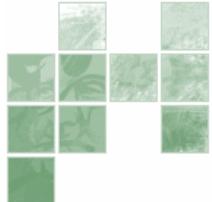


# Automatic Lexicon Acquisition for a Medical Cross-Language Information Retrieval System

Kornél Markó, Stefan Schulz, Udo Hahn

Freiburg University Hospital, Medical Informatics Department, Germany  
Jena University, Language & Information Engineering (JULIE), Germany





# Multilingual Text retrieval



„Korrelation von  
Hypertonie und  
Läsion der  
Weißen Substanz“

Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=pubmed: Home Bookmarks mozilla.org Latest Builds Search Print

□ 1: Stroke. 2004 May;35(5):1057-60. Epub 2004 Apr 1. Related Articles, Links

Comment in:

- [Stroke. 2004 May;35\(5\):1061-2.](#)

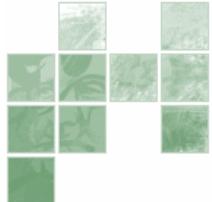
Full text article at  
[stroke.ahajournals.org](#)

**Interaction between hypertension, apoE, and cerebral white matter lesions.**

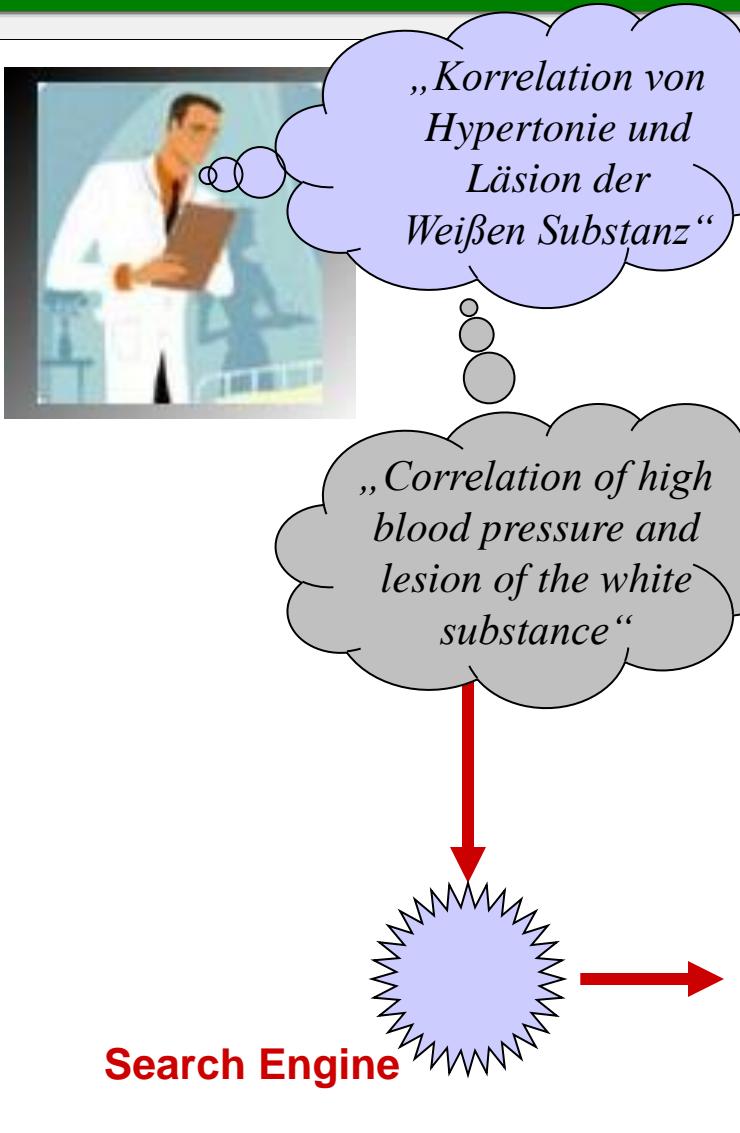
de Leeuw FE, Richard F, de Groot JC, van Duijn CM, Hofman A, Van Gijn J, Breteler MM.

Department of Epidemiology and Biostatistics, Erasmus Medical Center, Rotterdam, The Netherlands.

BACKGROUND AND PURPOSE: Cerebral white matter lesions (WMLs) are frequently found on magnetic resonance imaging scans in both cognitively intact and demented elderly persons. Vascular risk factors, especially hypertension, are related to their presence. However, not every person with vascular risk factors has WMLs, which suggests interaction with other determinants, eg, genetic factors. The epsilon4 allele



# Multilingual Text retrieval



„Korrelation von Hypertonie und Läsion der Weißen Substanz“

„Correlation of high blood pressure and lesion of the white substance“

Search Engine

Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=pubmed: mozila.org Latest Builds Search Print

Home Bookmarks mozilla.org Latest Builds

Entrez PubMed Overview Help | FAQ Tutorial New/Noteworthy E-Utilities

PubMed Services Journals Database WSH Database Single Citation Matcher Batch Citation Matcher Clinical Queries LinkOut My NCBI (Cubby)

Related Resources Order Documents NLM Catalog NLM Gateway

□ 1: Stroke. 2004 May;35(5):1057-60. Epub 2004 Apr 1. Related Articles, Links

Comment in:

- [Stroke. 2004 May;35\(5\):1061-2.](#)

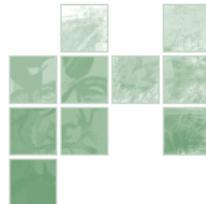
Full text article at stroke.ahajournals.org

**Interaction between hypertension, apoE, and cerebral white matter lesions.**

de Leeuw FE, Richard F, de Groot JC, van Duijn CM, Hofman A, Van Gijn J, Breteler MM.

Department of Epidemiology and Biostatistics, Erasmus Medical Center, Rotterdam, The Netherlands.

BACKGROUND AND PURPOSE: Cerebral white matter lesions (WMLs) are frequently found on magnetic resonance imaging scans in both cognitively intact and demented elderly persons. Vascular risk factors, especially hypertension, are related to their presence. However, not every person with vascular risk factors has WMLs, which suggests interaction with other determinants, eg, genetic factors. The epsilon4 allele



# Multilingual Text retrieval

„Korrelation von Hypertonie und Läsion der Weißen Substanz“

„Correlation of high blood pressure and lesion of the white substance“

Search Engine

Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop

Home Bookmarks mozilla.org Latest Builds

1: Stroke. 2004 May;35(5):1057-60. Epub 2004 Apr 1.

Comment in:

- Stroke. 2004 May;35(5):1061-2.

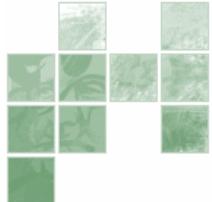
Full text article of stroke.ahajournals.org

Interaction between **hypertension**, apoE, and cerebral **white matter lesions**.

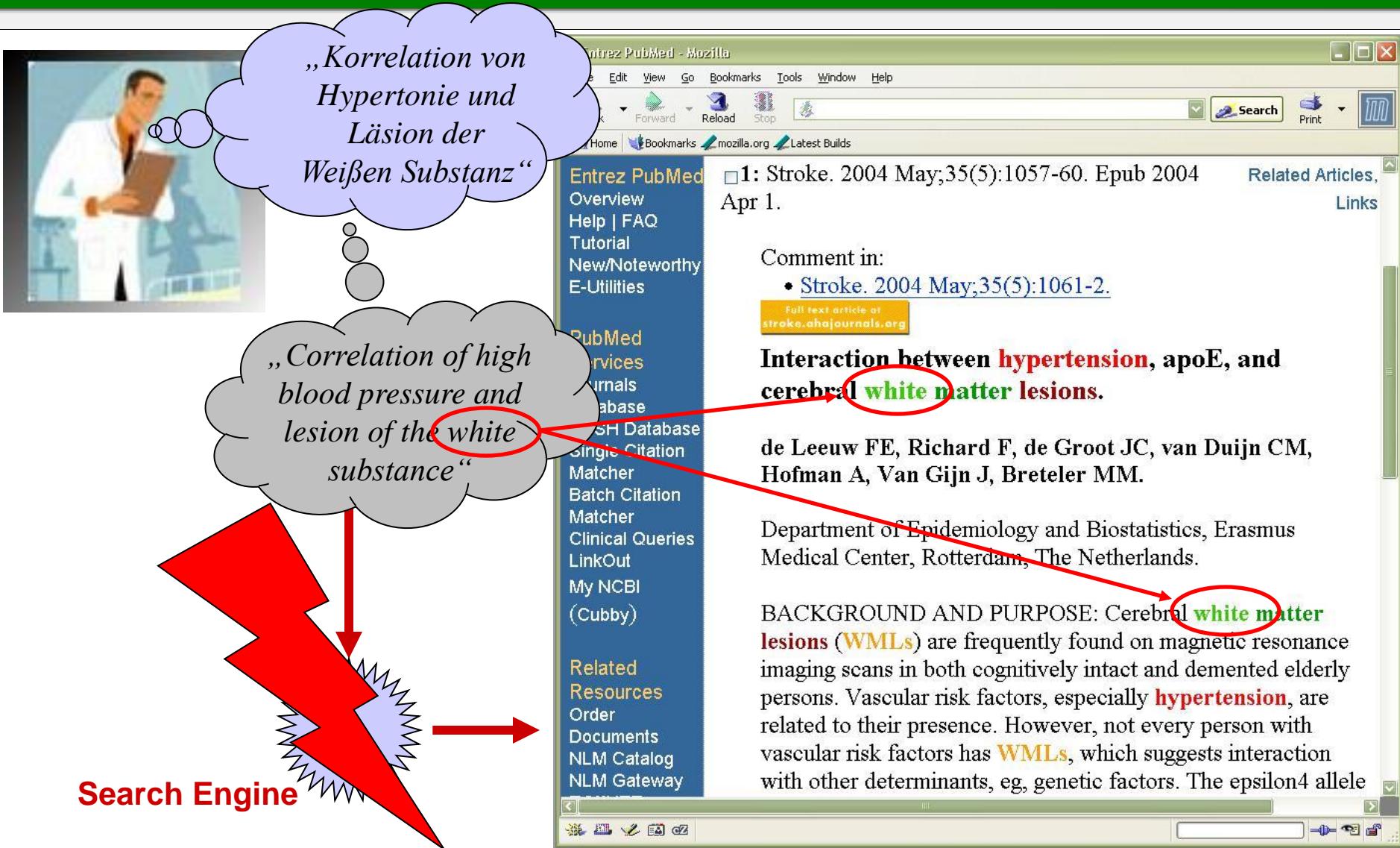
de Leeuw FE, Richard F, de Groot JC, van Duijn CM, Hofman A, Van Gijn J, Breteler MM.

Department of Epidemiology and Biostatistics, Erasmus Medical Center, Rotterdam, The Netherlands.

BACKGROUND AND PURPOSE: Cerebral **white matter lesions** (**WMLs**) are frequently found on magnetic resonance imaging scans in both cognitively intact and demented elderly persons. Vascular risk factors, especially **hypertension**, are related to their presence. However, not every person with vascular risk factors has **WMLs**, which suggests interaction with other determinants, eg, genetic factors. The epsilon4 allele



# Multilingual Textretrieval



The diagram illustrates the process of multilingual text retrieval. It starts with a medical concept represented by two thought bubbles: one in German ("Korrelation von Hypertonie und Läsion der Weißen Substanz") and one in English ("Correlation of high blood pressure and lesion of the white substance"). A red lightning bolt symbol labeled "Search Engine" points from these concepts to a screenshot of a PubMed search results page. The search results page shows a study titled "Interaction between hypertension, apoE, and cerebral white matter lesions." The word "white matter" is circled in red. The study is attributed to de Leeuw FE, Richard F, de Groot JC, van Duijn CM, Hofman A, Van Gijn J, Breteler MM. The background and purpose section of the study text also mentions "white matter lesions" (WMLs), which is also circled in red.

„Korrelation von Hypertonie und Läsion der Weißen Substanz“

„Correlation of high blood pressure and lesion of the white substance“

Search Engine

Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop

Home Bookmarks mozilla.org Latest Builds

Comment in:

- Stroke. 2004 May;35(5):1061-2.

Full text article of stroke.ahajournals.org

Interaction between **hypertension**, apoE, and cerebral **white matter** lesions.

de Leeuw FE, Richard F, de Groot JC, van Duijn CM, Hofman A, Van Gijn J, Breteler MM.

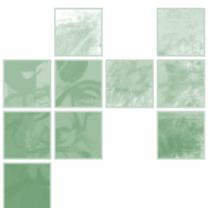
Department of Epidemiology and Biostatistics, Erasmus Medical Center, Rotterdam, The Netherlands.

BACKGROUND AND PURPOSE: Cerebral **white matter** lesions (**WMLs**) are frequently found on magnetic resonance imaging scans in both cognitively intact and demented elderly persons. Vascular risk factors, especially **hypertension**, are related to their presence. However, not every person with vascular risk factors has **WMLs**, which suggests interaction with other determinants, eg, genetic factors. The epsilon4 allele



# Linguistic Phenomena

- Morphological processes:
  - Inflection: *leukocyte* <> *leukozytes*,  
*appendix* <> *appendices*
  - Derivation: *leukocyte* <> *leukocytic*
  - Composition: *leukemia*, *parasympathectomy*,  
*Magen/schleim/haut/entzündung*
- Synonymy:
  - *ascorbic acid* <> *vitamin C*, *hemorrhage* <> *bleeding*
- Spelling variants:
  - *oesophagus* <> *esophagus*,
  - *Karzinom* <> *Carcinom* <> *Carzinom (carcinoma)*



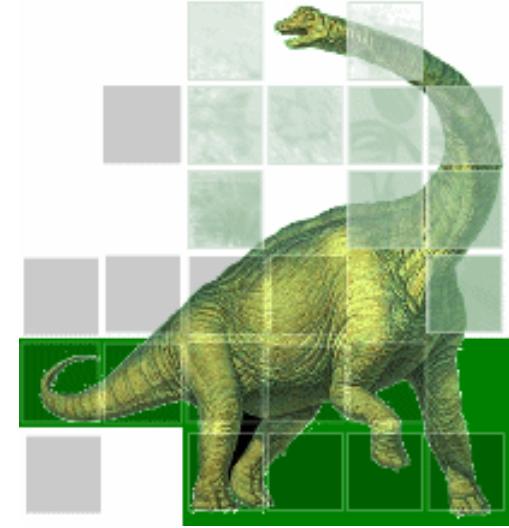
# Subword Approach

- Subwords are atomic, conceptual or linguistic units:
  - Stems: *stomach, gastr, diophys*
  - Prefixes: *anti-, bi-, hyper-*
  - Suffixes: *-ary, -ion, -itis*
  - Infixes: *-o-, -s-*
- Equivalence classes contain synonymous subwords and their translations in a thesaurus:

**#female = { woman, women, female, frau, weib, mulher }**



# Morphosaurus

- Subword-Lexicon:
    - Organizes subwords in several languages (English, German, Portuguese)
  - Subword-Thesaurus:
    - Groups synonymous subwords (within and between languages)
  - Subword-Segmenter:
    - Extraction of Subwords and Assignment of *Equivalence Classes*
- 
- Morphosaurus**
- Morphosaurus-  
Identifier (MID)**



# Example

High TSH values suggest the diagnosis of primary hypothyroidism ...

Erhöhte TSH-Werte erlauben die Diagnose einer primären Hypothyreose ...

## Original

## Interlingua

#up tsh #value #suggest  
#diagnost #primar #small  
#thyre

#up tsh #value #permit  
#diagnost #primar #small  
#thyre

## Orthographic Normalization

## Orthografic Rules

high tsh values suggest the diagnosis of primary hypothyroidism ...

erhoehte tsh-werte erlauben die diagnose einer primaeren hypothyreose ...

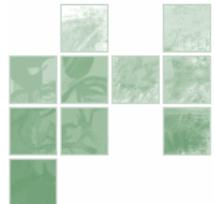
## Segmenter Subword Lexicon

## Semantic Normalization

## Subword Thesaurus

high tsh value s suggest the diagnos is of primar y hypo thyroid ism

er hoeh te tsh wert e erlaub en die diagnos e einer primaer en hypo thyre ose



# Example

High TSH values suggest the diagnosis of primary hypothyroidism ...

Erhöhte TSH-Werte erlauben die Diagnose einer primären Hypothyreose ...

**Original**

**Interlingua**

#up tsh #value #suggest  
#diagnost #primar #small  
#thyre

#up tsh #value #permit  
#diagnost #primar #small  
#thyre

**Orthographic Normalization**

**Orthografic Rules**

high tsh values suggest the diagnosis of primary hypothyroidism ...

erhoehte tsh-werte erlauben die diagnose einer primaeren hypothyreose ...

**Segmenter Subword Lexicon**

**Semantic Normalization**

**Subword Thesaurus**

high tsh value s suggest the diagnos is of primar y hypo thyroid ism

er hoeh te tsh wert e erlaub en die diagnos e einer primaer en hypo thyre ose



# Morphosaurus Search

Entrez PubMed - Mozilla Firefox

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=pubmed Search Print

Home Bookmarks mozilla.org Latest Builds

Entrez PubMed Overview Help | FAQ Tutorial New/Noteworthy E-Utilities

PubMed Services Journals Database MeSH Database Single Citation Matcher Batch Citation Matcher Clinical Queries LinkOut My NCBI (Cubby)

Related Resources Order Documents NLM Catalog NLM Gateway

□1: Stroke. 2004 May;35(5):1057-60. Epub 2004 Apr 1.

Comment in:  
• [Stroke. 2004 May;35\(5\):1061-2.](#)

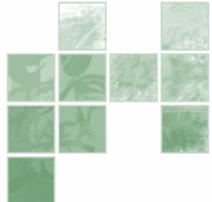
[Full text article at stroke.ahajournals.org](#)

**Interaction between hypertension, apoE, and cerebral white matter lesions.**

de Leeuw FE, Richard F, de Groot JC, van Duijn CM, Hofman A, Van Gijn J, Breteler MM.

Department of Epidemiology and Biostatistics, Erasmus Medical Center, Rotterdam, The Netherlands.

BACKGROUND AND PURPOSE: Cerebral white matter lesions (WMLs) are frequently found on magnetic resonance imaging scans in both cognitively intact and demented elderly persons. Vascular risk factors, especially hypertension, are related to their presence. However, not every person with vascular risk factors has WMLs, which suggests interaction with other determinants, eg, genetic factors. The epsilon4 allele



# Morphosaurus Search

Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop http://www.ncbi.nlm.nih.gov/entrez/query.fcgi?cmd=Retrieve&db=pubmed Search Print

Home Bookmarks mozilla.org Latest Builds

Entrez PubMed Overview Help | FAQ Tutorial New/Noteworthy E-Utilities

PubMed Services Journals Database MeSH Database Single Citation Matcher Batch Citation Matcher Clinical Queries LinkOut My NCBI (Cubby) Related Resources Order Documents NLM Catalog NLM Gateway

Comment in:  
• [Stroke. 2004 May;35\(5\):1061-2.](#)

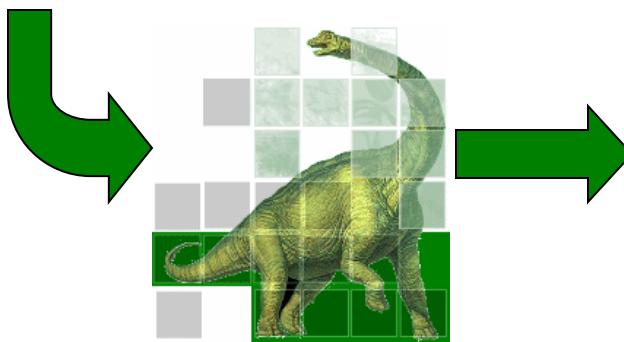
[Full text article at stroke.ahajournals.org](#)

**Interaction between hypertension, apoE, and cerebral white matter lesions.**

de Leeuw FE, Richard F, de Groot JC, van Duijn CM, Hofman A, Van Gijn J, Breteler MM.

Department of Epidemiology and Biostatistics, Erasmus Medical Center, Rotterdam, The Netherlands.

BACKGROUND AND PURPOSE: Cerebral white matter lesions (WMLs) are frequently found on magnetic resonance imaging scans in both cognitively intact and demented elderly persons. Vascular risk factors, especially hypertension, are related to their presence. However, not every person with vascular risk factors has WMLs, which suggests interaction with other determinants, eg, genetic factors. The epsilon4 allele



Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop mozilla.org Latest Builds

Entrez PubMed Overview Help | FAQ Tutorial New/Noteworthy E-Utilities

PubMed Services Journals Database MeSH Database Single Citation Matcher Batch Citation Matcher Clinical Queries LinkOut My NCBI (Cubby) Related Resources Order Documents NLM Catalog NLM Gateway

Comment in:  
• [Stroke. 2004 May;35\(5\):1061-2.](#)

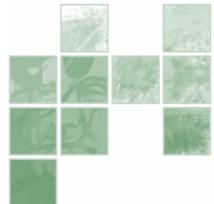
[Full text article at stroke.ahajournals.org](#)

**#interact #hyper #tens , apoE , #cerebr #whit #matter #lesion .**

**de leeuw fe , richard f , de groot jc , van duijn cm , hofman a , van gijn j , breteler mm .**

**#department #epidem #logic #bio #statist , erasmus #medic #centr , rotterdam , #dutch .**

**#back #ground #purpos : #cerebr #whit #matter #lesion ( wmls ) #frequent #find #magnet #resonanc #imag #scan #both #cognit #intact #dement #gero #human . #vascul #risk #factor , #special #hyper #tens , #relat #presenc . #not #total #human #vascul #risk #factor wmls ,# suggest #interact #other #determin , eg ,**



# Morphosaurus Search



„Korrelation von  
Hypertonie und  
Läsion der  
Weißen Substanz“

Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Forward Reload Stop

Home Bookmarks mozilla.org Latest Builds

Entrez PubMed

Overview Help | FAQ Tutorial New/Noteworthy E-Utilities

PubMed Services Journals Database MeSH Database Single Citation Matcher Batch Citation Matcher Clinical Queries LinkOut My NCBI (Cubby)

Related Resources Order Documents NLM Catalog NLM Gateway

□ 1: Stroke. 2004 May;35(5):1057-60. Epub 2004 Apr 1.

Comment in:

- [Stroke. 2004 May;35\(5\):1061-2.](#)

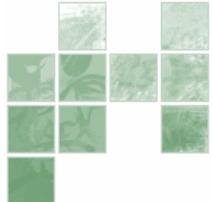
Full text article at stroke.ahajournals.org

#interact #hyper #tens , apoe , #cerebr #whit #matter #lesion .

de leeuw fe , richard f , de groot jc , van duijn cm , hofman a , van gijn j , breteler mm .

#department #epidem #logic #bio #statist , erasmus #medic #centr , rotterdam , #dutch .

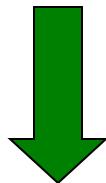
#back #ground #purpos : #cerebr #whit #matter #lesion ( wmls ) #frequent #find #magnet #resonanc #imag #scan #both #cognit #intact #dement #gero #human . #vascul #risk #factor , #special #hyper #tens , #relat #presenc . #not #total #human #vascul #risk #factor wmls ,# suggest #interact #other #determin, eg ,



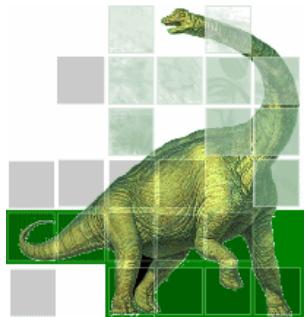
# Morphosaurus Search



„Korrelation von  
Hypertonie und  
Läsion der  
Weißen Substanz“



„#correl #hyper  
#tens #lesion #whit  
#matter“



Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Forward Reload Stop

Home Bookmarks mozilla.org Latest Builds

Entrez PubMed

Overview Help | FAQ Tutorial New/Noteworthy E-Utilities

PubMed Services Journals Database MeSH Database Single Citation Matcher Batch Citation Matcher Clinical Queries LinkOut My NCBI (Cubby)

Related Resources Order Documents NLM Catalog NLM Gateway

□ 1: Stroke. 2004 May;35(5):1057-60. Epub 2004 Apr 1.

Comment in:

- [Stroke. 2004 May;35\(5\):1061-2.](#)

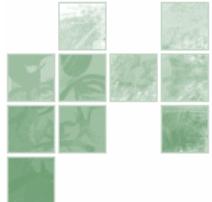
Full text article at stroke.ahajournals.org

#interact #hyper #tens , apoe , #cerebr #whit #matter #lesion .

de leeuw fe , richard f , de groot jc , van duijn cm , hofman a , van gijn j , breteler mm .

#department #epidem #logic #bio #statist , erasmus #medic #centr , rotterdam , #dutch .

#back #ground #purpos : #cerebr #whit #matter #lesion ( wmls ) #frequent #find #magnet #resonanc #imag #scan #both #cognit #intact #dement #gero #human . #vascul #risk #factor , #special #hyper #tens , #relat #presenc . #not #total #human #vascul #risk #factor wmls ,# suggest #interact #other #determin , eg ,



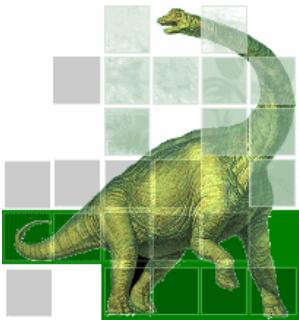
# Morphosaurus Search



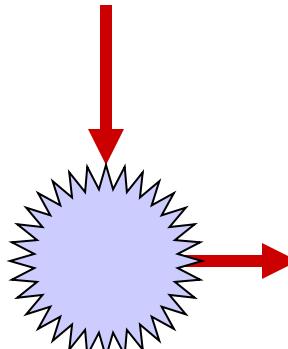
„Korrelation von Hypertonie und Läsion der Weißen Substanz“



„#correl #hyper  
#tens #lesion #whit  
#matter“



Search Engine



Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Forward Reload Stop

Home Bookmarks mozilla.org Latest Builds

1: Stroke. 2004 May;35(5):1057-60. Epub 2004 Apr 1.

Comment in:

- Stroke. 2004 May;35(5):1061-2.

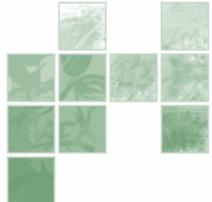
Full text article at stroke.ahajournals.org

#interact #hyper #tens , apoe , #cerebr #whit #matter #lesion .

de leeuw fe , richard f , de groot jc , van duijn cm , hofman a , van gijn j , breteler mm .

#department #epidem #logic #bio #statist , erasmus #medic #centr , rotterdam , #dutch .

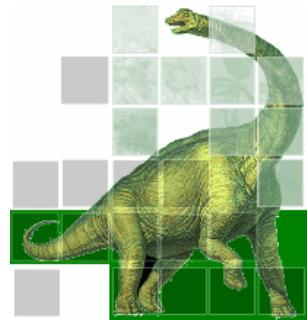
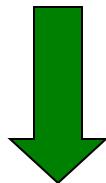
#back #ground #purpos : #cerebr #whit #matter #lesion ( wmls ) #frequent #find #magnet #resonanc #imag #scan #both #cognit #intact #dement #gero #human . #vascul #risk #factor , #special #hyper #tens , #relat #presenc . #not #total #human #vascul #risk #factor wmls , # suggest #interact #other #determin , eg ,



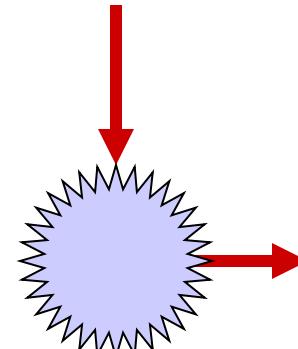
# Morphosaurus Search



„Korrelation von Hypertonie und Läsion der Weißen Substanz“



Search Engine



Entrez PubMed - Mozilla

File Edit View Go Bookmarks Tools Window Help

Back Forward Reload Stop

Home Bookmarks mozilla.org Latest Builds

Search Print

Related Articles, Links

□ 1: Stroke. 2004 May;35(5):1057-60. Epub 2004 Apr 1.

Comment in:

- [Stroke. 2004 May;35\(5\):1061-2.](#)

Full text article at [stroke.ahajournals.org](#)

#interact #hyper #tens , apoe , #cerebr #whit #matter #lesion .

de leeuw fe , richard f , de groot jc , van duijn cm , hofman a , van gijn j , breteler mm .

#department #epidem #logic #bio #statist , erasmus #medic #centr , rotterdam , #dutch .

#back #ground #purpos : #cerebr #whit #matter #lesion ( wmls ) #frequent #find #magnet #resonanc #imag #scan #both #cognit #intact #dement #gero #human . #vascul #risk #factor , #special #hyper #tens , #relat #presenc . #not #total #human #vascul #risk #factor wmls ,# suggest #interact #other #determin , eg ,

Entrez PubMed

Overview

Help | FAQ

Tutorial

New/Noteworthy

E-Utilities

PubMed Services

Journals

Database

MeSH Database

Single Citation Matcher

Batch Citation Matcher

Clinical Queries

LinkOut

My NCBI (Cubby)

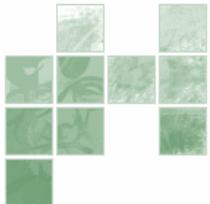
Related Resources

Order

Documents

NLM Catalog

NLM Gateway



# Automatic Language Acquisition

- Automatic Acquisition of **Spanish** and **Swedish** subword lexicons
- Step 1: Generation of *cognate* seed lexicons:
  - Automatic generation of cognate subword candidates
    - Spanish cognates from Portuguese
    - Swedish cognates from English and German
  - Selection of subword candidates
  - Semantic mapping (linkage to equivalence classes)
  - Validation of semantic mappings
- Step 2: Use cognate lexicons as a seed for iteratively learning *non-cognates*



# Step 1: Cognate Acquisition

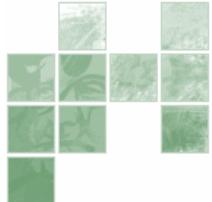
- **Resources** for cognate acquisition for **Spanish** (from Portuguese) and **Swedish** (from English and German):
  - Portuguese (~14,000 stems), German and English (~22,000 stems each) subword lexicons
  - Medical corpora for Portuguese, German, English, Spanish and Swedish acquired from the Web
  - Word frequency lists generated from these corpora
  - Manually created list of Spanish and Swedish affixes



# Generating Cognate Candidates

List of Portuguese  
Subwords  
(14,004 stems):

...  
**estomag**  
**mulher**  
...



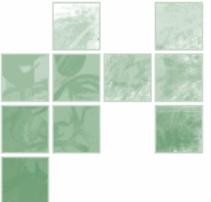
# Generating Cognate Candidates

List of Portuguese  
Subwords  
(14,004 stems):

...  
**estomag**  
**mulher**  
...

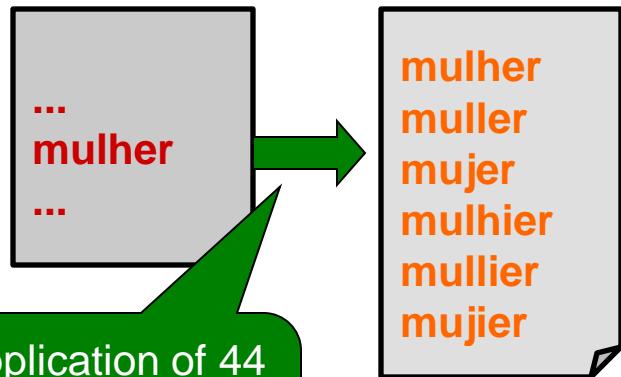
Application of 44  
string  
substitution  
rules

Rule (Port. » Span.)	Portuguese example	Spanish example	English equivalent
qua » cua	<b>quadr</b>	<b>cuadr</b>	<i>frame</i>
eia » ena	<b>veia</b>	<b>vena</b>	<i>vein</i>
ss » s	<b>fracass</b>	<b>fracas</b>	<i>fail</i>
lh » j	<b>mulher</b>	<b>mujer</b>	<i>woman</i>
l » ll	<b>lev</b>	<b>lllev</b>	<i>take</i>
i » y	<b>ensai</b>	<b>ensay</b>	<i>trial</i>
f » h	<b>formig</b>	<b>hormig</b>	<i>ant</i>
...	...	...	...

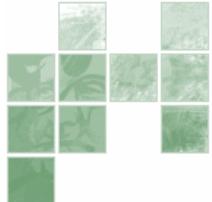


# Generating Cognate Candidates

List of Portuguese  
Subwords  
(14,004 stems):

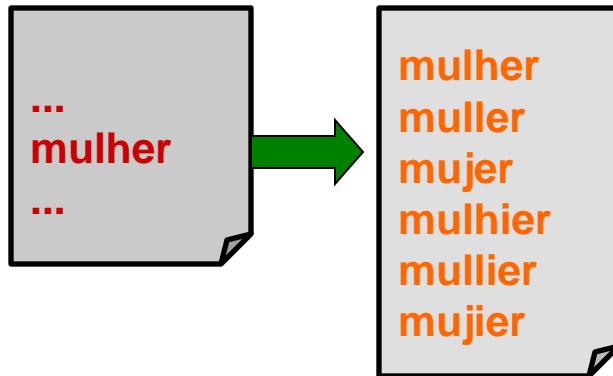


Application of 44  
string  
substitution  
rules

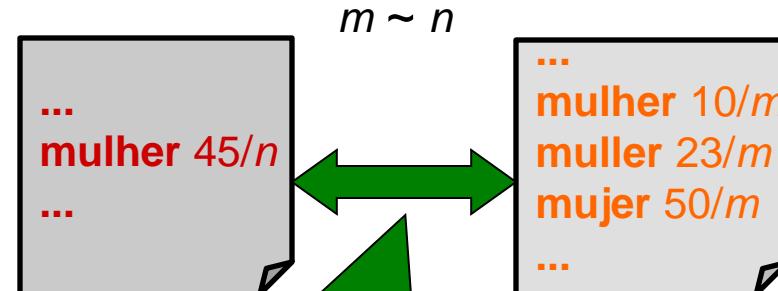


# Selecting Cognate Candidates

List of Portuguese  
Subwords  
(14,004 stems):

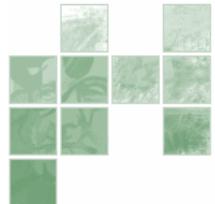


Word frequency lists  
derived from unrelated corpora:  
**Portuguese** (size =  $n$ )      **Spanish** (size =  $m$ )



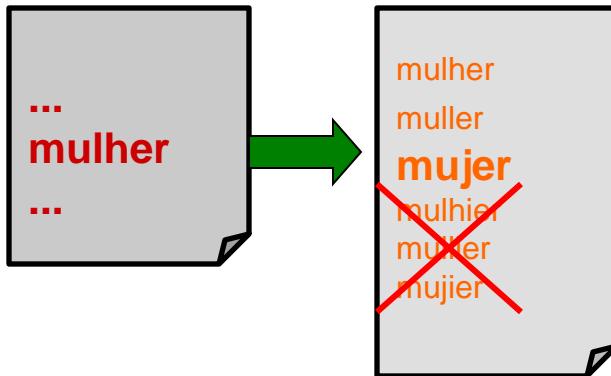
Comparison between word frequency lists:

- Elimination of non-matching subwords
- Choose that cognate alternative with the *most similar* corpus frequency

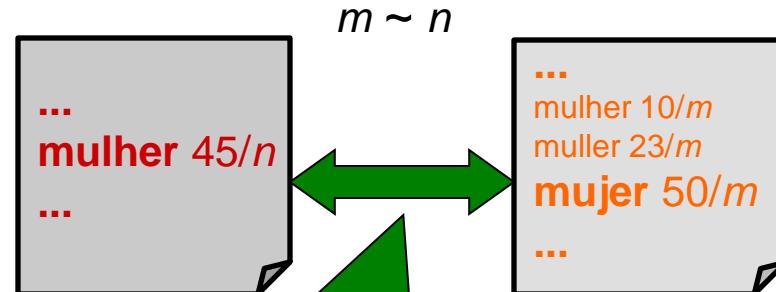


# Selecting Cognate Candidates

List of Portuguese  
Subwords  
(14,004 stems):



Word frequency lists  
derived from unrelated corpora:  
Portuguese (size =  $n$ )      Spanish (size =  $m$ )



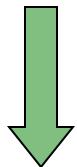
Comparison between word frequency lists:

- Elimination of non-matching subwords
- Choose that cognate alternative with the *most similar* corpus frequency



# Semantic Mapping

mulher → mujer

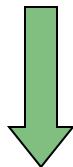


#female = { woman, women, female, frau, weib, mulher, mujer }



# Semantic Mapping

mulher → mujer



#female = { woman, women, female, frau, weib, mulher, mujer }

Language Pair	Source Lexicon	Selected Cognates
Portuguese-Spanish	14,004	8,644
German-Swedish	21,705	6,086
English-Swedish	21,501	



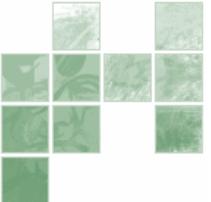
# Cognate Validation

- Use parallel corpora to identify *false friends*:
  -  Portuguese *crianc* (*child*) ↔ Spanish *crianz* (*breed*)
  - Portuguese *crianc* (*child*) ↔ Spanish *nin* (*child*)
  - Portuguese *criac* (*breed*) ↔ Spanish *crianz* (*breed*)
- UMLS Metathesaurus
  - Contains over 2M medical terms and phrases, aligned in various languages
  - English has the broadest coverage
    - English-Spanish: 60,526 alignments
    - English-Swedish: 10,953 alignments
- English-Spanish Examples
  - „Cell Growth“ ↔ „Crecimiento Celular“
  - „Heart transplantation, with or without recipient cardiectomy“ ↔ „Transplante cardiaco, con o sin cardiectomia en el receptor.“



# Cognate Validation

- Use generated cognate seed lexicons to process the UMLS alignments with the Morphosaurus system
- Whenever a MID co-occurs on both sides of an alignment unit, the lexicon entry that led to that particular MID is taken to be valid
- Candidates that never matched this procedure are discarded



# Cognate Validation

- Use generated cognate seed lexicons to process the UMLS alignments with the Morphosaurus system
- Whenever a MID co-occurs on both sides of an alignment unit, the lexicon entry that led to that particular MID is taken to be valid
- Candidates that never matched this procedure are discarded

„*Abdominal wall procedure*“:

abdomin all	#abdom
wall	#wall
proced ure	#operat



# Cognate Validation

- Use generated cognate seed lexicons to process the UMLS alignments with the Morphosaurus system
- Whenever a MID co-occurs on both sides of an alignment unit, the lexicon entry that led to that particular MID is taken to be valid
- Candidates that never matched this procedure are discarded

„*Abdominal wall procedure*“:

|abdomin|all|  
|wall|  
|proced|ure|

#abdom  
#wall  
#operat

„*Cirugia de la pared abdominal*“:

|cirug|ia|  
|pared|  
|abdomin|all|

#wall      (port. *pared*)  
#abdom    (port. *abdomin*)



# Cognate Validation

- Use generated cognate seed lexicons to process the UMLS alignments with the Morphosaurus system
- Whenever a MID co-occurs on both sides of an alignment unit, the lexicon entry that led to that particular MID is taken to be valid
- Candidates that never matched this procedure are discarded

„Abdominal wall procedure“:

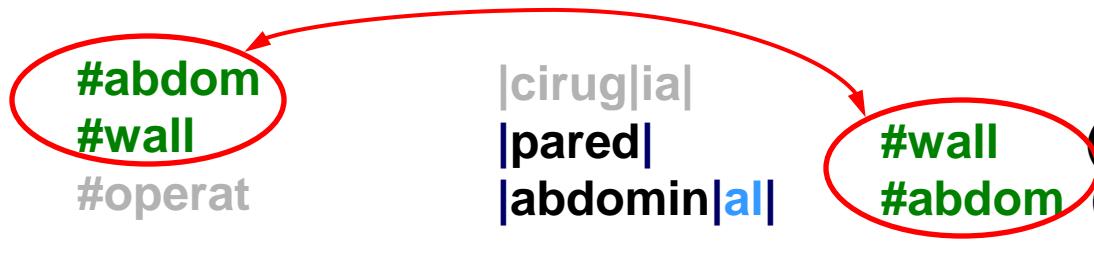
|abdomin|all  
|wall|  
|proced|ure|

#abdom  
#wall  
#operat

„Cirugia de la pared abdominal“:

cirug	ia
pared	l
abdomin	all

#wall (port. *pared*)  
#abdom (port. *abdomin*)





# Cognate Validation

- Use generated cognate seed lexicons to process the UMLS alignments with the Morphosaurus system
- Whenever a MID co-occurs on both sides of an alignment unit, the lexicon entry that led to that particular MID is taken to be valid
- Candidates that never matched this procedure are discarded

„Abdominal wall procedure“:

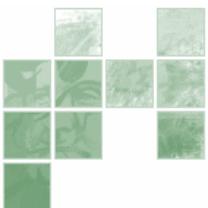
|abdomin|all  
|wall|  
|proced|ure|

#abdom  
#wall  
#operat

„Cirugia de la pared abdominal“:

cirug	ia
pared	
abdomin	all

#wall (port. *pared*)  
#abdom (port. *abdomin*)



# Cognate Validation

Language Pair	Hypotheses	Valid
English-Spanish	8,644	3,230 (37%)
English-Swedish	6,086	1,565 (26%)

„Abdominal wall procedure“:

|abdomin|all  
|wall|  
|proced|ure|

#abdom  
#wall  
#operat

„Cirugia de la pared abdominal“:

cirug	ia
pared	
abdomin	all

#wall (port. *pared*)  
#abdom (port. *abdomin*)



# Step 2: Bootstrapping

- Bootstrapping dictionaries using
  - validated cognate seed lexicons and
  - parallel corpora
  - for acquiring *non-cognates*

„*Abdominal wall procedure*“:

|abdomin|all|      #abdom  
|wall|                #wall  
|proced|ure|          #operat

„*Cirugia de la pared abdominal*“:

|cirug|ia|            #wall     (port. *pared*)  
|pared|                #abdom   (port. *abdomin*)  
|abdomin|all|



# Bootstrapping Algorithm

→ For every alignment in the UMLS do

→ „*Abdominal wall procedure*“: „*Cirugia de la pared abdominal*“:

|abdomin|all|      #abdom  
|wall|                #wall  
|proced|ure|        #operat

|cirug|ia|            #wall  
|pared|                #abdom  
|abdomin|all|        #abdom



# Bootstrapping Algorithm

For every alignment in the UMLS do

- If there is exactly one invalid segmentation in target language

„Abdominal wall procedure“:

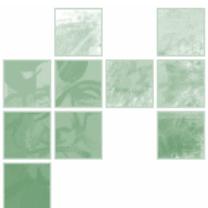
|abdomin|all|  
|wall|  
|proced|ure|

#abdom  
#wall  
#operat



„Cirugia de la pared abdominal“:

|cirug|ia|  
|pared|  
|abdomin|all|  
#wall  
#abdom



# Bootstrapping Algorithm

For every alignment in the UMLS do

- If there is exactly one invalid segmentation in target language
- If there is exactly one more MID in source language



„Abdominal wall procedure“:

|abdomin|all  
|wall|  
|proced|ure|

#abdom  
#wall  
#operat



„Cirugia de la pared abdominal“:

|cirug|ia|  
|pared|  
|abdomin|all|

#wall  
#abdom





# Bootstrapping Algorithm

For every alignment in the UMLS do

If there is exactly one invalid segmentation in target language

If there is exactly one more MID in source language

Take supernumerary MID and invalid segmentation from target



„Abdominal wall procedure“:

|abdomin|all  
|wall|  
**|proced|ure|**

#abdom  
#wall  
**#operat**

„Cirugia de la pared abdominal“:

|cirug|ia  
|pared|  
|abdomin|all  
#wall  
#abdom





# Bootstrapping Algorithm

For every alignment in the UMLS do

If there is exactly one invalid segmentation in target language

If there is exactly one more MID in source language

Take supernumerary MID and invalid segmentation from target

Restore invalid segmentation and strip off potential affixes



„Abdominal wall procedure“:

|abdomin|all|  
|wall|  
**|proced|ure|**

#abdom  
#wall  
**#operat**

„Cirugia de la pared abdominal“:

|**cirug|ia|** → |cirug|ia|  
|pared|  
|abdomin|all|

#wall  
#abdom



# Bootstrapping Algorithm

For every alignment in the UMLS do

If there is exactly one invalid segmentation in target language

If there is exactly one more MID in source language

Take supernumerary MID and invalid segmentation from target

Restore invalid segmentation and strip off potential affixes

Add new stem into target lexicon. Link it to source MID.



„Abdominal wall procedure“:

|abdomin|all  
|wall|  
**|proced|ure|**

#abdom  
#wall  
**#operat**

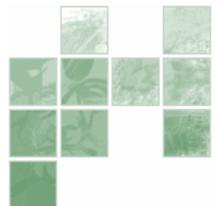
„Cirugia de la pared abdominal“:

|cirug|ia|  
|pared|  
|abdomin|all|

→ |cirug|**ia**|  
#wall  
#abdom

**#operat = { proced, surgery, operat, prozess, operier, proced, process, metod, cirug }**





# Bootstrapping Algorithm

For every alignment in the UMLS do

If there is exactly one invalid segmentation in target language

If there is exactly one more MID in source language

Take supernumerary MID and invalid segmentation from target

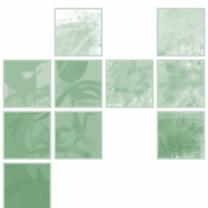
Restore invalid segmentation and strip off potential affixes

Add new stem into target lexicon. Link it to source MID.

→ Repeat all until quiescence

,,Abdominal wall procedure“:

,,Cirugia de la pared abdominal“:



# Bootstrapping Algorithm

For every alignment in the UMLS do

If there is exactly one invalid segmentation in target language

If there is exactly one more MID in source language

Take supernumerary MID and invalid segmentation from target

Restore invalid segmentation and strip off potential affixes

Add new stem into target lexicon. Link it to source MID.

Repeat all until quiescence

„Abdominal wall procedure“:

„Skin operations“:

|skin|

|operat|ions|

#derma

#operat

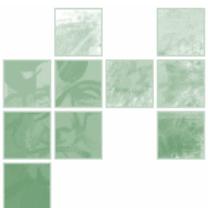
„Cirugia de la pared abdominal“:

„Cirugia de piel“:

|cirug|ia|

|piel|

#operat



# Bootstrapping Algorithm

For every alignment in the UMLS do

If there is exactly one invalid segmentation in target language

If there is exactly one more MID in source language

Take supernumerary MID and invalid segmentation from target

Restore invalid segmentation and strip off potential affixes

Add new stem into target lexicon. Link it to source MID.

Repeat all until quiescence

„Abdominal wall procedure“:

„Skin operations“:

|skin|

|operat|ions|

#derma

#operat

„Cirugia de la pared abdominal“:

„Cirugia de piel“:

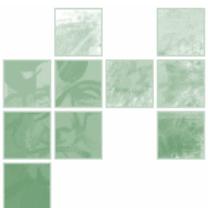
|cirug|ia|

|piel|

#operat

|piel|

#derma = { **derm**, **cutis**, **skin**, **haut**, **kutis**, **pele**, **cutis**, **piel** } ←



# Bootstrapping Algorithm

For every alignment in the UMLS do

If there is exactly one invalid segmentation in target language

If there is exactly one more MID in source language

Take supernumerary MID and invalid segmentation from target

Restore invalid segmentation and strip off potential affixes

Add new stem into target lexicon. Link it to source MID.

Repeat all until quiescence

„Abdominal wall procedure“:

„Skin operations“:

„Skin abnormalities“:

„Cirugia de la pared abdominal“:

„Cirugia de piel“:

„Malformacion de la piel“:

|skin| #derma  
|abnorm|alities| #anomal

|malformation| #derma  
|piel|



# Bootstrapping Algorithm

For every alignment in the UMLS do

If there is exactly one invalid segmentation in target language

If there is exactly one more MID in source language

Take supernumerary MID and invalid segmentation from target

Restore invalid segmentation and strip off potential affixes

Add new stem into target lexicon. Link it to source MID.

Repeat all until quiescence

„Abdominal wall procedure“:

„Skin operations“:

„Skin abnormalities“:

|skin|

#derma

|abnorm|alities| #anomal

„Cirugia de la pared abdominal“:

„Cirugia de piel“:

„Malformacion de la piel“:

|malformation|

|piel|

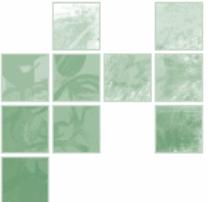
→

#derma

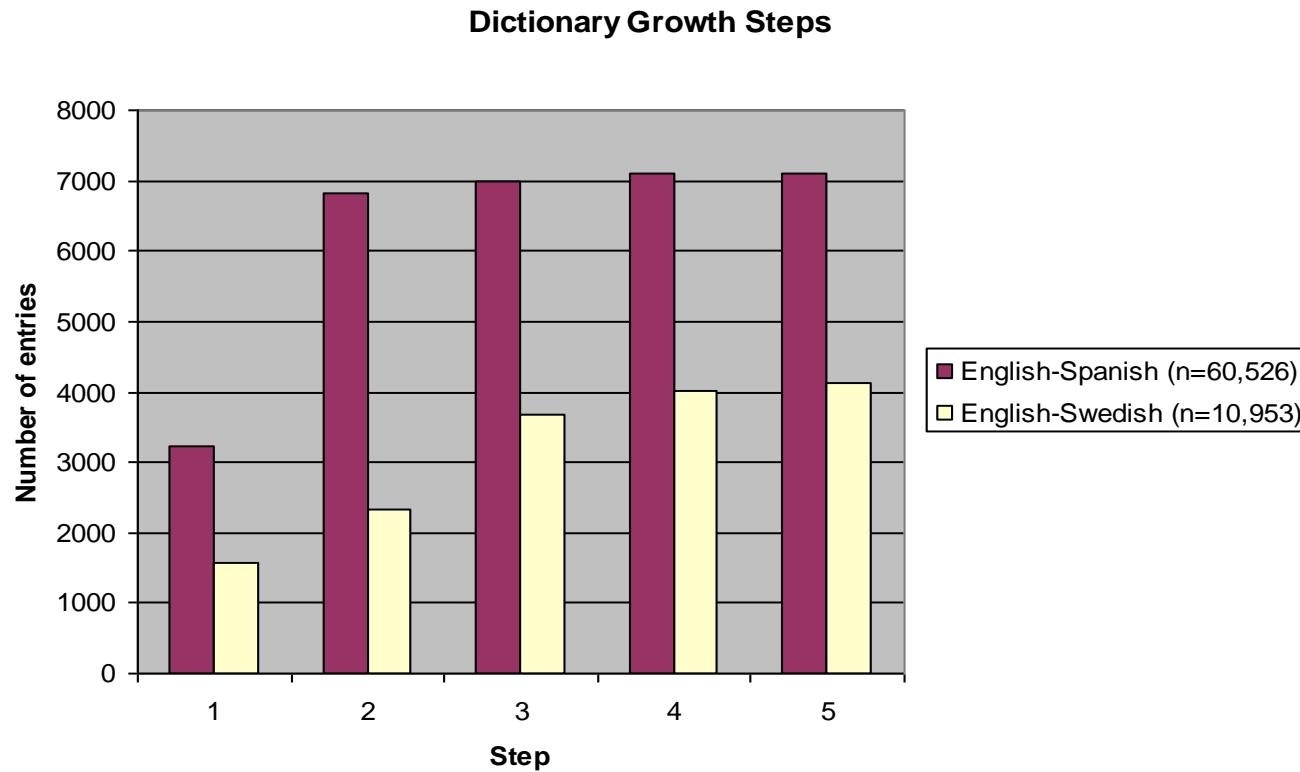
|malform|ation|

↓

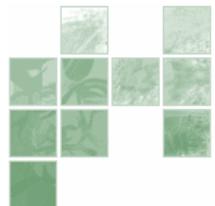
#anormal = { abnorm, anomal, abnorm, anomal, abnorm, anomal, malform }



# Bootstrapping Results

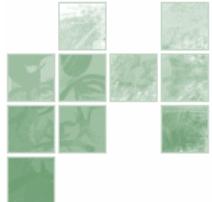


Total: 7,154 Spanish and 4,148 Swedish entries acquired

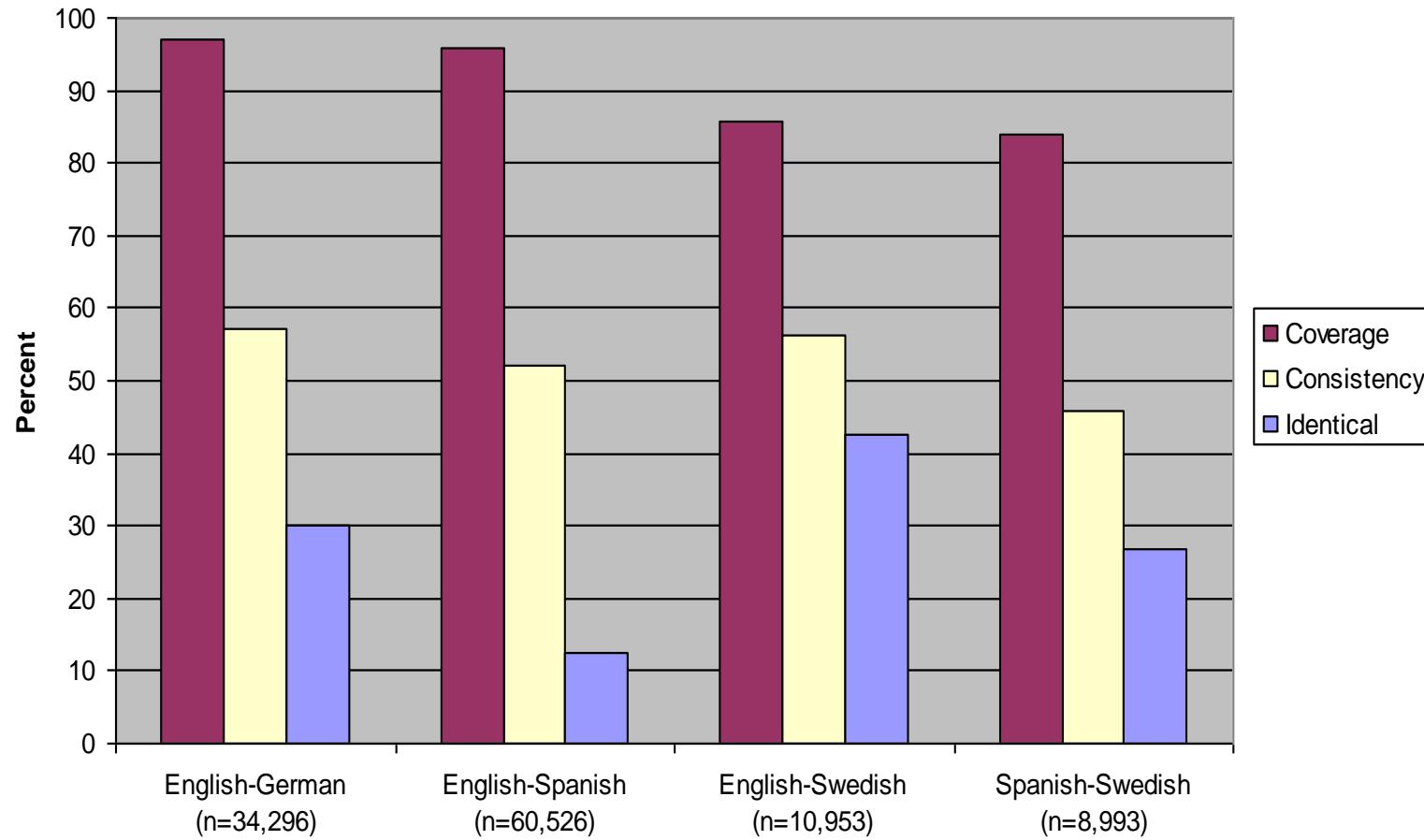


# Evaluation

- Process the English-Spanish and English-Swedish UMLS alignments with the Morphosaurus system
- Additionally process Spanish-Swedish UMLS alignments
- Measures:
  - Coverage: At least one MID co-occurs on both sides
  - Consistency 
$$C_{AU(i)} = \frac{(100 * A)}{(A + N + M)}$$
    - $A$ : Number of MIDs co-occurring on both sides
    - $N, M$ : Number of MIDs occurring on only one side
  - Identical Indexes



# Results



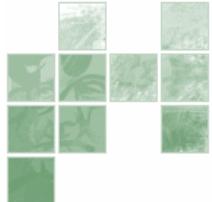


# Conclusion

- Cross-Language Document Retrieval based on the matching of search/document terms on a language-independent, interlingual layer.
- Significant amount of English, German, and Portuguese subwords can be mapped to Spanish and Swedish cognates using simple string substitution rules.
- These cognate seed lexicons are further enlarged by subword translations which are *not* cognates by bootstrapping and using parallel corpora.
- Methodology proved to be useful in a standardized CLIR experimental setting
- Generality of the approach:
  - Need of large, aligned corpora
  - Eurodicautom (12 languages, 5M entries), Eurovoc (13 languages), OECD, UNESCO, AGROVOC, etc.

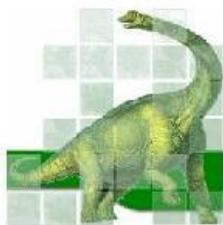


**[www.morphosaurus.net](http://www.morphosaurus.net)**



# Morphosaurus Search

## Search Engine



Search:  (English/German/Portuguese)

Language:

Information Level:

MorphoSaurus:

### Search results for 'darmkrebsrisiko':

### Best 20 Documents of 2095 matches (298 msec)

#### Entrez PubMed (94%) Expert information

Keywords: MEDLINE, NCBI, National Center for Biotechnology Information, National Library of Medicine, NLM, PubMed

Description: PubMed is the National Library of Medicine's search service that provides access to over 11 million citations in MEDLINE, PreMEDLINE, and other related databases, with links to participating online journals.

... almost unaltered **risk** of all other **cancers** (SIR, 1.2; 95% CI, 1.0-1.4), including nonelevated **risks** for several **gastrointestinal** tract **cancers**. At 10 years of follow up, the absolute **risk** of liver **cancer** was 6% among men and 1.5% among women. With 21 liver **cancers** and 508 nonhepatobiliary **cancers**, first degree ... (Cached)

[http://supreme.coling.uni-jena.de/~coling/search\\_engine\\_docs/original/medline/14724826.html](http://supreme.coling.uni-jena.de/~coling/search_engine_docs/original/medline/14724826.html)

#### Deutsches Ärzteblatt (91%) Expert information

... DEUTSCHES ÄRZTEBLATT PRINT wErhöhtes **Risiko für gastrointestinale Karzinome** nach Cholezystektomie Deutsches Ärzteblatt 99, Ausgabe 26 vom 28.06.2002, Seite A 1824 B 1541 C 1437 MEDIZIN: Referiert Eine Cholezystektomie führt möglicherweise über toxische Effekte des alkalischen Refluats auf die Speiseröhrenschleimhaut zu einem mäßiggradigen Anstieg des **Adenokarzinomrisikos** der Speiseröhre (**Barrett Karzinom**). Aber ... (Cached)

[http://supreme.coling.uni-jena.de/~coling/search\\_engine\\_docs/original/aerzteblatt/artikeldruck.asp?3id=32191.html](http://supreme.coling.uni-jena.de/~coling/search_engine_docs/original/aerzteblatt/artikeldruck.asp?3id=32191.html)

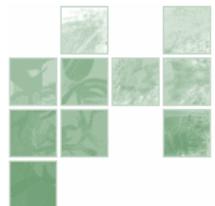
#### Dickdarm und Mastdarmkrebs Kolorektales Karzinom (91%) Patient information

Keywords: Dickdarm, Mastdarmkrebs

Description: Dickdarm und Mastdarmkrebs ist eine bösartige Schleimhautwucherung im Dickdarm oder Mastdarm. Ärzte bezeichnen diese Krebsart auch als kolorektales Karzinom (von griechisch kolon, Darm und lateinisch intestinum rectum, Enddarm). Häufig wird der Dickdarm und Mastdarmkrebs auch nur als Kolonkarzinom bezeichnet, obwohl ein KOLONkarzinom im eigentlichen Sinne nur Dickdarmkrebs ist.

... Werbung Sponsoring NetDoctor.com Dickdarm und **Mastdarmkrebs** (Kolorektales **Karzinom**) Prof. Dr. med. Stefan Endres, Facharzt für Innere Medizin und **Gastroenterologie** Was ist Dickdarm und **Mastdarmkrebs**? Dickdarm und **Mastdarmkrebs** ist eine **bösartige** Schleimhautwucherung im Dickdarm oder Mastdarm. Ärzte bezeichnen diese **Krebsart** auch als kolorektales **Karzinom** (von griechisch kolon, **Darm** und lateinisch **intestinum** ... (Cached)

[http://www.netdoktor.de/krankheiten/Fakta/dickdarm\\_mastdarmkrebs.htm](http://www.netdoktor.de/krankheiten/Fakta/dickdarm_mastdarmkrebs.htm)



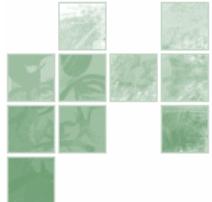
# Evaluation

- OHSUMED-Corpus (Hersh et al., 1994)
  - Subset of MEDLINE
  - ~233,000 English documents
  - 106 English user queries, additionally translated to German, Portuguese, Spanish and Swedish by medical experts
  - query-document pairs have been manually judged for relevance
- Search Engine: Lucene
  - <http://lucene.apache.org/>

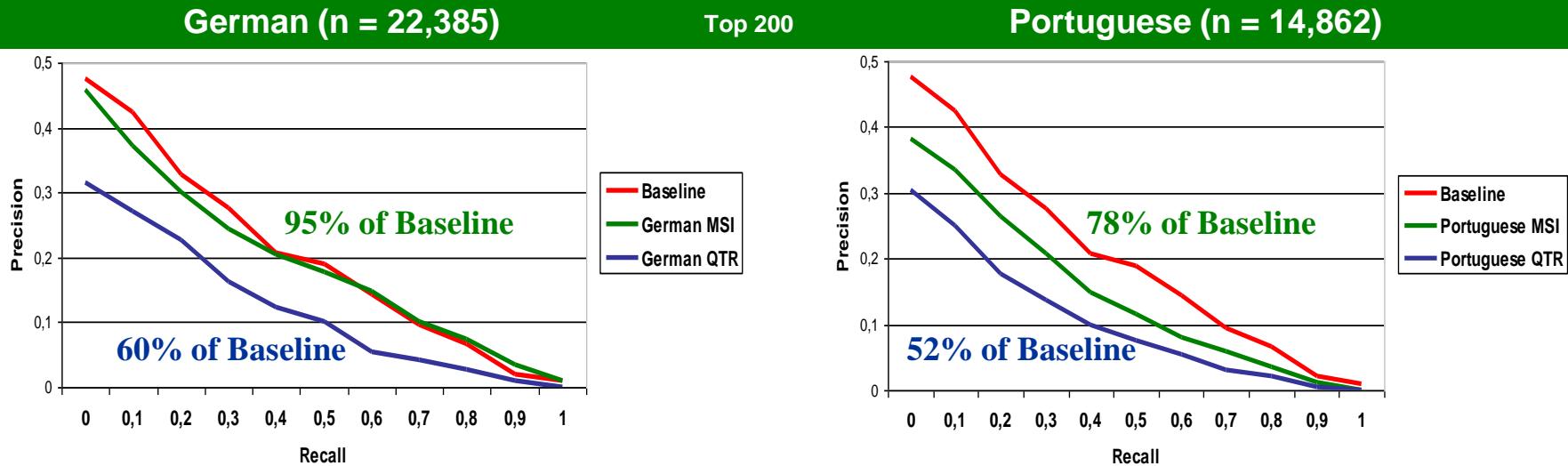


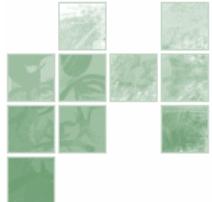
# Evaluation

- Baseline: monolingual text retrieval
  - (stemmed) English user queries
  - (stemmed) English texts
- Query translation (QTR)
  - Google translator
  - Multilingual dictionary compiled from UMLS
- Morphosaurus Indexing (MSI)
  - Interlingual representation of both user queries and documents

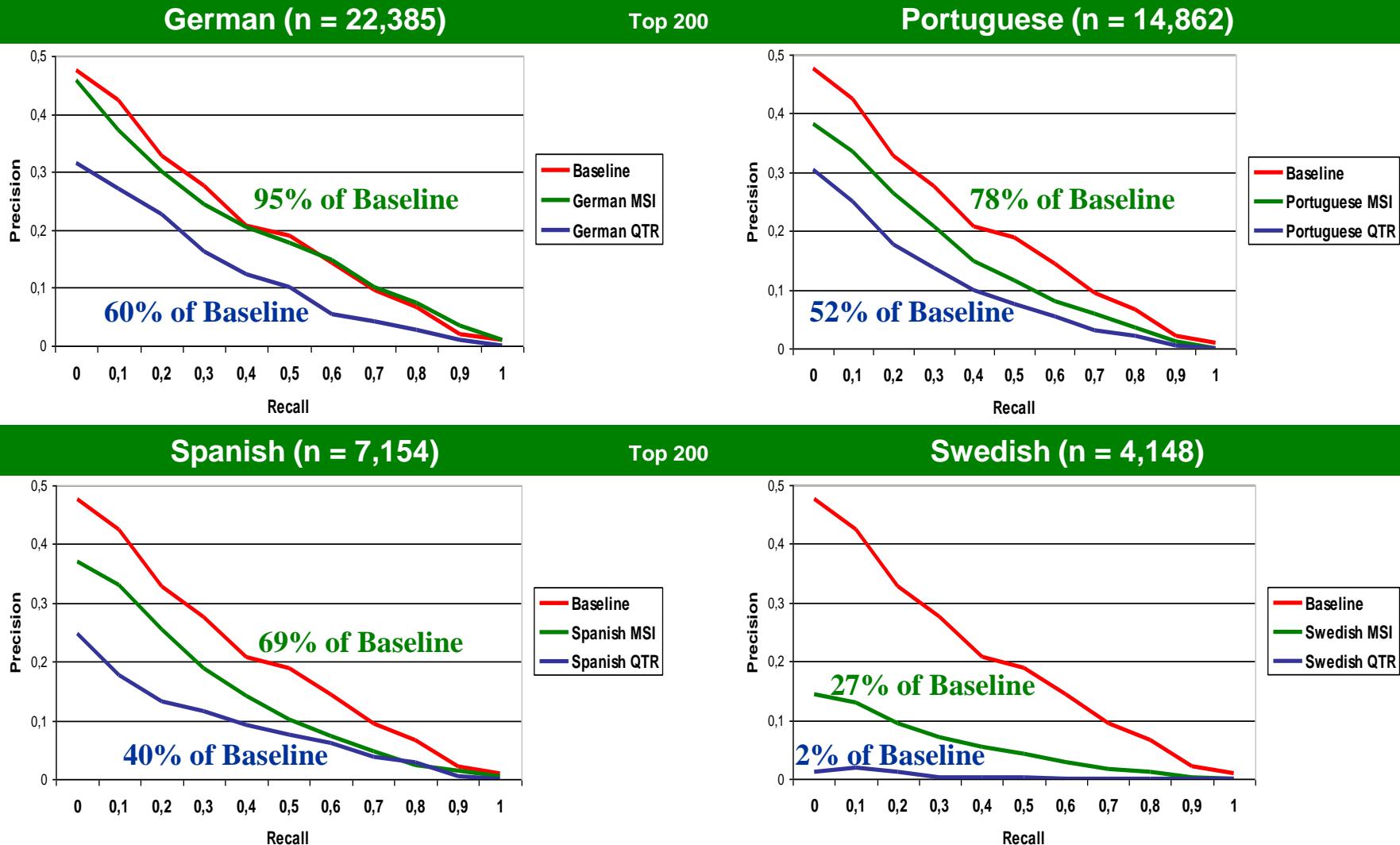


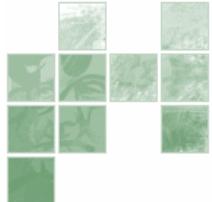
# Evaluation Results





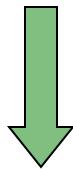
# Evaluation Results





# Semantic Mapping

mulher → mujer



#female = { woman, women, female, frau, weib, mulher, mujer }

Language Pair	Source Lexicon	Selected Cognates	Linked MIDs
Portuguese-Spanish	14,004	8,644	6,036
German-Swedish	21,705	4,249	3,308
English-Swedish	21,501	4,140	3,208
Combined Swedish Evidence (set union)		6,086	4,157