

Furthering Natural Language Processing in Bulgaria

Svetla Koeva

Institute for Bulgarian, Bulgaria
svetlal@dcl.bas.bg

META-FORUM
Budapest, Hungary, 2011-06-27-28



Co-funded by the 7th Framework Programme of the European Commission through the contract T4ME, grant agreement no.: 249119.



Co-funded by the ICT PSP Programme of the European Commission through the contract CESAR, grant agreement no.: 271022.



General facts

- Republic of Bulgaria
- Area - 110, 993. 6 km²
- Population - 7 351 633
- Bulgarian -
9 million
native speakers



General facts

- Republic of Bulgaria
- Area - 110, 993. 6 km²
- Population - 7 351 633
- Bulgarian -
9 million
native speakers



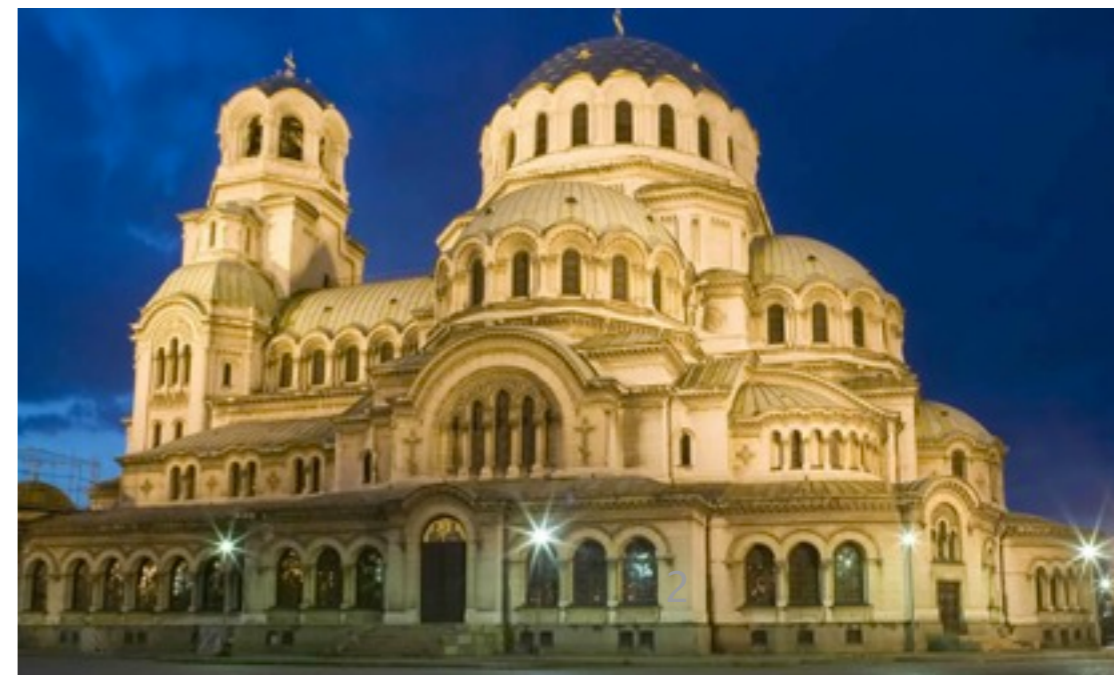
General facts

- Republic of Bulgaria
- Area - 110, 993. 6 km²
- Population - 7 351 633
- Bulgarian -
9 million
native speakers



General facts

- Republic of Bulgaria
- Area - 110, 993. 6 km²
- Population - 7 351 633
- Bulgarian -
9 million
native speakers



General facts

- Republic of Bulgaria
- Area - 110, 993. 6 km²
- Population - 7 351 633
- Bulgarian -
9 million
native speakers



General facts

- Republic of Bulgaria
- Area - 110, 993. 6 km²
- Population - 7 351 633
- Bulgarian -
9 million
native speakers



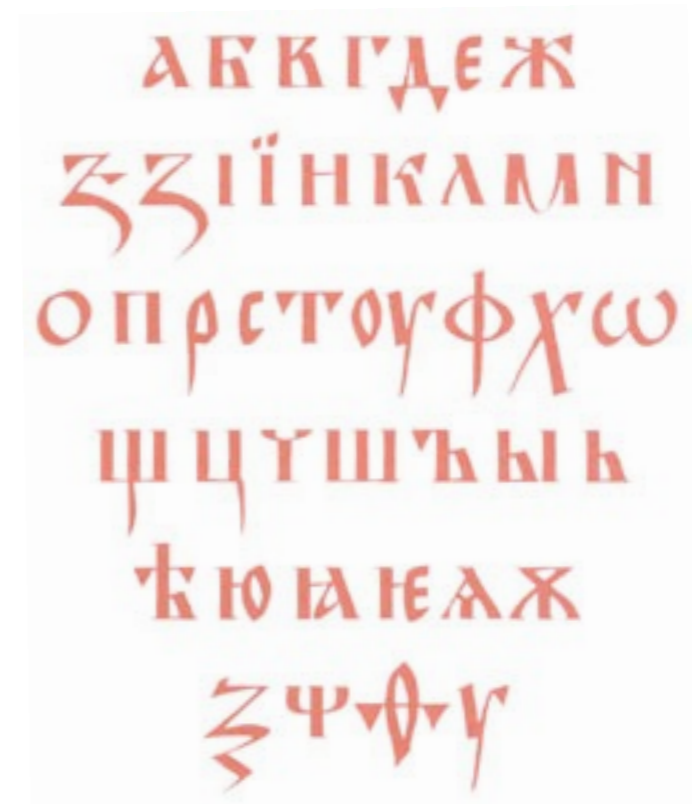
General facts

- Republic of Bulgaria
- Area - 110, 993. 6 km²
- Population - 7 351 633
- Bulgarian -
9 million
native speakers



General facts

- The official alphabet is Cyrillic.
- Официалната азбука е кирилица.
- Cyrillic became the third official alphabet of the European Union, following the Latin and Greek alphabets.



Furthering NLP in Bulgaria

Research



Furthering NLP in Bulgaria

Research



Furthering NLP in Bulgaria

Research



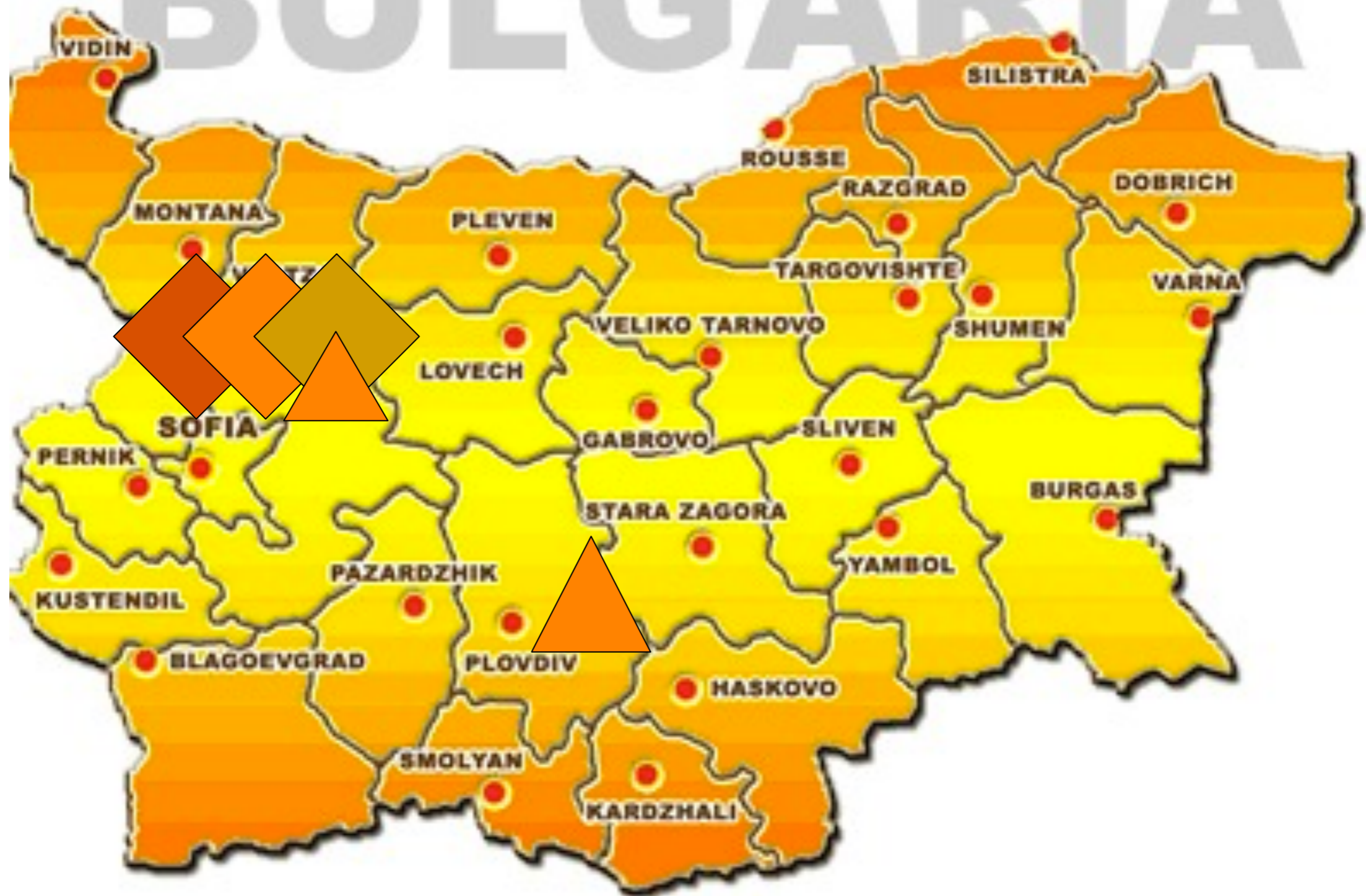
Furthering NLP in Bulgaria Research



Furthering NLP in Bulgaria Research



Furthering NLP in Bulgaria Research



Furthering NLP in Bulgaria Research



Furthering NLP in Bulgaria Research



Furthering NLP in Bulgaria

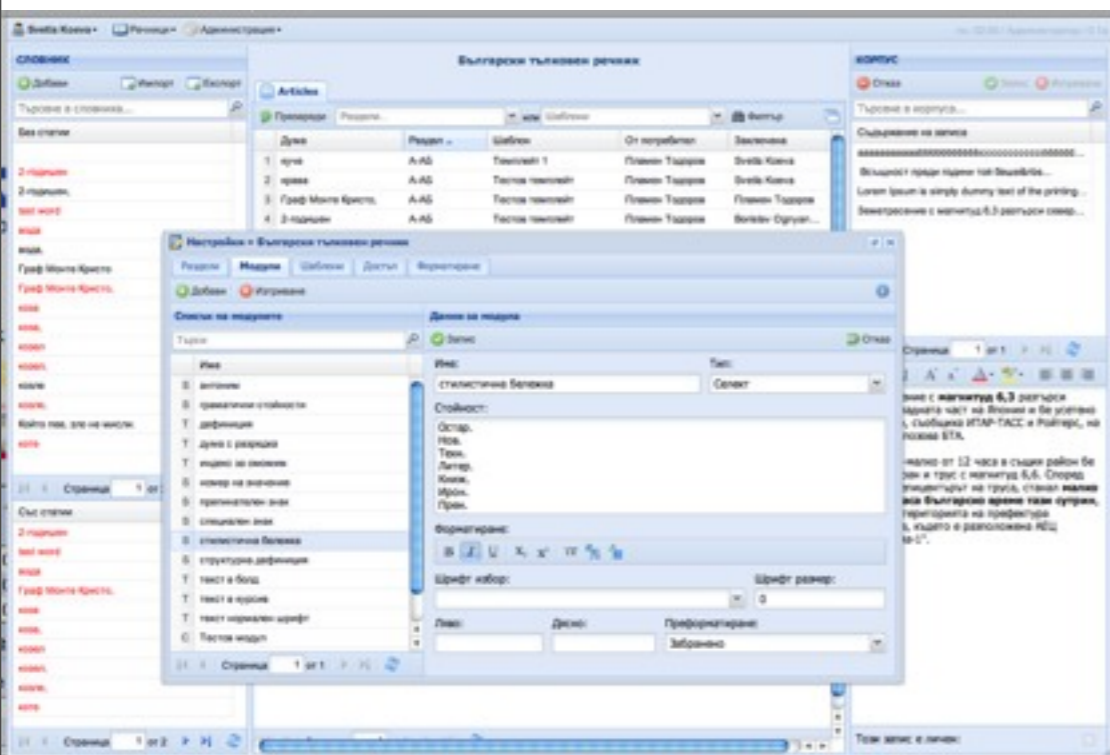
BLARK

- Over the past decade a number of important language resources and tools have been developed.

Furthering NLP in Bulgaria

BLARK

- Over the past decade a number of important language resources and tools have been developed.



Furthering NLP in Bulgaria

BLARK

- Over the past decade a number of important language resources and tools have been developed.



DEPARTMENT OF COMPUTATIONAL LINGUISTICS

The Bulgarian Semantically Annotated Corpus

The Bulgarian Semantically Annotated Corpus is part of the Bulgarian Brown Corpus. It consists of 95119 lexical units, annotated with the most appropriate synonymous set from the Bulgarian wordnet.

Word: Lemma: лице Search

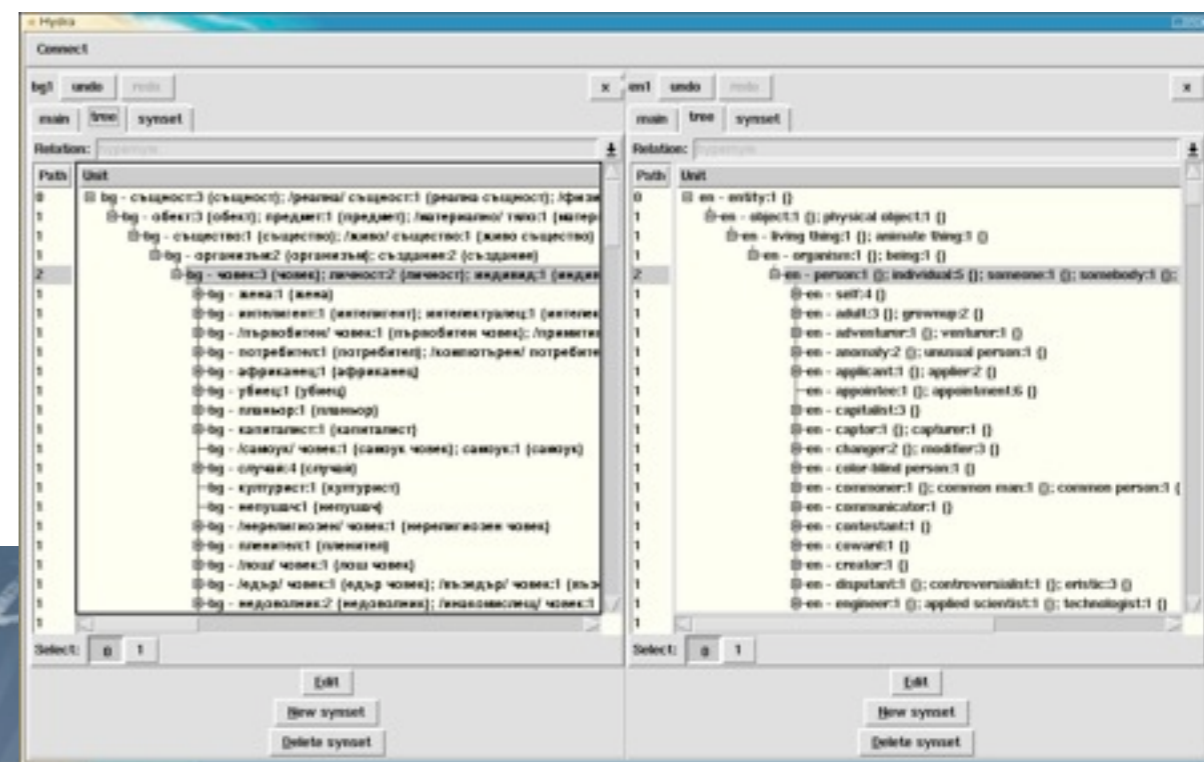
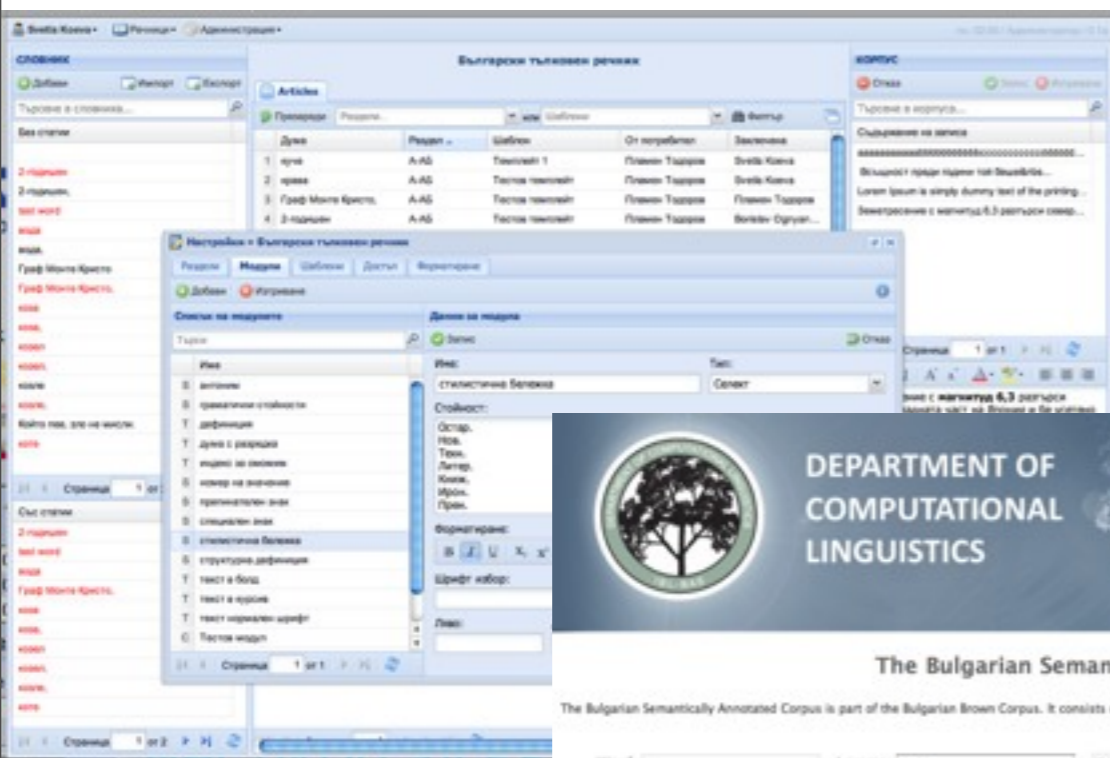
Results:

- Occurrences: 43, sense definition: предната по-плоска част на главата на човек, ограничена от челото и челюстите, заедно с разположените върху нея органи - уста, очи, нос
Суботният провал изгласи широки усмивки на удовлетворение по **лицата** на хората, които през последните две седмици открито се обвяха срещу регистрацията на партията по принцип и срещу начина, по който организаторите подготвят несъстоялото се събитие.
- Occurrences: 16, sense definition: представител на едноименния и единствен вид на същото семейство висши бозайници (Hominidae), различаващ се от останалите животни по силно развития мозък, съзнание, абстрактно мислене, членоразделна реч; движи се с изправено тяло, произвежда оръдия на труда и други артефакти
- Бъдете спокоен, господни полковник, в мое **лице** ще имате добър водач.
- Occurrences: 4, sense definition: метонимично название за отделен човек

Furthering NLP in Bulgaria

BLARK

- Over the past decade a number of important language resources and tools have been developed.



 DEPARTMENT OF COMPUTATIONAL LINGUISTICS

The Bulgarian Semantically Annotated Corpus

The Bulgarian Semantically Annotated Corpus is part of the Bulgarian Brown Corpus. It consists of 95119 lexical units, annotated with the most appropriate synonymous set from the Bulgarian wordnet.

Word: Lemma: лице Search

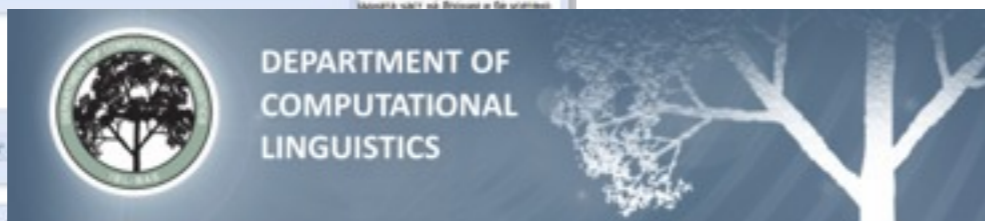
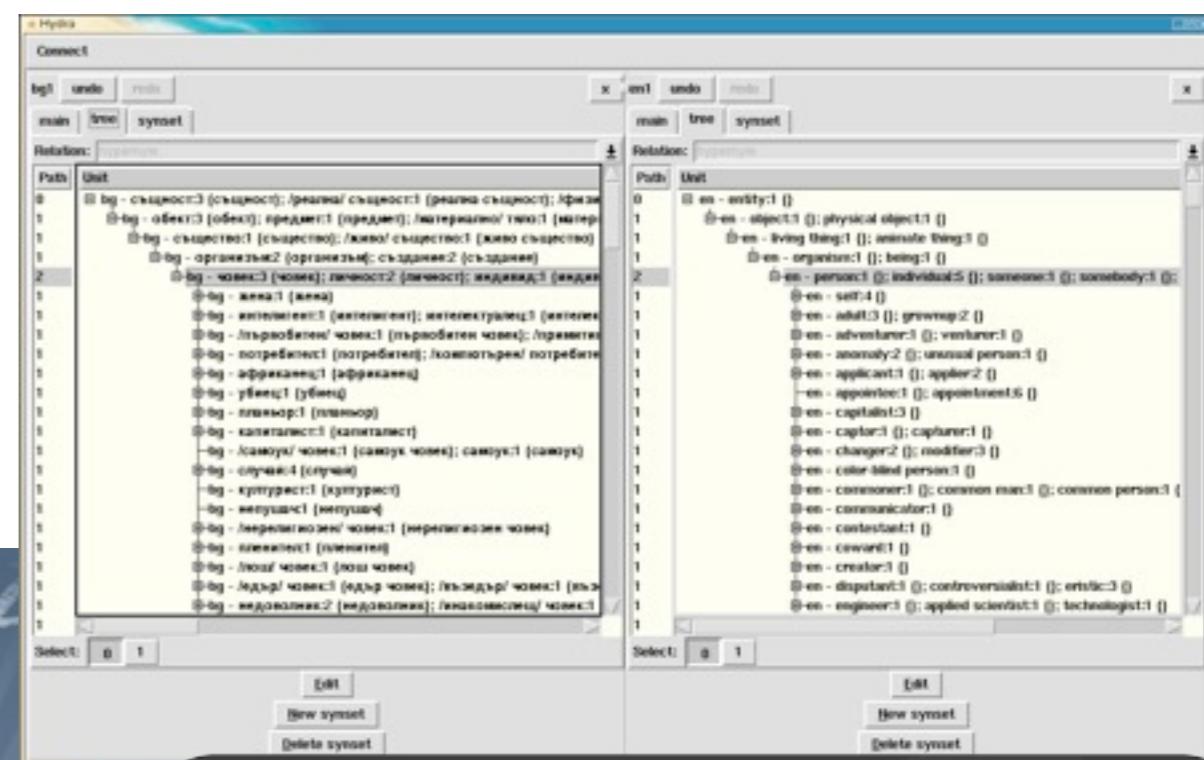
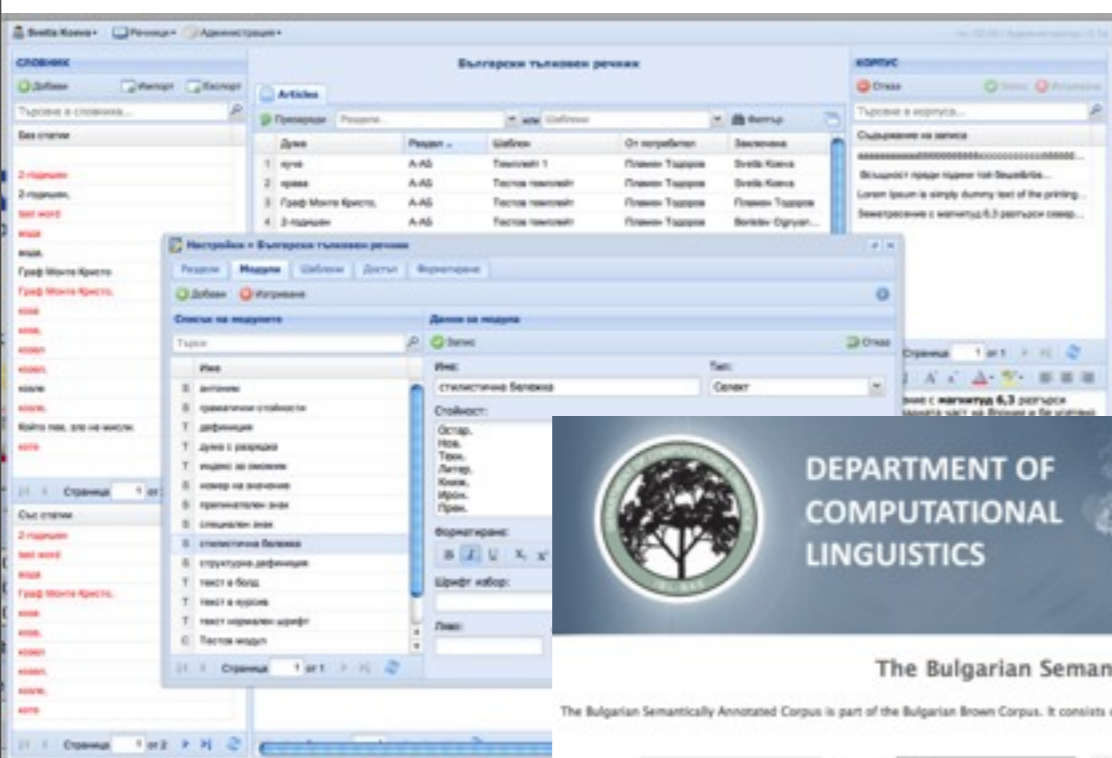
Results:

- Occurrences: 43, sense definition: предната по-голяма част на главата на човек, ограничена от челото и челюстите, заедно с разположените върху нея органи - уста, очи, нос
Съботният провал изгласи широки усмивки на удовлетворение по **лицата** на хората, които през последните две седмици открито се обвяха срещу регистрацията на партията по принцип и срещу начина, по който организаторите подготвят несъстоялото се събитие.
- Occurrences: 16, sense definition: представител на едноименния и единствен вид на същото семейство висши бозайници (Hominidae), различаващ се от останалите животни по силно развития мозък, съзнание, абстрактно мислене, членоразделна реч; движи се с изправено тяло, произвежда оръдия на труда и други артефакти
- Бъдете спокоен, господни полковник, в мое **лице** ще имате добър водач.
- Occurrences: 4, sense definition: метонимично название за отделен човек

Furthering NLP in Bulgaria

BLARK

- Over the past decade a number of important language resources and tools have been developed.



The Bulgarian Semantically Annotated Corpus

The Bulgarian Semantically Annotated Corpus is part of the Bulgarian Brown Corpus. It consists of 95119 lexical units, annotated with the most appropriate synonymous set from the Bulgarian

Word: Lemma: лице Search

Results:

- Occurrences: 43, sense definition: предната по-голяма част на главата на човек, ограничена от челото и челюстите, заедно с разположените върху нея органи - уста, очи, нос
Съботният провал изгласи широки усмивки на удовлетворение по лицата на хората, които през последните две седмици открито се обвяха срещу регистрацията на партията по принцип и срещу начина, по който организаторите подготвят несъстоятелното се събитие.
- Occurrences: 16, sense definition: представител на едноименния и единствен вид на същото семейство висши базилени (Hominidae) различаващ се от останалите животни по силно развития мозък, съзнание, абстрактно мислене, членоразделна реч движи се с изправено тяло, произвежда оръдия на труда и други артефакти
- Бъдете спокоен, господин полковник, в мое лице ще имате добър водач.
- Occurrences: 4, sense definition: метонимично название за отделен човек

DCL - Department of Computational Linguistics (IBL -BAS)

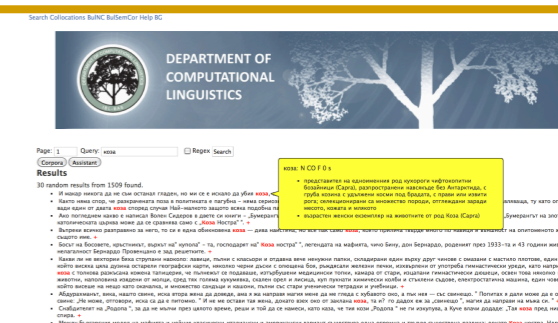
MacEst

Bulgarian Spell Checker for Mac OS

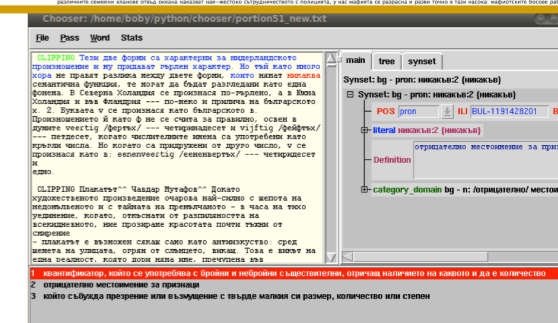
Furthering NLP in Bulgaria

BLARK

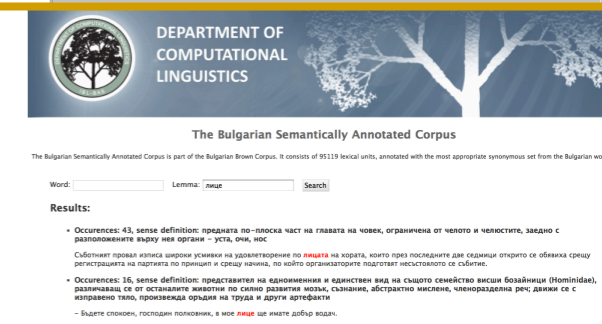
Bulgarian National Corpus - app.
500M words



Bulgarian POS-annotated
Corpus



Bulgarian Sense-annotated Corpus



Dependency part of BulTreeBank



Furthering NLP in Bulgaria

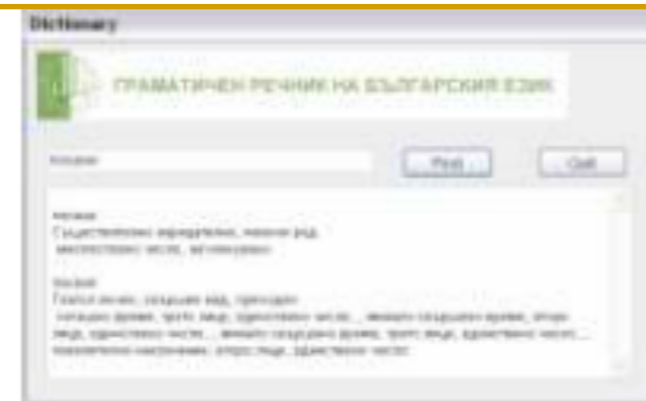
BLARK

- ❑ SEE-ERA.net Administrative and Literally Corpus
- ❑ Bilingual collection of cultural texts in Greek and Bulgarian
- ❑ Bulgarian-Polish-Lithuanian Corpus
- ❑ Bulgarian-English-X language parallel corpus - app. 100M words for Bulgarian
- ❑ ...

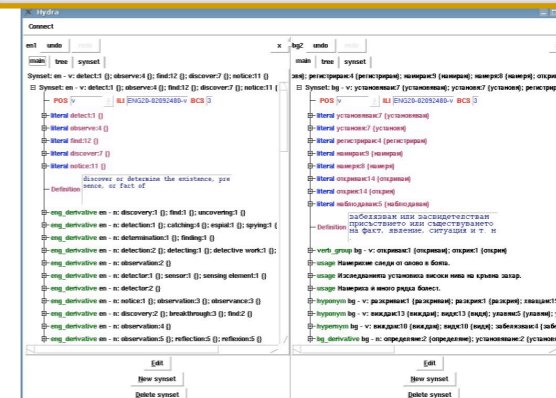
Furthering NLP in Bulgaria

BLARK

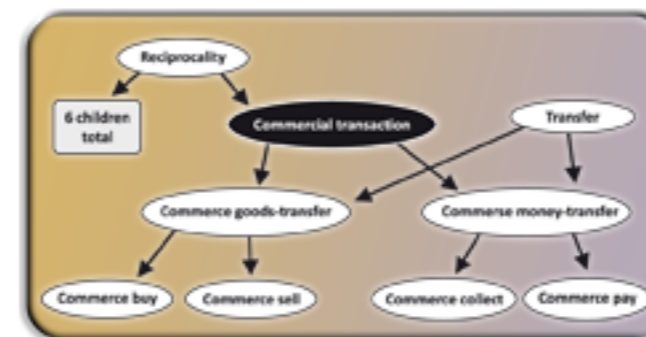
Several large inflectional dictionaries



Bulgarian WordNet



Bulgarian FrameNet



Furthering NLP in Bulgaria Companies



- ❑ Development of tools and solutions based on semantic technologies
- ❑ Ontology design
- ❑ Data integration, management and publishing

Furthering NLP in Bulgaria Companies



 <p>cubist Your Business Intelligence</p>	 <p>LARKC LARGE KNOWLEDGE COLLIDER</p>	 <p>NoTube</p>
 <p>insemtives Incentives for Semantics</p>	 <p>MOLTO</p>	 <p>RENDER Reflecting Knowledge Diversity</p>
 <p>khresmoi MEDICAL INFORMATION ANALYSIS & RETRIEVAL</p>		 <p>SOA Semantics Web2.0 Web Services Context ALL</p>

Furthering NLP in Bulgaria

Companies



- ❑ Web applications (dynamic web content)
- ❑ Content Management Systems (CMS)
- ❑ Tools for web site content management
- ❑ Multilingual tools and services for natural language processing

Furthering NLP in Bulgaria Companies



Library Home



My items



Add new Book



Shared Area



Public domain



My profile



Logout

My Items

Search

Author: ARTHUR CONAN DOYLE ✕

All entries

By title

By author

By year ↑

Recent entries

By year ↓

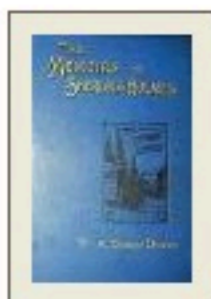


The adventures of Sherlock Holmes

Author Arthur Conan Doyle

Year of publication 1892

File  pg1661.txt



The Memoirs of Sherlock Holmes

Author Arthur Conan Doyle

Year of publication 1894

File  pg834.txt



Furthering NLP in Bulgaria

Companies

SkyCode

- ❑ WebTrance - a translation software package from English, French, German, Spanish, Italian and Turkish to Bulgarian and vice versa.
- ❑ SkyCode is one of the partners of *iTranslate4*.

Approximate status

Technology	Median
Tokenization, Morphology	4
Parsing	3
Information Retrieval	2
Speech Synthesis	2
Text semantics	2
Information extraction	2
Summarization, QA	2
Machine translation	2
Language generation	1

Resources	Median
Reference Corpora	4
Thesauri, WordNets	4
Lexicons, Terminologies	3
Semantic corpora	3
Parallel Corpora, TM	2
Syntax-Corpora	2
Discourse-Corpora	1
Multimedia/multimodal data	1

Approximate status

Technology	Median
Tokenization, Morphology	4
Parsing	3
Information Retrieval	2
Speech Synthesis	2
Text semantics	2
Information extraction	2
Summarization, QA	2
Machine translation	2
Language generation	1

Resources	Median
Reference Corpora	4
Thesauri, WordNets	4
Lexicons, Terminologies	3
Semantic corpora	3
Parallel Corpora, TM	2
Syntax-Corpora	2
Discourse-Corpora	1
Multimedia/multimodal data	1

Approximate status

Technology	Median
Tokenization, Morphology	4
Parsing	3
Information Retrieval	2
Speech Synthesis	2
Text semantics	2
Information extraction	2
Summarization, QA	2
Machine translation	2
Language generation	1

Resources	Median
Reference Corpora	4
Thesauri, WordNets	4
Lexicons, Terminologies	3
Semantic corpora	3
Parallel Corpora, TM	2
Syntax-Corpora	2
Discourse-Corpora	1
Multimedia/multimodal data	1

Approximate status

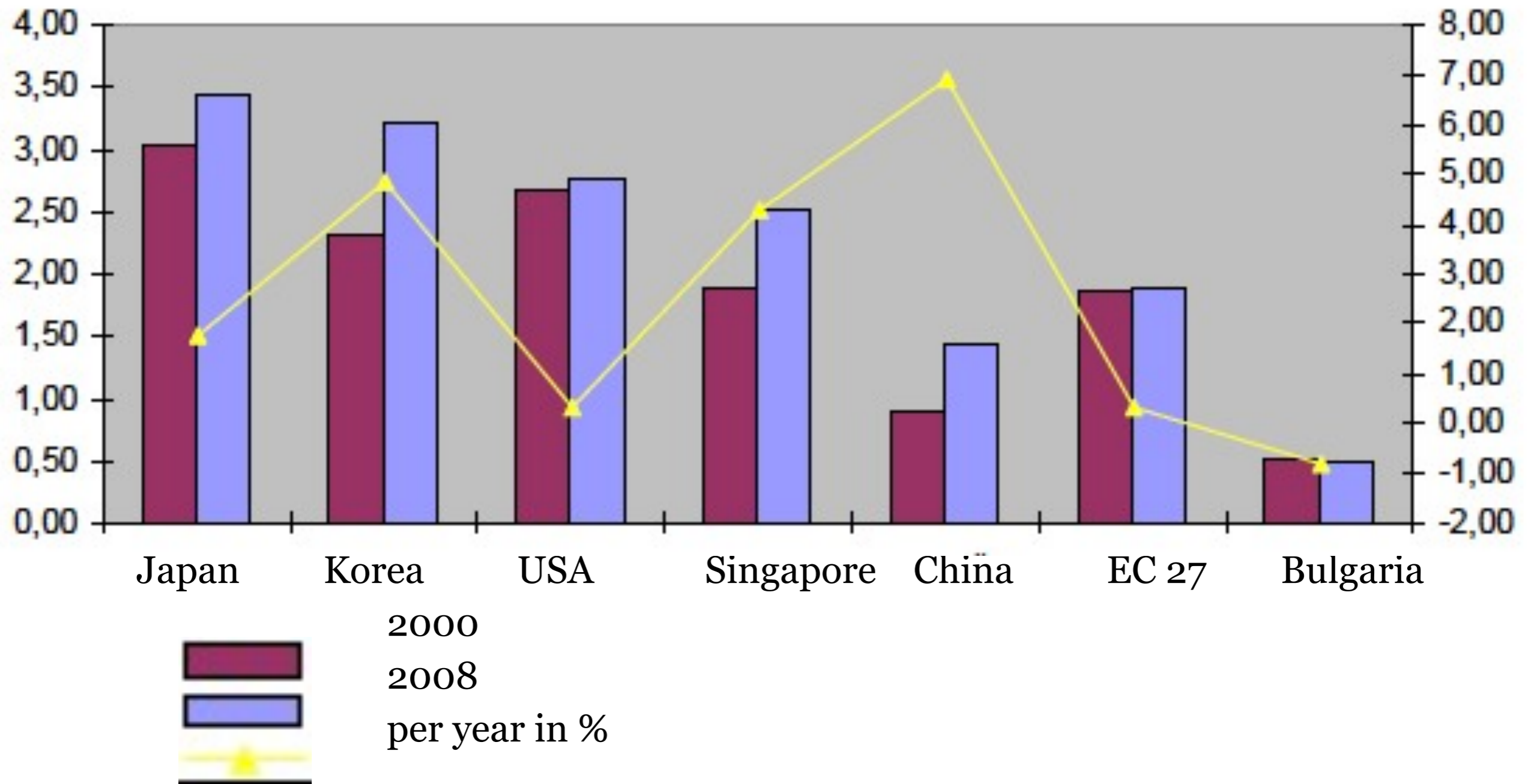
Technology	Median	Resources	Median
Tokenization, Morphology	4	Reference Corpora	4
Parsing	3	Thesauri, WordNets	4
Information Retrieval	2	Lexicons, Terminologies	3
Speech Synthesis	2	Semantic corpora	3
Text semantics	2	Parallel Corpora, TM	2
Information extraction	2	Syntax-Corpora	2
Summarization, QA	2	Discourse-Corpora	1
Machine translation	2	Multimedia/multimodal data	1
Language generation	1		

Approximate status

	Quantity	Availability	Quality	Coverage	Maturity	Sustainability	Adaptability
Technology	2	2	2.5	2.5	2	2	2.5
Resources	2	2.5	3	3.5	2.5	2.5	2.5
Total	2	2	2.5	3	2	2	2.5

Furthering NLP in Bulgaria

State contribution to R&D



Furthering NLP in Bulgaria

State contribution to R&D

- ❑ Strategic research agenda:
- ❑ Cultural-historical heritage - language being a central part of it
- ❑ ICT as a horizontal instrument

- ❑ **META, Multilingual Europe Technology Alliance**
- ❑ Institute for Bulgarian, BAS, a member of META^{NET}
- ❑ Institute for Literature, BAS
- ❑ Institute of Information and Communication Technologies, BAS
- ❑ Sofia University St. Kliment Ohridski
- ❑ University of Plovdiv
- ❑ Ontotext, Bulgaria
- ❑ Musala Soft, Bulgaria
- ❑ Tetracom Interactive Solutions, Bulgaria
- ❑ TransGlobe International Ltd., Bulgaria



Conclusions

- Several factors are mutually related for the success:
 - clear formulation of target goals and strategies for their accomplishment
 - stable financing
 - effective management of the resources
 - beneficial relations between education - research - business - end users
 - networking
- META-NET as a concerted, substantial, continent-wide effort in language technology research and engineering is relevant for all of these factors.



□ Thank you very much for your attention.