequency limit:

lemma	Freq	
p/n cél	94663	<div></div>
p/n csal	5179	<div></div>
p/n csaj	4103	<div></div>
p/n csel	1008	<div></div>
p/n col	165	
p/n cal	40	
p/n dzsal	28	
p/n cul	6	
p/n csol	6	
p/n cil	6	
p/n csáj	5	
p/n cel	4	
p/n csul	2	
p/n csej	1	
p/n cely	1	



Concordance Word List		Collocation candidates				
?		Page <input type="text" value="1"/> <input type="button" value="Go"/> Next >				
		Freq	T-score	MI	logDice	
Save		p/n lámpa	292	17.066	9.624	8.804
View		p/n kockás	122	11.039	10.752	7.920
concordance		p/n színű	132	11.469	9.170	7.823
Sample		p/n betűs	101	10.046	11.182	7.687
Filter		p/n kék	186	13.566	7.559	7.596
Frequency		p/n rózsá	124	11.098	8.217	7.518
Node tags		p/n lap	325	17.864	6.779	7.383
Node forms		p/n süt	111	10.496	8.040	7.352
Doc IDs		p/n zászló	109	10.394	7.822	7.264
ConcDesc		p/n sárga	126	11.155	7.328	7.194
?		p/n folt	85	9.190	8.298	7.138
		p/n alma	73	8.524	8.733	7.041
		p/n zöld	148	12.049	6.704	6.973
		p/n elefánt	63	7.925	9.304	6.925
		p/n szín	123	10.978	6.626	6.817
		p/n ceruza	58	7.601	9.017	6.789
		p/n sapka	60	7.723	8.393	6.749
		p/n arcú	58	7.597	8.680	6.749
		p/n jelzés	71	8.371	7.249	6.661



történt volna . Így eshetett meg , hogy **piros lámpánál** az út közepén futottam , s éppen autótolvajok új stratégiát követnek . Így például a **piros lámpánál** hátulról belehajtanak a kiszemelt
Módszerük az , hogy az autópályáról levezető **piros lámpánál** megálló gépkocsi mellé berobogó
Felelőtlenség a gyerekeknek rossz példát mutatni , **piros lámpánál** átszaladni , a szabályokat nem
<p> Mindennapi tapasztalatot elevenít fel a **piros lámpánál** megrekedt autósról szóló jelenet
össze a motoros rendőrökkel . </p><p> Ha még a **piros lámpán** is áthajtott , és a kerékpár felszereltsége
vezető minden alkalommal megjelent , hatalmas **piros** elemes **lámpát** tartott a kezében , és elmondta
arcom . De az biztos , ha a stúdióban ég a **piros lámpa** , akkor mindenki egy kicsit összekapja
és kíséretét szállító konvoj minden egyes **piros lámpánál** megállt Lahorban - a szemtanúk
gépkocsioszlop szabályosan megállt minden egyes **piros lámpánál** , nem kis megdöbbenést okozva
beavatkozásokat eszközölni , ők csupán várakoznak a **piros lámpánál** és adományokat gyűjtenek az autósoktól
" technika . </p><p> Az előző módszerem a **piros lámpánál** álló vezetőt kikényszerítik autójából
például annak , hogy amikor Pesten járt , és a **piros lámpánál** egy autóbusz utolérte , a sofőr
kikapcsolása . (A három futó állapotát egy sor **piros lámpa** is jelzi . Csakhogy Szaloniki felett



<http://hnc.nytud.hu>



Thank you for your attention