

The Hungarian Gigaword Corpus

The Hungarian Language in the Digital Age

Budapest, Jan 18, 2013

Csaba Oravecz

oravecz.csaba@nytud.mta.hu

Research Institute for Linguistics

Hungarian Academy of Sciences

Overview

- Origins
- Motivation
- Objectives
- Preparation
- Outcome

Overview

- Origins
- Motivation
- Objectives
- Preparation
- Outcome

Overview

- Origins
- Motivation
- Objectives
- Preparation
- Outcome

Overview

- Origins
- Motivation
- Objectives
- Preparation
- Outcome

Overview

- Origins
- Motivation
- Objectives
- 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

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

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

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

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

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

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

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

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

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

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.

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.

Objectives

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.

Objectives

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.

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.

Text collection

- clean IPR issues
- extensive metadata (simple webcrawling not sufficient)
- ease of processing → no pdf, no OCR

Preprocessing, normalization

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

Text collection

- clean IPR issues
 - extensive metadata (simple webcrawling not sufficient)
 - ease of processing → no pdf, no OCR

Preprocessing, normalization

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

Text collection

- clean IPR issues
- extensive metadata (simple webcrawling not sufficient)
- ease of processing → no pdf, no OCR

Preprocessing, normalization

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

Text collection

- clean IPR issues
- extensive metadata (simple webcrawling not sufficient)
- ease of processing → no pdf, no OCR

Preprocessing, normalization

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

Text collection

- clean IPR issues
- extensive metadata (simple webcrawling not sufficient)
- ease of processing → no pdf, no OCR

Preprocessing, normalization

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

Text collection

- clean IPR issues
- extensive metadata (simple webcrawling not sufficient)
- ease of processing → no pdf, no OCR

Preprocessing, normalization

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

Text collection

- clean IPR issues
- extensive metadata (simple webcrawling not sufficient)
- ease of processing → no pdf, no OCR

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
- XML annotation compatible with international standards

Corpus query engine

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

Analysis and annotation

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

Corpus query engine

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

Analysis and annotation

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

Corpus query engine

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

Analysis and annotation

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

Corpus query engine

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

Analysis and annotation

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

Corpus query engine

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

Analysis and annotation

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

Corpus query engine

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

Analysis and annotation

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

Corpus query engine

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

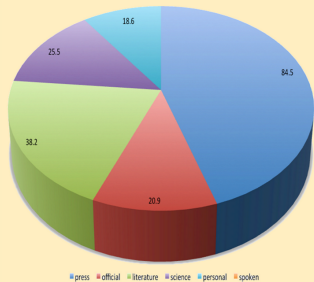
Outcome

HNC: 187 m.

HGC (+HNC): 1091 m.

Outcome

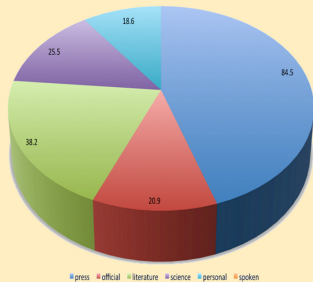
HNC: 187 m.



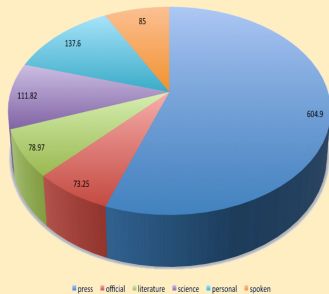
HGC (+HNC): 1091 m.

Outcome

HNC: 187 m.



HGC (+HNC): 1091 m.



"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)

"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)

"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)

"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)

"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)

"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)

"piros ..."

Concordance Word List		Collocation candidates					
?		Page	<input type="text" value="1"/>	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	szinű	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	rózsa	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	falt	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	szin	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

"piros lámpa"

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 **piros lámpánál** megrekedt autósról szóló jelenet **piros lámpánál** össze a motoros rendőrökkel . **piros lámpánál** Ha még a **piros lámpánál** is áthajtott , és a kerékpár felszereltsége vezető minden alkalommal megjelent , hatalmas **piros lámpánál** elemes **piros lámpánál** tartott a kezében , és elmondta arcom . De az biztos , ha a stúdióban ég a **piros lámpánál** , akkor mindenki egy kicsit összekapja és kíséretét szállító konvoj minden egyes **piros lámpánál** megállt Lahorban - a szemtanúk **piros lámpánál** , nem kis megdöbbenést okozva **piros lámpánál** beavatkozásokat eszközölni , ők csupán várakoznak a **piros lámpánál** és adományokat gyűjtenek az autósoktól " technika . **piros lámpánál** Az előző módszeremél a **piros lámpánál** álló vezetőt kikényszerítik autójából példádul annak , hogy amikor Pesten járt , és a **piros lámpánál** egy autóbusz utolérte , a sofőr kikapcsolása . (A három futó állapotát egy sor **piros lámpánál** is jelzi . Csakhogy Szaloniki felett

Web interface

<http://mnsz.nytud.hu>

The end

Thank you for your attention