

Bilingual Lexicography: Relevance of Corpora

The example of German-Basque as language pair

David Lindemann

PhD cand.

UPV-EHU University of the Basque Country

david.lindemann.soraluze@gmail.com



Motivation for this paper

The Problem

Translation Studies, Philology, Lexicography vs. Computational Linguistics

- ▼ Experts on both sides, lack of comprehensive discipline-linking publications with focus on practical issues

Low or Medium Density Languages and more or less 'exotic' Language Pairs

- ▼ DE-EN is not DE-EU is not FI-EU

The Intention

Useful information for lexicographers that start from scratch

- ▼ as computational linguist
 - ▼ as corpus linguist
- ▼ with no or little work done on their language pair

The Example

- ▼ German-Basque

Overview

- ▼ 1. Lexicography and Corpus Linguistics
- ▼ 2. Bilingual Lexicography and Corpus Linguistics
 - ▼ 2.1 Corpus based lemma list retrieval (German)
 - ▼ 2.1.1 Comparing corpus-based vs. editorial-only lemmalist
 - ▼ 2.2 Corpora and Corpus based lemma list retrieval (Basque)
 - ▼ 2.3 Usage examples from monolingual corpora
 - ▼ 2.4 Translation examples from parallel corpora
- ▼ 3. New DE-EU dictionary: Editing, Corpora, Tools
 - ▼ 3.1 New DE-EU dictionary: frontend design
 - ▼ 3.2 New DE-EU Literary Corpus
 - ▼ 3.3 Verb entry
 - ▼ 3.4 Noun entry
 - ▼ 3.5 Adjective and Adverb
- ▼ 4. Automatic pairing of Translation Equivalents or SADD
 - ▼ Corpus-based methods, Wikimedia, WordNet
- ▼ 5. Conclusions

1. Lexicography and Corpus Linguistics

▼ Making dictionaries from corpora

- Introductions: Atkins & Rundell 2008: 53-96; Svensén 2009:43-58
- “Corpus Linguistics and Lexicography” (Teubert 1999, 2002, 2007)
- “Corpus Lexicography” (Kilgarriff & Tugwell 2002, Kilgarriff 2012)
- “Korpusgestützte Lexikographie” (Klosa 2007)

▼ Quantity: Large electronic corpora, much larger than in pre-computer age

▼ Quality: Human made index cards vs. e-corpora: The rare rather than the common

▼ Automatic retrieval of frequency lists / lemma lists

▼ Lexicographers' Documentation process when editing entries

▼ Concordances (KWIC)

▼ Cooccurrence / Collocation Statistics / Multiword LU

▼ 'Wordsketch' incl. grammar (eg for verbs: Subjects, objects, prepos., conjunct.)

▼ “Corpus-based” vs. “corpus-driven” lexicography: Trust Corpora for WSD?

▼ Evidence from parallel corpora: TE proposed by biling. dics vs. translators' choice

▼ Data from large corpora in electronic dictionary publishing

▼ On-the-fly KWIC generation as part of dictionary search result page

- as shown on some (still not many) dictionary websites

2. Bilingual Lexicography and Corpus Linguistics

▼ Monolingual and Bilingual Lexicography:

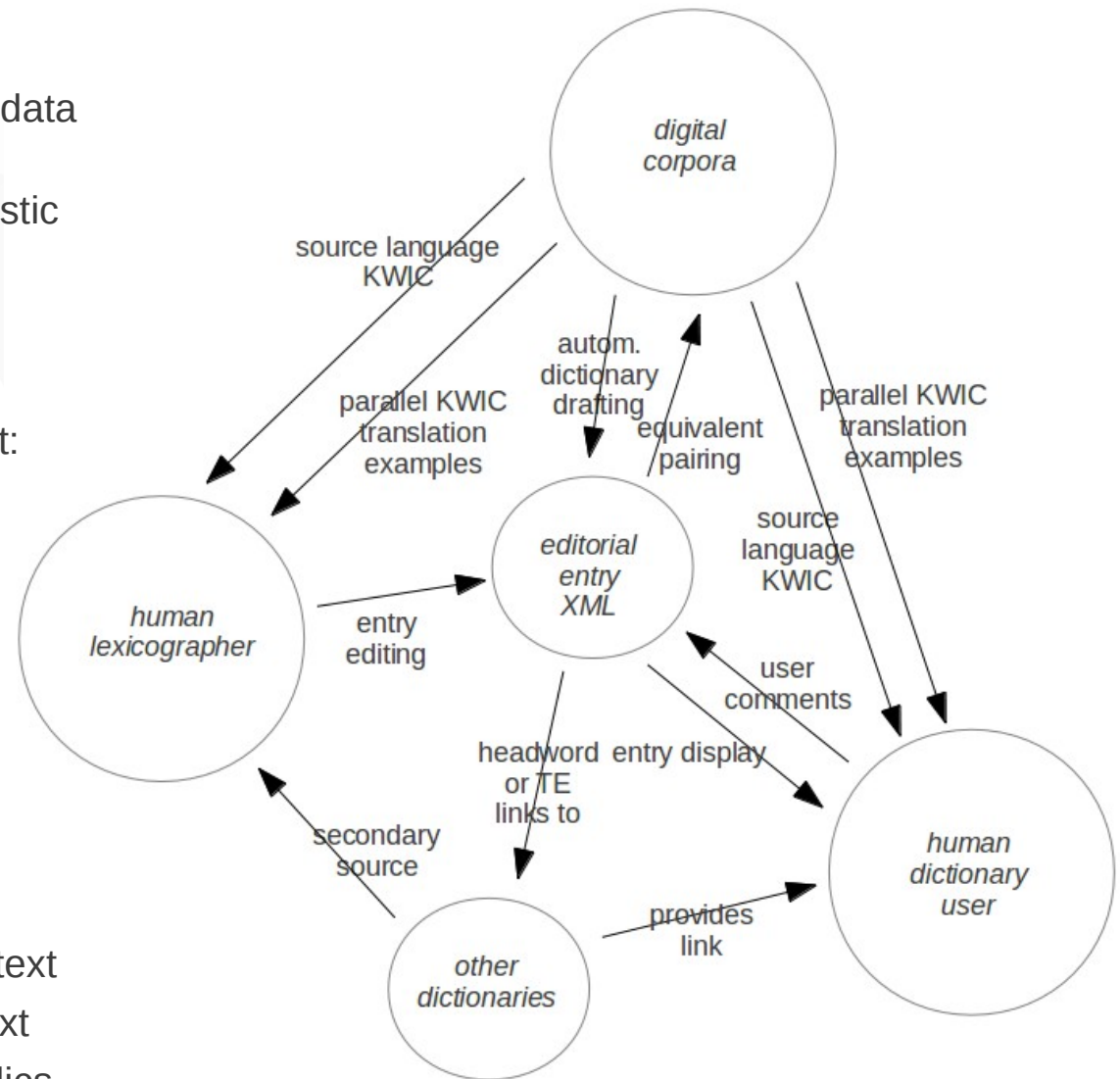
- ▼ Corpus-based lemma-list retrieval (frequency lists)
 - ▼ Frequency list becomes lemma list combining automated and hand-craft methods (IDS 2009)
 - ▼ Corpus frequency = look-up frequency (in the top few thousand) (De Schryver 2010)
- ▼ Usage examples from monolingual corpora
 - ▼ Commonplace in Lexicography: Useful for dictionary entry editing
 - ▼ Corpus Concordances shown on monolingual dictionary websites (e.g. <http://www.dwds.de>)

▼ Bilingual Lexicography:

- ▼ Automatic pairing of Translation Equivalents
 - ▼ Different strategies and software tools (examples Héja 2010, Nazar 2012)
- ▼ Translation examples from parallel corpora
 - ▼ Useful for bilingual entry editing
- ▼ Parallel KWIC as part of the dictionary search-result page
 - ▼ Teubert, *op. cit.*; first mention in: Atkins 1996, Dickens & Salkie 1996.
 - ▼ Recent, maybe leading example: <http://www.linguee.com/> (no German-Basque)
 - ▼ Mainly based on software localization files: <http://en.glosbe.com/de/eu> (also German-Basque)

Corpora in bilingual Lexicography

- Lexicographer or Dictionary editor
 - “reasoned condensation” of linguistic data provided by computers
 - “introspection”, use of their own linguistic competence
- Digital Corpora
 - monolingual: automatically built
 - parallel corpus building and alignment:
 - automatically
 - statistical methods
 - starting from (reliable) bilingual lexicon
 - semi-automatic or by hand
- Dictionary user
 - condensed in one search result:
 - editorial entry
 - headword in source language context
 - headword and TE in parallel context
 - headword or sense links to other dics



Density: Availability of digitally stored material (Varga et al. 2005)

▼ Approximation indicators

- ▼ Number of speakers of a language
- ▼ Web size of a language

▼ Parallel Corpus Building methods

▼ Web mining

- ▼ Fine for monolingual corpora even for low density languages
- ▼ For parallel corpora more difficult:

“When the other language is much lower density, the actual number of automatically detectable parallel pages is considerably smaller”

▼ Other sources

- ▼ Literary text, religious text, international law, movie dubs/subtitles, software internationalization, bilingual magazines, annual reports, corporate home pages
- ▼ Apply to German-Basque / Don't apply
- ▼ UPV-EHU: German-Basque literature corpus (81 books, 2 million tokens)

	German	Basque
Speakers	98 million	0,8 million
Biggest Corpus	5,4 billion	0,12 billion
Wikipedia Articles	1,6 million	0,15 million
ELRA Products	444	6

2.1 Corpus based lemma list retrieval (DE)

- ▼ German: DeReWo (IDS 2009)
 - ▼ Based on DeReKo (5,4 bill. words)
 - ▼ DeReWo 40.000: Lemmatized
 - ▼ Comp. w. DUDEN (2011 [240.000])

First sample: abc-first 2000

61 sorted out

23 not found in DUDEN

(29 found in another form)

20 female human gender

15 proper names (orgs, geogr.)

4 other cases

35 added

Homonyms: 8

Infinitives (corr. Part. or *-ung*): 27

2000 - 61 + 35 = 1974 (98,7%)

76 added acc. to lexicographer's subjective relevance criteria

2000 - 61 + 35 + 76 = 2050 (102,5%)

Most freq.	Least freq.	abc-first	abc-last
<u>der/die/das</u>	Kindheitstraum	á	<u>Zwischenspiel</u>
<u>und</u>	<u>Gulaschsuppe</u>	à	<u>Zwischenspurt</u>
<u>sein</u>	<u>verschreibungs</u> <u>pflichtig</u>	Aal	<u>zwischenstaatlich</u>
<u>in</u>	Superminister	Aar	Zwischenstand
<u>ein</u>	Seniorentreffen	ab	<u>Zwischenstopp</u>
<u>werden</u>	<u>Schlampe</u>	<u>abändern</u>	<u>Zwischenzeit</u>
<u>von</u>	<u>Rinnsal</u>	<u>abarbeiten</u>	<u>zwischenzeitlich</u>
<u>mit</u>	<u>Kronprinzessin</u>	Abbau	<u>Zwist</u>
<u>haben</u>	<u>Hockeyspieler</u>	<u>abbauen</u>	<u>zwitschern</u>
<u>zu</u>	Heimaufsicht	<u>abbekommen</u>	<u>zwölf</u>
<u>für</u>	Führungsmitgli ed	<u>Abberufung</u>	<u>Zwölf</u>
<u>er/sie/es/sie</u>	<u>Wegweisung</u>	<u>Abbestellung</u>	<u>Zwölfer</u>
<u>auf</u>	Umbaukosten	<u>abbezahlen</u>	<u>zwölfjährig</u>
<u>nicht</u>	<u>Straßenbahnha</u> <u>ltestelle</u>	<u>abbiegen</u>	zwölfköpfig
<u>eine</u>	<u>Stahlrohr</u>	<u>Abbild</u>	<u>zwölfmal</u>

DeReWo-40.000	EAH-2007	LDS	DeReWo Pos.	DeReWo Freq	Leipzig Freq.
Befehl	Befehl	Befehl	5926	13	12
befehlen	befehlen	befehlen	10262	14	16
	Befehlsform	Befehlsform			21
	Befehlsgewalt				16
Befehlshaber	Befehlshaber	Befehlshaber	18747	16	14
		Befehlsverweigerung			16
befestigen	befestigen	befestigen	8204	14	15
befestigt			13643	15	13
Befestigung	Befestigung	Befestigung	21221	16	15
	befeuchten	befeuchten			18
befeuern			33303	17	16
befinden	befinden	befinden	569	9	10
Befinden	Befinden	Befinden	24672	16	15
befindlich		befindlich	5499	13	19
Befindlichkeit		Befindlichkeit	13423	15	15
		befingern			21
		beflaggen			21
	beflecken	beflecken			18
beflissen			39294	17	17
beflügel		beflügel	8895	14	13
beflügelt			25788	16	13
		befohlen			14
befolgen	befolgen	befolgen	12262	15	14
	Befolgung				17
	Beförderer				20
befördern	befördern	befördern	4647	13	12
befördert			12377	15	12
Beförderung	Beförderung	Beförderung	9731	14	13
	Beförderungsmittel	Beförderungsmittel			19
		befrachten			20
befragen	befragen	befragen	2829	12	12
befragt			14584	15	11
Befragte			4367	13	14
Befragung	Befragung	Befragung	5226	13	11
befreien	befreien	befreien	2234	11	11
Befreier	Befreier	Befreier	30667	17	15
befreit			8821	14	11
Befreiung	Befreiung	Befreiung	4917	13	11
Befreiungsarmee			23679	16	16
Befreiungsbewegung		Befreiungsbewegung	30167	17	16
		Befreiungskampf			17
Befreiungsschlag			14791	15	13
Befreiungstiger			33275	17	14
	Befreiungsweg				21
befremden		befremden	15373	15	15
		Befremden			15
befremdlich		befremdlich	19219	16	14
	befreunden	befreunden			19
befreundet		befreundet	5158	13	13

2.1.1 Comparing corpus-based list with dictionary lemmalists

- ▶ Not on freq.list but freq. (according to other sources)
- ▶ Not on freq.list, less freq.
- ▶ Not on freq.list, not freq.
- ▶ On freq.list, not on DUDEN lemmalist, to be erased (proper names)
- ▶ Gaps from a human lexicographer's point of view

...on a freq.list sample of 33:

- 2 to be erased
- 3 frequent to be included
- 4 less freq. to be included

German Lemma List: Problems (1)

→ Schnorr 1991
→ Löttsch 1991

▼ Partizipien

▼ Partizip I mit und ohne adjektivischer Bedeutung

- ▼ Eigenes Lemma, wenn im Duden als "Adjektiv": *alleinerziehend*
- ▼ Infinitiv auf Derewo, Part I. nicht, aber Part. I im Duden, Verbinfinitiv nicht: Geänderte Aufnahme (*andersdenken>andersdenkend*)

▼ Partizip II mit und ohne adjektivischer Bedeutung

- ▼ Wie im Duden als "Adjektiv" und/oder "Partizip"
- ▼ Link vom Partizip zum Verbinfinitiv (der bei Fehlen aufgenommen wird)

▼ Weibliches Human Gender

- ▼ Zusammen mit der männlichen Form als Lemma-Nebenform

▼ "Transparente Komposita" (Mehrwortsyntagmen)

▼ *Abschlusstabelle, Abschleppwagen, Holzhaus, Küchentisch*

- ▼ Wenn auf DEREWO und im Duden-230.000, dann als Lemma aufnehmen (auch hier: Frequenz als Kriterium)
- ▼ "Postulat der optimalen Platzverwendung" gilt nicht
- ▼ Erscheinen im WB für DE L2 Produktion sinnvoll, über EU-Register suchbar.

German Lemma List: Problems (2)

- ▼ Substantiv-Nennform immer Singular?
 - ▼ (Quasi-)Plurale Tantum
 - ▼ Aufgenommen wie im Duden, DEREWO-Form bei Bedarf geändert
 - ▼ *Abbrucharbeiten, Achtzigerjahre, Arbeitsbedingungen*
- ▼ Nicht auf DEREWO enthaltene Lemmata
 - ▼ Gewonnen aus Abgleich mit den Lemmastrecken der drei Sekundärquellen
 - ▼ *Aas, Aasfresser, abbinden, abbüßen, Abendrot, Abendschule, abgefeimt, abgekartet, abgepackt, abgeschmackt, abgesehen, abgespannt, abgestanden, abgetakelt, abgetreten, abgewetzt, abgezehrt, abgießen, Abguss, abhärten, abkacken, abkehren, Abklatsch, Abkömmling, abküssen, ablecken, abmessen, Abrede, abreiben, Abreibung, Abreise, Abriss, abrüsten, Abschleppwagen, Absolutismus, Absonderlichkeit, absondern, Abspann, abspülen, abstillen, abstinent, abstumpfen, abteilen, abwertend, abwracken, achtsam, Adoptivbruder, Adoptivkind, Adoptivschwester, Adoptivsohn, Adoptivtochter, Afroamerikaner, Ähre, Akkordarbeit, Aktiv, Alabaster, Albaner, Alias, allwissend, altklug, am, amen, Amphibie, Andalusier, andersartig, anfeuchten, anflehen, Angeberei, angeberisch, Angemessenheit, ankern, Anschnallgurt, anspitzen, Antiquitätenhändler*

German Lemma List: Problems (3)

▼ Nicht übernommene DEREWO-Einträge:

- ▼ **Nicht im DUDEN:** *á, abend, Abfahrts, abgründen, ablaufend, aborigines, Abstiegsrunde, Abstimmungskampf, abwartend, Abwässer, accessoires, achtelfinal, Achtelfinals, achtzehnt, Acrylbild, Adventkonzert, Adventsbasar, adventsingend, afro, Aktienkursus, alleinbleiben, allround, Allstar, allzuoft, allzusehr, allzuviel, Alpenliga, Alpenstraße, Alpin, Altenheimen, Alterssiedlung, altertum, Altertümer, älterwerden, Altstadtfest, Altstadtfreund, Ampelregelung, Anders, Anerkennungspreis, Angeln, anglo, Angriffsbemühung, anlaufend, Anwaltskosten, Arbeits, Arbeitsamtsbezirk, Arbeitsjahr, Arbeitsmarktdatum, Arbeitsmarktexperte, Arbeitsmarktlage, Arbeitspartei, Arbeitsschwerpunkt*
- ▼ **Spez. Eigennamen:** *Achensee, Ackermann, Akademietheater, Alaska, Alb, Allen, Alster, Am, Amnesty, AOK, Aphrodite, Apostelkirche, Appenzell, Appenzeller, Apulien*

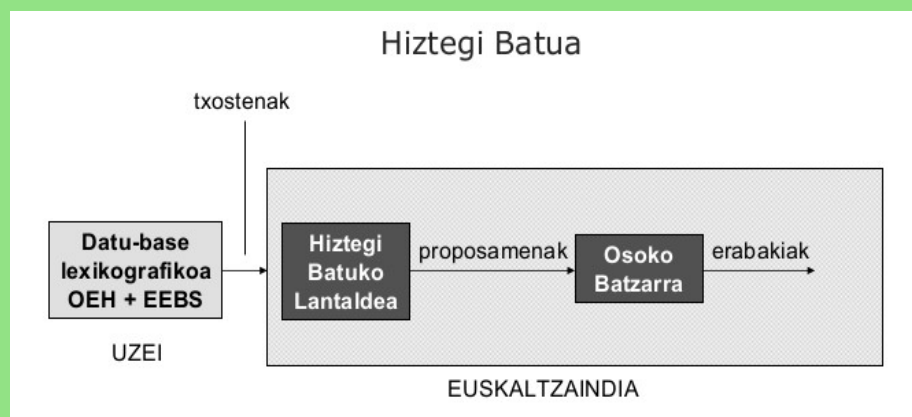
German Lemma List: Prediction on [2000] sample

- ▼ German Lemmata:
 - ▼ [DEREWO-40.000] \cap [DUDEN-230.000] \approx 38.700
 - ▼ + 5-6% from Secondary Sources
 - ▼ results in approx. 42.000 Lemma signs
- ▼ Basque Lemma Sources (future)
 - ▼ *Hiztegi Batua* 2011 (52.615 lemmata)
 - ▼ freely available XML (GPL)
 - ▼ *Basque WordNet* (26.727 lemmata)
 - ▼ ELRA resource (not free)
 - ▼ Elhuyar/EHU Science & Technology Corpus frequency list (78.600 lemmata)

2.2 Corpora and Corpus based lemma list retrieval (EU)

Euskaltzaindia (Language Academy): Lexicography

- ▼ OEH Corpus: All (304) printed books (religious and secular literature) until 1970: 5 mill. words
- ▼ **UZEI 20th century corpus**: 4,6 mill. words
- ▼ *Automatic filtering of lemma candidates, final decision one-by-one by editorial board*
- ▼ *Hiztegi Batua*: 52.615 entries



UPV-EHU/Elhuyar

- ▼ **Science and Technology Corpus**
 - ▼ 1990-2002 scientific pub.
 - ▼ 8,5 mill. words
 - ▼ filtered lemmalist: 78.600 words
- ▼ **EPG Contemporary Prose Corpus**
 - ▼ from 2000 onwards
 - ▼ 13,1 mill. words from 287 books
 - ▼ 12 mill. words from press
 - ▼ frequency lists or lemmalists?
- ▼ **XUXEN spell-checker/EDBL**
 - ▼ 120.000 lemmata, flexion-morphemes and irregular word forms

2.3 Usage examples from monolingual corpora

- German **DWDS monolingual**: Entries in several editorial dictionaries, wordsketch, KWIC of several corpora, hit statistics. Appearance fully customisable.
- For German, leading in dictionary publishing



Ressourcen - Erschließung - Projekt - Aktuelles -

Begriff Retrodigt. +Ressourcen

Der Tagesspiegel

Treffer: 8981

1	...haften Muster, das sie mit den Begriffen "Feierabend, Kühlschrank, Wein...
2	... Stam (1899 - 1986) wiesen die Begriffe noch auf die soziale Verantwort...
3	...hr auffällt. Nehmen wir den Begriff Ein-Euro-Jobber. Er ist zwa...
4	...r Sowjetunion gewählt, und die Begriffe Glasnost und Perestrojka waren...
5	... Er benutzt gern militärische Begriffe , äußert sich zu Verteidigungs...
6	...bis 9. Juli). Und würde der Begriff Nonkonformist noch in aller an...
7	...deutung dieser sehr abstrakten Begriffe zu vermitteln? Wie jedes ne...
8	... Von Elke Windisch, Moskau Der Begriff "strategische Partnerschaft" i...
9	...meisten jüngeren Deutschen ein Begriff , doch wie der Bundesrat funk...
10	...mt das Haus der Kulturen große Begriffe in einer globalisierten Welt a...

DDC-Query | Darstellung | Suchfilter

Korpusfrequenzen

Korpus	Hits	Hits [ppm]	Korpusgröße [Mill. Token]
Berliner Tagesspiegel	8981	54.3685	165.19
Berliner Zeitung	10628	41.9471	253.37
C4-Korpus	5052	63.1500	80.00
Compact Memory Corpus	3101	118.1451	26.25
DDR-Korpus	738	85.2524	8.66
Die ZEIT & ZEIT Online	71350	155.1087	460.00
DWDS-Kernkorpus	18196	180.8730	100.60
DWDS-Korpus21	145	77.5968	1.87
Juillard-Korpus	120	240.0000	0.50

Kernkorpus 21

Treffer: 145

1	...ges Mal wurde von Plessner der Begriff der Kulturnation im emphatisch...
2	...in sentimental-oppositioneller Begriff des zu Höherem bestimmten Mens.
3	...der DDR-Gesellschaft hatte den Begriff der Politik, so schien es, auf...
4	...richt lange darauf, sowohl den Begriff der » Deutschen Nation « als a...
5	...nd Haltungen. Amerikanische Begriffe und Haltungen sind es nicht. ...
6	...st eben, als könnte man diesen Begriff nicht weit genug fassen, ihn n...
7	...nen Debatte machte ein anderer Begriff Karriere. Wegen ihrer Weige...
8	...vortretend gab ihm ganz andere Begriffe , wohin die Menschheit gelange...
9	...und Genie= Pathologie, das den Begriff Deutschum zusammen mit je...
10	...ntelligenz vor den politischen Begriffen . Das Zoon politikon, diese...

DDC-Query | Darstellung | Suchfilter

DWDS-Wörterbuch

Begriff **mask.;** -s/-es; -e **Detailsansicht**

Aussprache: ▶

Zusammensetzungen: Elementar**begriff**, Gattungsb**egriff**, Grund**egriff**, ... ↓

- wesentliche Merkmale einer Sache oder einer Gruppe von Erscheinungen, die zu einer gedanklichen Einheit zusammengefasst sind ↓
- Vorstellung ↓
 - Auffassung, Meinung ↓
- schwer, langsam von **Begriff** sein + schwer, langsam auffassen **umgangssprachlich** ↓
- im **Begriff** sein, stehen + etw. tun wollen, gerade anfangen, etw. zu tun ↓

dwdsvb-0.3.41

Openthesaurus

Synonymgruppen für **Begriff**

- Anschauung, Auffassung, **Begriff**, Denkweise, Idee, Vorstellung
- Begriff**, Benennung (**umgangssprachlich**), Benennung, Bezeichnung, Nomenklatur, Notation

OpenThesaurus 2012-10-29 [OpenThesaurus Webseite](#)

Wortprofil 2012 für **Begriff**: 5578 signifikante Verbindungen

Substantiv **Dice** MI Freq. Anzahl:

Attribut

abstrakten dehnbaren englischen gebrauchten gelaufigen
geographischer geprägten gängigen philosophischen relativer
schillernden schwammigen vagen verwendeten
zentralen

hat Genitivattribut
ist Akkusativobjekt von
ist Aktivsubjekt von

Tabellensicht | Darstellungsoptionen | Suchfilter

Kernkorpus 20 (eingeschränkte Version)

Treffer: 831

Im Kernkorpus wird zwischen Groß- und Kleinschreibung unterschieden.

1	...e Haus durchnummeriert. Der Begriff geht vielmehr auf die Bandung...
2	...assen. Die Besetzung war im Begriff , eines der zahlreichen Prohib...
3	...eichen Prozess durchlaufen die Begriffe , Trends, Wörter der » Design...
4	...her nach dem anderen dieselben Begriffe vorträgt, auf die, völlig unab...
5	...tt hängt. Da stehen jetzt Begriffe wie » Retro «, » Balance «, » ...
6	...iskutierte zu bündeln. Alle Begriffe vorne an der Tafel werden um d...
7	...Walther von Hollander noch ein Begriff . Mit 66, also schon ziemli...
8	...klassisches Erbe", und um den Begriff von der "deutschen Kulturatio...
9	... Argument, dass die Mathematik Begriffe des Unendlichen konstruiert un...
10	...eist nur noch als romantischer Begriff existiert - die "Gesamtkraft" ...

nach Anmeldung 18196 von 20794 Treffern anzeigbar DDC-Query | Darstellung | Suchfilter

Etymologisches Wörterbuch

Begriff, ferner: **Begriff**, **begreiflich**, **begrifflich**, **begriffsstutzig**

begreifen Vb. 'verstehen, geistig voll erfassen, umfassen, einschließen'. Ahd. *bigrifan* (8. Jh.), mhd. *begrifen*, mnd. *begripen*, nl. *begrijpen* bedeuten ursprünglich 'ergreifen, betasten', auch 'umfassen, enthalten' (zur Herkunft s. *greifen*). Schon in ahd. Zeit, wo es lat. *comprehendere* 'begreifen' übersetzt, später vor allem bei den Mystikern, erfährt das Wort eine Bedeutungsweiterung, indem körperliches 'Greifen, Fassen' auf geistige Aneignung, 'mit dem Verstande erfassen, verstehen', ausgedehnt wird. – **Begriff** m., mhd. *begrif* ist zu *begreifen* im Sinne von 'umfassen, enthalten' gebildet und bedeutet 'Umfang, Bezirk, Umfang und Inhalt einer Vorstellung', ferner 'Zusammenfassung, kurzer Auszug'. **Begriff** kommt durch die philosophischen Aufklärer Wolff und Thomasius allgemein in Gebrauch; **Begriff** und *Vorstellung* werden in der philosophischen Terminologie bald gegeneinander abgegrenzt, so daß **Begriff** die heute vorherrschende Bedeutung 'wesentliche Merkmale einer Sache oder einer Gruppe von Erscheinungen, die zu einer gedanklichen Einheit zusammengefaßt sind' erhält.

eywmb-1.0.27

„Begriff“	
kontzeptu	18
ideia	11
hitza	2
adigai	1
burutapena	1
ezagutza	1
gai	1
hitz	1
ikusmolde	1
pentsakera	1
termino	1
ulerbide	1
ulerkera	1

„im Begriff sein“			
-tzera joan	5	-tzekotan izan	1
-tzeko zorian	4	-tzera apoderatzen izan	1
-tzear egon	3	ekarri	1
-tzeko asmoa izan	3	gutxi falta + subjkt.	1
-tzeko asmotan egon	2	hasia izan	1
-tzea pentsatu	1	inf. nahi izan	1
-tzear izan	1	inf. nahian	1
-tzeko duda egin	1	egin behar izan	1

LU and TE candidates in Parallel Corpus (SkE). Example: “Begriff” (noun)

kein Begriff sein	ez ezagutu	1
kein Begriff sein	horren entzuterik ere ez izan	1
keinen andern Begriff haben als baino ez pentsatu	1
nicht allzu schnell von Begriff sein	ez oso azkarra izan	1
schwer von Begriff sein	burugogorra izan	1
sich einen stillen Begriff machen	gutxi gora behera irudika ahal izan	1
sich kaum einen Begriff machen	ozta-ozta ideia bat izan	1
sich keinen Begriff machen	ez jakin	1
sich keinen Begriff machen	ezin imajinatu ere egin	1
über alle Begriffe	ezin esan bezain	1
Der Begriff X	X	1
einen Begriff geben	aditzera eman	1

2.4 Translation examples from parallel corpora

Parallel KWIC from lemmatized and POS-tagged DE-EU corpus using SkE. Example: “Begriff” (noun)

Corpus: DE_1
Hits: 86 (37.3 per million) [EU_1] [EU_2]
Page 1 of 5 Go Next | Last

#95733	war frech , davon machen Sie sich keinen Begriff . Denken Sie sich , ich hab meinen Überzieher
EU_2:	Zoragarria zen ... Maiatza ... Bero egiten zuen ... Eta ni lotsagabe sentitzen nintzen , ez dakizu ondo .
#110678	Gedächtnis , weil ich gerade zu dieser Zeit im Begriffe war , die Reisebeschreibung des berühmten
EU_2:	Oroimean geratu zitzaidan ; izan ere , hain zuzen garai hartan , izen bereko ibiltari ezagunaren bidaiari buruzko kontakizunak irakurtzekotan nintzen .
#136122	sonnverbrannte Wüste geht . Und schon war Andreas im Begriffe , den Helden des Films sympathisch und
EU_2:	Eta Andreas hasia zenean heroiari sinpatia eta atxikimendua hartzen , filmea uste ez bezala zorian aldera makurtu , handik pasatzen zen zientifiko espedizio batek desertuan galdutako gizona salbatu , eta Europako zibilizazioaren gerizpera ekartzen zuen atzera .
#136185	den Helden des Films . Und schon war er im Begriff , sich zu erheben , als auf der Leinwand
EU_2:	Eta jeikitzeko zorian zegoela , hara non pantaHan agertu zen lehentxeago , baran zegoela , tabernariaren bizkar atzean ikusitako eskola lagunaren irudia .
#139853	Hier stieß er auf einen Mann , der eben im Begriffe war , die Treppe hinaufzusteigen , und
EU_2:	Han , eskaileretan gora hastear zen gizon batekin topo egin zuen . Oso ezaguna egin zitzaion .
#151414	aufgewachsen und hatte von dorthier eine Menge von Begriffen und Schablonen beibehalten . Er hatte theoretisch
EU_2:	Gainera hezimolde burges txikiari hasia zen , eta orduanik kontzeptu eta kate andana gordetzen zituen .
#154557	den Dichtungen des alten Indien ist dieser Begriff ganz unbekannt , die Helden der indischen
EU_2:	India zaharreko poesietan kontzeptu hori erabat ezezaguna da , Indiako epopeietako heroeak ez dira pertsonak , pertsona matzakak , enkarnazio sailak baizik .
#155264	Leben . Was die Menschen jeweils unter dem Begriff " Mensch " verstehen , ist stets nur eine
EU_2:	Gizonek " gizaki " kontzeptuaren bidez gehienetan ulertzen dutena , akordio burges iragankor bat besterik ez da betiere .
#160170	und erkannte mich erst , als ich schon im Begriff war , an ihm vorüberzugehen . Er stürzte
EU_2:	Parez pare zetorkidan jakituna , zurrin eta nahiko begilauso , eta jaramonik egiteke albotik pasatzea pentsatzen nuen ulean ezagutu ninduen .
#202482	Ihre Arbeit zu Ende zu führen . Wenn der Begriff der Pflicht Ihnen unbekannt ist ... " Er
EU_2:	Egitekoaren ideia ezezaguna bazaizu ...
#202528	Warten Sie ! " sagte Gustav höflich . " Den Begriff der Pflicht allerdings kenne ich nicht
EU_2:	- Egitekoaren ideia egiaz ez dut ezagutzen , ez gehiago .
#202602	Gegenteil getan . Aber wenn ich auch den Begriff der Pflicht nicht mehr kenne , so kenne
EU_2:	Baina eginbeharraren ideia gehiago ezagutzen ez badut ere , ongi ezagutzen dut ordea kulparena ; agian biak gauza bera izango dira .
#203637	Menschen , einst ein hohes Ideal , ist im Begriff , zu einem Klischee zu werden . Wir Verrückten
EU_2:	Gizaiardia , garai batean ideal altu bat , azken batean klitxe bihurtzeko zorian dago .
#204784	lückenhafte Seelenlehre der Wissenschaft durch den Begriff , den wir Aufbaukunst nennen . Wir zeigen
EU_2:	Guk horregatik , akatsez beneriko jakintza horren animari buruzko irakaskintza osatu egiten dugu , eraikuntz deritzogu n ideien bidez .
#392066	der Herr Vorsteher schwieg . Ich war im Begriff , ihm irgendeine grobe Beleidigung ins
EU_2:	Irain satsuren bat aurpegira botatzeaz egon nintzen .

Begriff m. kontzeptu | adigai ▶ ich stand im Begriff, das Haus zu verlassen: *Etxea uzteko abian nengoan*

editorial EAH 2007

14.03.13

Bilingual Lexicography DE-EU

2.4.1 Summary: Functions

- ▼ Parallel KWIC display: real translation examples
 - ▼ TEs for the headword that are not listed in the editorial entry
 - ▼ TEs in context
 - ▼ domain, register, regional
 - ▼ TEs for MWE and collocations
 - ▼ MWE: “im Begriff sein”
 - ▼ phraseological coll.: “üblicher Begriff”
 - ▼ idiom: “schwer von Begriff”
 - ▼ POS-transposition
 - ▼ in editorial entry POS is maintained even if translators don't
 - ▼ other adaptations done by human translator

berücksichtigen (V)	consider (V)	aintzat hartu (PosP+V) kontuan izan (PosP+V)
Berücksichtigung (N)	consideration (N)	aintzat hartze?? (PosP+N) kontuan izate?? (PosP+N)
unter Berücksichtigung +Gen. (Prep+N+Gen.)	considering sth.; taking sth. into account (Part. Phrases)	zbt. aintzat harturik; zbt. kontuan izanik (Part. Phrases)

Abfrageergebnis für Ihre Anfrage nach »Begriff«		
30 Treffer aus dem deutsch->englisch Lexikon:		#Belege
Begriff	notion	6
Begriff	concept	3
Begriff	term	2
begriff	comprehended	2
Begriff	idea	2
begriff	understanded	2
Begriff	sense	2
Begriff	realization	2
Begriff	acquaintance	2
Begriff	knowledge	2
Begriff	item	2
Begriff	apprehension	1
Begriff	seidea	1
Begriff	conception	1
Begriff	perception	1
im Begriff sein	to be about to	1
üblicher Begriff	commonly used term	1
aktrakter Begriff	abstract idea	1
im Begriff sein zu	be about to	2
abstrakter Begriff	abstract term	2
im Begriff sein zu	to be on the point of	2
im Begriff sein zu	to be about to	2
abstrakter Begriff	abstract concept	1
schwer von Begriff	slow on the uptake	1
im Begriff sein zu	to be going to	1
im Begriff sein zu	gonna	1
numerischer Begriff	numeric character	1
langsam von Begriff	slow-witted	1
der Begriff 'Pferd'	the concept 'horse'	1
mehrdeutiger Begriff	ambiguous concept	1

dict.uni-leipzig.de: DE-EN parallel corpus TE extraction and cooccurrence statistics as on-the-fly generated dictionary entry (under development; still noisy)

Editorial Dictionary:

- Begriff** noun, masculine
 - term n
 - concept n
 - notion n
 - idea n
 - item n
 - conception n
 - specialist term n
 - perception n
 - begriff**
 - understood
 - recognized
 - comprehended
 - saw
 - savvied
 - "er/sie/es begriff"* could be 3rd person simple past
 - begreifen** verb
 - understand v
 - comprehend v
 - see v
 - grasp v
 - realize v
 - conceive v
 - realise v [BE]
 - recognize v
 - apprehend v
 - figure out v
 - catch on v
 - grok v [sl]
 - savvy v
 - Begreifen** noun, neuter
 - comprehension n
 - apprehension n
 - uptake n
- Examples:
- im Begriff sein zu**
 - be going to
 - be on the point of
 - gonna** [AE] [colloq] [sl]
 - dieser Begriff bezieht sich auf**
 - this term refers to
 - abstrakter Begriff** noun, masculine
 - abstract term n
 - zu einem Begriff werden** verb
 - become a household name v
 - schwer von Begriff**
 - slow on the uptake adj
 - allgemeiner Begriff** noun, masculine
 - general term n
 - universal concept n
 - Übergeordneter Begriff** noun, masculine
 - superordinate n
 - hypernym n
 - Üblicher Begriff** noun, masculine
 - commonly used term n
 - Ich/er/sie begriff**
 - I/he/she understood

Translation examples from external sources for 'Begriff':

German	English
[...] Syntax der Marke sowie ihre Eignung als unmittelbar und direkt beschreibender Begriff für die fraglichen Waren und Dienstleistungen falsch beurteilt habe, (ii) nicht [...]	[...] mark, as well as its aptness or otherwise as an immediate and direct descriptive term for the goods and services in question ; (ii) failed to establish facts of its [...]
Artikel 1 Absatz 1 des Übereinkommens vom 27. September 1968 über die gerichtliche Zuständigkeit und die Vollstreckung gerichtlicher Entscheidungen in Zivil- und Handelssachen in der Fassung der Übereinkommen vom 9. Oktober 1978 über den Beitritt des Königreichs Dänemark, Irlands und des Vereinigten Königreichs Großbritannien und Nordirland und vom 25. Oktober 1982 über den Beitritt der Republik Griechenland ist dahin auszulegen, dass der Begriff "Zivilsache" eine Rückgriffsklage umfasst, mit der eine öffentliche Stelle gegenüber einer Privatperson die Rückzahlung von Beträgen verfolgt, die sie als Sozialhilfe an den geschiedenen Ehegatten und an das Kind dieser Person gezahlt hat, soweit für die Grundlage dieser Klage und die Modalitäten ihrer Erhebung die allgemeinen Vorschriften über Unterhaltsverpflichtungen gelten.	The first paragraph of Article 1 of the Convention of 27 September 1968 on Jurisdiction and the Enforcement of Judgments in Civil and Commercial Matters, as amended by the Convention of 9 October 1978 on the Accession of the Kingdom of Denmark, Ireland and the United Kingdom of Great Britain and Northern Ireland and by the Convention of 25 October 1982 on the Accession of the Hellenic Republic , must be interpreted as meaning that the concept of 'civil matters' encompasses an action under a right of recourse whereby a public body seeks from a person governed by private law recovery of sums paid by it by way of social assistance to the divorced spouse and the child of that person, provided that the basis and the detailed rules relating to the bringing of that action are governed by the rules of the ordinary law in regard to maintenance obligations.
Es folgten - um nur einige Werke zu nennen - das Buch "Der hörende Mensch", in dem er seine Forschungen auf die harmonikalischen Strukturen in der Chemie, der Astronomie, der organischen Natur, in Licht und Farbe und in der Architektur ausdehnt; die "Harmonia Plantarum", eine Entwicklung der Morphologie der Pflanzen aus harmonikalischen Grundelementen; das Büchlein "Akroasis", das eine allgemeinverständliche Einführung in die Gedankenwelt der Harmonik geben will, und das auch heute noch erhältlich ist (der Begriff "Akroasis" = Anhörung wurde von Kayser als anderer Name für die Harmonik geprägt, da der Begriff "Harmonik" oft falsch verstanden wurde) und als umfassendstes Werk das "Lehrbuch der Harmonik", in dem er die Ergebnisse seiner Forschungen zusammenfaßt und sie in Form eines Lehrbuches präsentiert.	This was followed - to mention but a few works - by Der hörende Mensch (The Hearing Human), in which he expanded his research into the harmonic structures to be found in chemistry, astronomy, organic nature, light and colour, and architecture; the Harmonia Plantarum, a development of the morphology of plants out of harmonic basic elements; the short book Akroasis, which was designed as an intelligible-to-all introduction to the intellectual world of harmonics and that is still available today (the word Akroasis = hearing was coined by Kayser as another name for harmonics, which term was often misunderstood) and the comprehensive Lehrbuch der Harmonik (Textbook of Harmonics), in which he presented the results of his investigations in textbook format.
Der Begriff "persönliches Interesse", wie er im vorstehenden Absatz verwendet wird, findet keine Anwendung auf jedwede Beziehung und jedwedes Interesse, die nur [...]	The term "personal interest" as used in the above section shall not apply to any relationship and any interest which arise only due to the fact that the legal transaction [...]
[...] unterschied zu den üblichen dokumentationen, bei denen wert auf gewicht und beständigkeit gelegt wird, zeigt sich dieses druckwerk besonders fragil und vergänglich: auf über 600 seiten sind nur die namen der künstler zu lesen. sie sind auf extrem dünnes papier gedruckt. es ist durchscheinend, wie die geräusche, die in der ausstellung zu hören sind. der charakter des tons schlägt sich nieder in der verletzlichen gestalt. der buchblock hat keinen schützenden einband, so dass bereits das erste durchblättern spuren hinterlässt. das	the catalogue, with its graphic design and sound transcriptions, is itself a direct interpretation of the transient medium of sound. unlike conventional exhibition documentation, which believes in weight and durability, this publication is highly fragile and ephemeral: its 600+ pages contain nothing but the artists names, printed on extremely thin paper. the pages are transparent, like the sounds that can be heard in the exhibition. the fragile form captures the character of sound. the body of the book has no

Editorial dictionary entry and parallel KWIC combination as e-dictionary search result (linguee.com)

Noise: word highlighted without clear reason

"vote up" / "vote down" buttons: crowd votes for good examples (affects display order)

Noise: 3 TEs for 2 appearances of the HW: "word", "term" plus "intellectual world" (TE for "Gedankenwelt"="Begriffswelt")

3.0 New DE-EU dictionary: Editing, Corpora, Tools

- Editing Software:
 - TshwaneLex (TLex)**
(de Schryver & Joffe 2005, 2010)
 - Freely designable DTD, Styles
 - Output as XML, html, rtf, tsldict
 - Attribute Lists, Label Sets
 - Import, Reverse, Compare, Merge
 - CrossRefs, Hyperlinks
- Software for Corpus Tagging and KWIC:
 - SketchEngine (SkE)**
(Kilgarriff et al. 2004)
 - Queries in large annotated corpora for DE and EU
 - Wordsketch, GDEX
 - Parallel Concordances

The screenshot displays the TshwaneLex (TLEX) software interface. The main window shows a dictionary entry for the German verb 'absetzen'. The entry is structured as follows:

- Lemma:** absetzen
- Sense 2:** EremuSemant: zama, bidalaria
- TE:** utzi
- OsagLexPre:** (empty)
- OsagLexPost:** (empty)
- ArgKasu:** -
- TE:** laga

The right-hand pane shows the English translation and usage examples for 'absetzen':

- absetzen** *brz*
 - I *tr.* -h > 1 (txapela, betaurrekoak) **kendu**; **erantzi** 2 (zama, bidalaria) **utzi**; **laga** 3 **kargugabetu** 4 **EkON** **saldu** 5 von der Steuer **absetzen** zergatik **kendu** 6 (sendagaia) (hartzeari) **utzi** 7 **ANTZERK ZINEMA** (karteletik) **kendu**
 - II *intr.* -h > (hitz egiten, idazten) **eten**; **gelditu**
 - III *rfIx.* -h > 1 (gart. > a **ihes** egin; **hanka** egin **b** **erbestera** 2 **jalki**; **sedimentatu**
- Absetzung**
 - iz. f.* > 1 **kargugabetze**; **kargutik kentze** 2 **ANTZERK ZINEMA** (karteletik) **kentze**
- absichern** *brz*
 - I *tr.* -h > **segurtatu**; **bermatu**
 - II *rfIx.* -h > **ziurtatu**; **segurtatu**
- Absicht**
 - iz. f.* > **asmo**
- absichtlich**
 - I *adb.* > **nahita**
 - II *adj.* > **nahita egindako**
- Absichtserklärung**
 - iz. f.* > **asmo-aitorpen**
- absinken** *intr. brz*
 - intr.* -sn > 1 **hondoratu**; **beheratu** 2 **gutxiagotu** 3 **gainbehera joan**; **moteldu**; **makaldu**
- absitzen** *intr. brz*
 - I *intr.* -sn > (zaldi, biziklo) (gaimetik) **eraitsi**; **jaitsi**
 - II *tr.* -h > 1 **jarg.** **denbora pasa** egon 2 (kartzela zigorra) **bete**
- absolut**
 - I *adj.* > 1 **erabete** 2 **FIS HIZK MATEM** **absolutu**
 - II *adb.* > **erabat**
- Absolutismus**
 - iz. m.* > **HIST POL** **absolutismo**
- Absolvent, Absolventin**
 - iz. bizid. m.f.* > **ikasle** **ohi**; **gradudun**
- absolvieren** -ge
 - tr.* -h > **HEZK** **amaitu**; **gainditu**; **gradua lortu**
- absonderlich**

3.1 New DE-EU Literary Corpus

- Created at UPV-EHU (Sanz Villar 2011, Zubillaga 2012)
- 81 Digital and OCR-ed literary DE originals and official EU versions
- about 2 million tokens per language
- 146.457 sentence pairs
- Sentence alignment
 - one-to-many alignment revision by hand
 - Improved sentence alignment tools like *hunalign* (Varga et al. 2005) require seed lexicons (DE-EU still not available)
- Exported to TMX format, imported in the SketchEngine
- DE: Lemmatized with TreeTagger
- EU: Tokenized, not lemmatized (at the moment)

<u>lemma</u>	<u>Freq</u>
d	145113
und	57846
sein	54259
ich	46686
er	43723
ein	39059
haben	25846
sie	25470
in	22647
zu	22589
nicht	22374
es	19506
er es sie	18760
mit	15777
auf	15423
du	14433
ihr	13580
wie	11903
von	11372
an	10944
so	10854
als	10526
aber	10518
sagen	10474
wir	10217
werden	9393
Sie sie sie	9335
noch	8271
was	7530

<u>word</u>	<u>Freq</u>
eta	58220
zuen	29551
ez	28415
zen	21353
bat	16811
ere	16025
egin	14295
da	10926
bere	9961
zion	9338
izan	8328
du	6907
baina	6592
behar	6149
egiten	6032
nahi	5746
Eta	5647
esan	5407
Ez	5143
bezala	5089
ziren	4948
zegoen	4688
Baina	4664
nuen	4493
hori	4374
dut	4182
zituen	4134
nire	4130
zuten	3887

3.2 New DE-EU dictionary: frontend design

Begriff ~s, ~e

Substantiv m. s. **Verbinf.** [begreifen](#) ▶ **1** kontzeptu; ideia; hitz; ezagutza **2** ● im [Begriff](#) sein -tzera joan; -tzeko zorian egon; -tzeko asmoa izan **3** ● sich (k)einen [Begriff](#) (von etw.) machen (ez) jakin

Article has cross-references to >>

begreifen *unregelmäßig* **begreift**, **begriff**, **begriffen**

Verb Transitiv -haben ▶ **1** = [verstehen](#) **ulertu**; **aditu**; **barrendu** **2** [*als etw.*] -tzat hartu

Hyperlinks Block:

- German monolingual dics
DWDS, Leipzig, Openthesaurus...
- Wiktionary & Wikipedia
- Comment form (user comments)

Corpus: **deTenTen**
Hits: **280464** (98.6 per million)

Page 1 of 7012 [Next](#) | [Last](#)

#2706173999	denjenigen vernünftig reden, die unseren	Begriff	von Öffentlichkeit nicht teilen? Nachvollziehen
#2709650368	Leistungen sind aber nur durchschnittlich. Der	Begriff	Zeitgeschichte ist aus dem Informationsinteresse
#2271543246	wiederm hinterließ dem 18. Jahrhundert den	Begriff	einer dynamischen, offenen Geschichte.
#2100412987	ausgewählte Forum wird dann nach den eingegebenen	Begriffen	durchsucht. Gib einfach den Benutzernamen
#645973149	Zumindest gegen Ihr Verständnis dieses	Begriffs	. Ein sportlich gehaltener Kinderrodel in
#2719021251	Projektbeschreibung taucht immer wieder der	Begriff	Modul auf. Ein Modul ist letztlich ein
#1864856840	Beipackzettel finden Sie dies unter dem	Begriff	Gegenanzeigen. Auch weitere Hinweise zu
#843383439	nicht zusammenpassen. Und entfernte den	Begriff	?gesund? von der Verpackung. Im Internet
#1488077412	Kommunikationspartner einander sind, desto mehr werden auch	Begriffe	aus dem umgangssprachlichen Bereich verwendet
#284107654	Expedition@bo.dr.s.de bestellt werden. Der	Begriff	Narzissmus ist so vieldeutig, dass ihm
#1453968627	Text::Abbrev Erzeugt aus einer Liste von	Begriffen	eine Abkürzungstabelle. Text::Balanced
#1055825351	Friedrichstraße liefern. Präsenzpflicht Der	Begriff	"City-Filiale" bürgerte sich ein. Hier
#1545778620	programmatische Vorgabe für die nächsten Jahre. Der	Begriff	ist ehrlich. Denn Talent bedeutet nicht
#626917325	eines tragenden Bauteils nach § 24 . Der	Begriff	"raumabschließend" ist somit nicht im landläufigen
#998849916	Tschernobyl hat die tödliche Dimension des	Begriffs	Restrisiko spüren lassen. Nicht Überheblichkeit
#2484747439	spirits" nur am Rande auftaucht und der	Begriff	behavioral finance" nicht einmal Erwähnung
#1865327519	Kulturforum: der Melancholie. Obwohl als	Begriff	durchaus gebräuchlich, dürften die meisten
#1066195727	Eine Firmware? Ein ROM? Ich hoere immer die	Begriffe	Firmware und ROM, wie soll man sich das
#931180247	nichts Ernstes. Nie, sagt sie, würde ihr der	Begriff	"Provinz" über die Lippen kommen. Nicht
#151488829	einer einfachen Suche. Dokumente, welche den	Begriff	'keine' enthalten, werden unverändert gewichtet
#1315067624	meint: (17.2.2008 um 23:29) Antworten Der "	Begriff	WLAN-Kabel" ist ein Oxymoron, wie ich neulich
#1896232155	Weder trifft es bei mir zu noch ist der	Begriff	"Bumser" in meinen Augen sozialadäquat.
#2635746810	Vollends absurd wird eine Orientierung an	Begriffen	, wenn die benutzten Begriffe mit den abgebildeten
#2202616937	Erkrankung ist der medizinisch korrekte	Begriff	, unter dem über 900 derzeit bekannte Krankheitsformen
#169557921	von Absatz 1 dient der Legaldefinition des	Begriffs	"internationale Koproduktionen ". Zu §
#1536755453	Zugangsdaten und Standortdaten. Bezüglich des	Begriffes	"Zugangsdaten" wird auf § 92 Abs. 3 Z 4
#558406229	geforderten, betrieblichen Altersversorgung. Der	Begriff	steht für eine Lebensversicherung, die
#176324671	Entwicklungsumgebungen und Werkzeuge nimmt der	Begriff	der Produktivitätssteigerung eine zentrale
#112512653	eine Krümmung dieser Raumzeit zurück. Der	Begriff	der Gravitationskraft wird dabei ersetzt
#90941217	allerdings nur, wenn man einen sehr weiten	Begriff	von Techno heranzieht. Noch stärker als
#2761581102	neue Bedeutung der Biologie z.B. durch den	Begriff	der Bionik einer größeren Öffentlichkeit
#2654703085	diese Weise sollte es möglich sein, den	Begriff	der Regelstudienzeit ad absurdum zu führen
#2598941734	natürliche oder künstliche Allergene aus der	Begriff	der Umwelt ist geprägt durch die anthropogene
#1001387296	Leben der Menschen und die Welt durch den	Begriff	der Gerechtigkeit zurückgebunden an Gott
#74499997	ich die Politik. Nehmen wir den Marxschen	Begriff	: hier ist Politik die Herrschaft des Menschen
#1040568325	zwischenmenschliche Kommunikation mit den	Begriffen	aus der Signalübertragung beschrieben wird
#2074189024	Diskussion so bedeutsamen aber auch schwierigen	Begriff	des Habitus zu erläutern und in seinen

Corpus: **DE_1**
Hits: **86** (37.3 per million) [[EU_1](#)] [[EU_2](#)]

Page 1 of 5 [Next](#) | [Last](#)

#95733	war frech , davon machen Sie sich keinen	Begriff	. Denken Sie sich , ich hab meinen Überzieher
	EU_2: Zoragarria zen ... Maiatza ... Bero egiten zuen ... Eta ni lotsagabe sentitzen nintzen , ez dakizu ondo .		
#110678	Gedächtnis , weil ich gerade zu dieser Zeit im	Begriffe	war , die Reisebeschreibung des berühmten
	EU_2: Oromenean geratu zitaldan ; izan ere , hain zuzen garai hartan , izen bereko ibiltari ezagunaren bidaiei buruzko kontakizunak irakurtzekotan nintzen .		
#136122	sonnverbrannte Wüste geht . Und schon war Andreas im	Begriffe	, den Helden des Films sympathisch und
	EU_2: Eta Andreas hasia zenean heroiari sinpatia eta atxikimendua hartzen , filmea uste ez bezala zorian aldera makurtu , handik pasatzen zen zientifiko espedizio batek desertuan galdutako gizona salbatu , eta Europako zibilizazioaren gerizpera ekartzten zuen atzera .		
#136185	den Helden des Films . Und schon war er im	Begriff	, sich zu erheben , als auf der Leinwand
	EU_2: Eta jeikitzeko zorian zegoela , hara non pantaHan agertu zen lehentxeago , barran zegoela , tabemariaren bizkar atzean ikusitako eskola lagunaren irudia .		
#139853	Hier stieß er auf einen Mann , der eben im	Begriffe	war , die Treppe hinaufzusteigen , und
	EU_2: Han , eskalieretan gora hastear zen gizon batekin topo egin zuen . Oso ezaguna egin zitzaion .		
#151414	aufgewachsen und hatte von dorthier eine Menge von	Begriffen	und Schablonen beibehalten . Er hatte theoretisch
	EU_2: Gainera hezimolde burges txikian hazia zen , eta orduanik kontzeptu eta kate andana gordetzen zituen .		
#154557	den Dichtungen des alten Indien ist dieser	Begriff	ganz unbekannt , die Helden der indischen
	EU_2: India zaharrek poesietan kontzeptu hori erabat ezezaguna da , Indiako epopeietako heroek ez dira pertsonak , pertsona matatak , enkamazio sailak baizik .		
#155264	Leben . Was die Menschen jeweils unter dem	Begriff	" Mensch " verstehen , ist stets nur eine
	EU_2: Gizonek " gizaki " kontzeptuaren bidez gehienetan ulertzen dutena , akordio burges iragankor bat besterik ez da betiere .		
#160170	und erkannte mich erst , als ich schon im	Begriff	war , an ihm vorüberzugehen . Er stürzte
	EU_2: Parez pare zetorkidan jakituna , zurrun eta nahiko begilauso , eta jaramonik egiteke albotik pasatzea pentsatzen nuen unean ezagutu ninduen .		
#202482	Ihre Arbeit zu Ende zu führen . Wenn der	Begriff	der Pflicht Ihnen unbekannt ist ... " Er
	EU_2: Egitekoaren ideia ezezaguna baizatu		
#202528	Warten Sie ! " sagte Gustav höflich . " Den	Begriff	war , an ihm vorüberzugehen . Er stürzte
	EU_2: -Egitekoaren ideiaa egiaz ez dut ezagutzen , ez gehiago .		
#202602	Gegenteil getan . Aber wenn ich auch den	Begriff	der Pflicht nicht mehr kenne , so kenne
	EU_2: Baina eginbeharraren ideia gehiago ezagutzen ez badut ere , ongi ezagutzen dut ordea kulparena ; agian biak gauza bera izango dira .		
#203637	Menschen , einst ein hohes Ideal , ist im	Begriff	, zu einem Klischee zu werden . Wir Verrückten
	EU_2: Gizaindia , garai batean ideal altu bat , azken batean klitxe bihurtzeko zorian dago .		
#204784	lückenhafte Seelenlehre der Wissenschaft durch den	Begriff	, den wir Aufbaukunst nennen . Wir zeigen
	EU_2: Guk horregatik , akatsez beteriko jakintza horen animari buruzko irakaskintza osatu egiten dugu , erakuntz deritzogu n ideiaen bidez .		
#392066	der Herr Vorsteher schwieg . Ich war im	Begriff	, ihm irgendeine grobe Beleidigung ins
	EU_2: Irain satsuren bat aurpegira botatzear egon nintzen .		

3.3 XML element hierarchy

```
<entry>
  <sense>
    <gramGrp>
    </gramGrp>
  </sense>
</entry>
```

- ▼ Bilingual entry (TEI)
 - ▼ syntactical information below WSD
 - ▼ TEI Standard
 - ▼ most common
 - ▼ redundancies kept on the syntax level

```
<entry>
  <synt>
    <sense>
    </sense>
  </synt>
</entry>
```

- ▼ Bilingual entry (2)
 - ▼ WSD below syntactical disambiguation
 - ▼ Redundancies
 - ▼ Common in German
 - ▼ monolingual for learners (PONS GW DaF)
 - ▼ bilingual (Langenscheidt)
 - ▼ never different POS
 - ▼ Useful for Basque(?)
 - ▼ possible to group different POS below the same headword

```
<synset>
  <LexUnit>
  </LexUnit>
</synset>
```

- ▼ WordNet
- ▼ Lexical Headword below WSD

abkühlen sep.
 I *tr./h* ▷ *hoztu* ^{+du}
 II *intr./sn* ▷ *hoztu* ^{+da}
 III *rlx./h* ▷ *hoztu* ^{+da}

heldu
 I *Verb Transitiv* -*dio* ▶ 1 fassen; greifen; halten; festhalten 2 ergreifen; anfangen 3 (Tier) beißen; stechen
 II *Verb Transitiv* -*du* = *bidali* ▶ 1 schicken 2 • *nbi*, *he|arazi* jdm. zukommen lassen
 III *Verb Intransitiv* -*da* ▶ 1 ankommen; kommen; herkommen 2 kommen; geschehen 3 reifen
 IV *Adjektiv* ▶ 1 reif 2 erwachsen 3 vernünftig
 V *Substantiv* ▶ 1 (eines Tiers) Biss; Stich

bake

verb UK US /beɪk/

Definition



- [I or T] **to cook inside a cooker, without using added liquid or fat**
I made the icing while the cake was baking.
a baked potato
freshly baked bread
Bake at 180°C for about 20 minutes.
Bake for 5-7 minutes in a preheated oven.
a baking dish/tin/tray
- [I or T] **to make something such as earth or clay hard by heating it, usually in order to make bricks**
- [I] **INFORMAL to be or become very hot**
It's baking outside.
You'll bake in that fleece jacket!

(Definition of bake verb from the Cambridge Advanced Learner's Dictionary & Thesaurus © Cambridge University Press)

3.4 Verb entry

```
<Lemma id="19375" LemmaSign="befinden" Source="Derewo" DerewoRang="569" Frequency="9"
  FlexMorfTag="unregelmäßig, ohne ge">
  <Section id="211053" SectionNumber="1" POS="Verb Transitiv" GramInfo="+haben">
    <Sense id="211054" SenseNumber="1" SyntaxHinweis="jdn./etw. für +adj.">
      <TE id="211055" TE="-tzat eman"/>
      <TE id="289097" TE="iritzi"/>
      <MWE id="289101" MWE="etw. für gut ~" Translation="ontzat eman; -i oniritzi"/>
    </Sense>
  </Section>
  <Section id="289102" SectionNumber="2" POS="Verb Intransitiv" GramInfo="+haben">
    <Sense id="289103" SenseNumber="1">
      <TE id="289104" TE="erabaki"/>
      <TE id="289105" TE="erabakia hartu"/>
    </Sense>
  </Section>
  <Section id="289098" SectionNumber="3" POS="Verb +&apos;sich&apos;" GramInfo="+haben">
    <Sense id="289099" SenseNumber="1">
      <TE id="289106" TE="izan"/>
      <TE id="289100" TE="egon"/>
      <TE id="289107" TE="aurkitu"/>
    </Sense>
  </Section>
</Lemma>
```

befinden Derewo 569 - *unregelmäßig, ohne ge.*

- I *Verb Transitiv* +haben ▶ 1 [jdn./etw. für +adj.] -tzat eman; iritzi
 - etw. für gut befinden ontzat eman; -i oniritzi
- II *Verb Intransitiv* +haben ▶ 1 erabaki; erabakia hartu
- III *Verb + 'sich'* +haben ▶ 1 izan; egon; aurkitu

3.5 Noun entry

Ausdruck¹ Derewo 1918 ~s, Ausdrücke

Substantiv *m. s.* *Verbinf.* [ausdrücken](#) ▶ 1 adierazpen 2 termino; esamolde; esapide
3 ezaugarri; adierazgarri 4 adierazkortasun; adierazgarritasun

Ausdruck² Derewo 1918 ~s, Ausdrücke

Substantiv *m. s.* *Verbinf.* [ausdrücken](#) ▶ 1 inprimatze; inprimaketa

3.6 Adjective and Adverb

anständig Derewo 8054

I *Adjektiv* ▶ 1 zintzo

II *Adjektiv adverbial* ▶ 1 zintzo; zintzoki; jendetasunez

anteilig Derewo 20466

I *Adjektiv* ▶ 1 zatikako; proportzional

II *Adjektiv adverbial* ▶ 1 zatika; proportzionalki

arg Derewo 3262

I *Adjektiv* ▶ 1 latz; gaitz 2 *veraltend* gaizto; maltzur

II *Adjektiv adverbial* ▶ 1 *ugs. südd.* oso

anders Derewo 596

I *Adverb* ▶ 1 ezberdin; bestela; beste era batean

II *Adverb präd./att.* ▶ 1 ezberdin; bestelako; beste era bateko

4.0 Automatic TE pairing for Bilingual Dictionary Drafting

- ▼ Corpus based
 - ▼ 1. GIZA++
 - ▼ 2. Bifid
 - ▼ 3. Elhuyar Pivot Dictionaries
- ▼ Others
 - ▼ 4. Wikimedia
 - ▼ 5. WordNet

4.1 Corpus-based strategies: GIZA++

▼ Bilingual Parallel Corpus, GIZA++

- ▼ English-Swedish 700.000 sentence pairs, 82,3% precision, 73,3& recall (Holmqvist 2010)
- ▼ Lithuanian 1.765.000 tokens, Hungarian 2.121.000 tokens, 4026 TE-candidates (Héja 2010)
- ▼ DE-EU:
 - ▼ DE: 2.303.307 tokens,
 - ▼ EU: 1.948.504 tokens,
 - ▼ Extremely low recall:
 - ▼ 573 TE-candidates,
 - ▼ 153 after revision
 - ▼ Low precision
 - ▼ First attempt, Gorka Labaka, UPV-EHU
 - ▼ Second attempt: After lemmatizing EU corpus

abend	<i>gauean</i>	aber	<i>ordea</i>
abend	<i>gau</i>	aber	<i>-baina</i>
abend	<i>arrats</i>	abgesehen	<i>utzita</i>
abend	<i>arratsean</i>	abgesehen	<i>aparte</i>
abend	<i>iluntzean</i>	abgrund	<i>amildegira</i>
abend	<i>arratsaldean</i>	abgrund	<i>amildegia</i>
abend	<i>gauerako</i>	abgrund	<i>amildegia</i>
abend	<i>arratsalde</i>	abgrund	<i>amildegi</i>
abend	<i>gauetz</i>	abholen	<i>bila</i>
abendessen	<i>afaria</i>	abrupt	<i>-batean</i>
abendessen	<i>afaltzeko</i>	abschied	<i>agur</i>
abends	<i>gauean</i>	abschied	<i>adio</i>
abends	<i>iluntzean</i>	absicht	<i>asmoa</i>
abends	<i>gauetz</i>	absicht	<i>nahita</i>
abends	<i>arratsean</i>	absicht	<i>asmo</i>
abends	<i>gauetan</i>	absicht	<i>asmorik</i>
abends	<i>arratsalde</i>	absichten	<i>asmo</i>
abenteuer	<i>abentura</i>	absichtlich	<i>nahita</i>
abenteuer	<i>abenturak</i>	absichtlich	<i>propio</i>
aber	<i>baina</i>	abstand	<i>tartea</i>

4.2 Corpus based strategies: Bifid

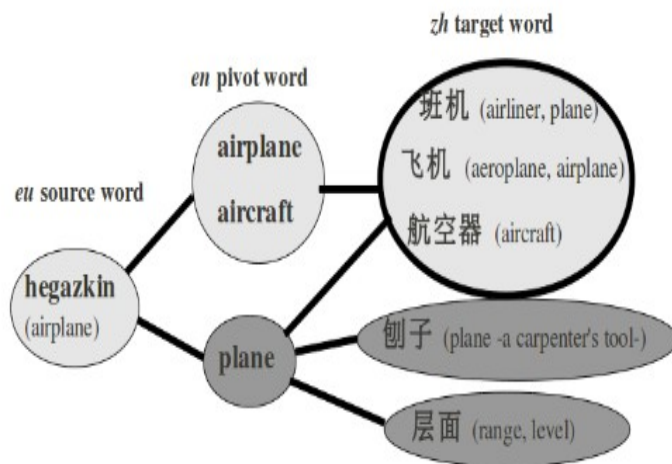
▼ Bilingual Corpus, **Bifid** (cf. Nazar 2012)

- ▼ First attempt done: 2089 pairings,
- ▼ Works much worse than with other language-pairs, evaluation in course
- ▼ Insufficient corpus size
- ▼ Second attempt: After lemmatizing EU corpus

wills gesicht	<i>willen aurpegia</i>
bildhauer	<i>eskultorea</i>
einem tiefen	<i>hasperen sakon</i>
individuum	<i>gizabanako</i>
eines hundes	<i>txakur baten</i>
mensch ärgere	<i>partxisean</i>
rechts links	<i>eskuin ezker</i>
schranktür	<i>armairuko</i>
sagte irene	<i>irenek</i>
wolfgang	<i>amadeo</i>
fünfunddreißig	<i>hogeita hamabost</i>
fußgänger	<i>oinezko</i>
tante	<i>izeba</i>

4.3 Corpus-based strategies: Elhuyar Pivot Dictionaries

- ▶ Pairing by combining pivot and corpus based methods
 - ▶ **Research in course** at Elhuyar (Saralegi et al., 2012): Pivot-based bilingual dictionary building
 - ▶ Pairings obtained by two reference dictionaries with English as a pivot
 - ▶ Inverse Consultation (IC) and Distributional Similarity (DS) in bilingual comparable corpora
 - ▶ Problems resolving polysemy (words with ambiguities are skipped)
 - ▶ Done for Basque and Chinese, German, Hindi, Swahili, Arabic
 - ▶ German-Basque: In the actual version, 13.000 DE lemmata with 1-25 basque TE each
 - ▶ Evaluation in course. Still noisy.



eu | es | en | de | hi | sw | ar | zh
Babeslea:

e
hiztegien
ataria

PIBOTAJE BIDEZ SORTUTAKO HIZTEGI AUTOMATIKOAK

EUSKO JAURLARITZA
GOBIERNO VASCO

Sarrera | Laguntza | Deskargak | Eztabaidagunea |

Hiztegi Kontsulta

Alemana → Euskara

Begriff

Bilatu

Emaiza: **1** sarrera topatu dira hiztegian.

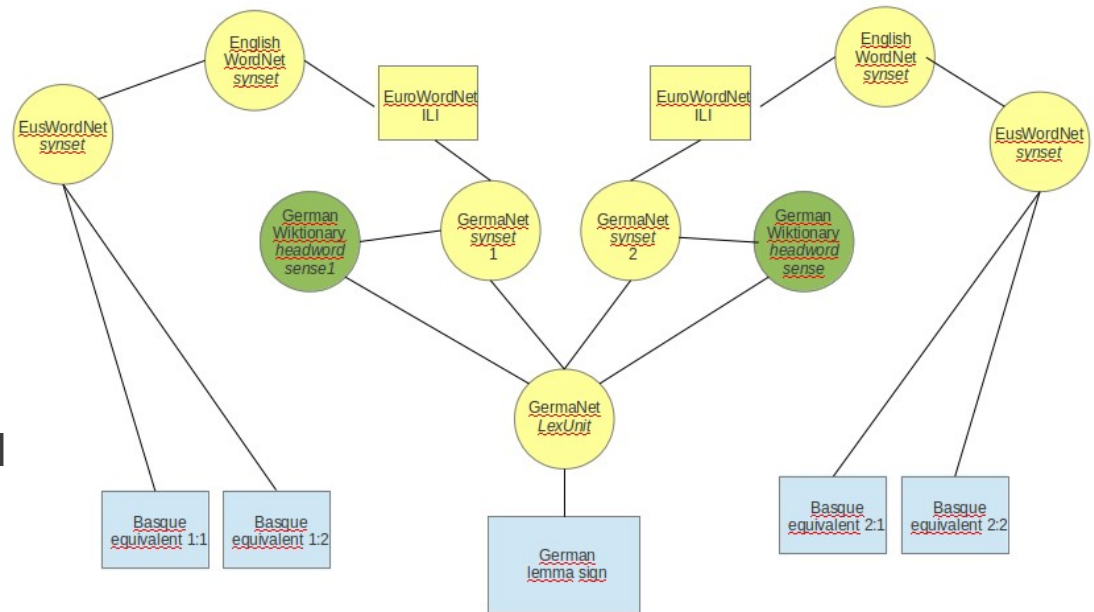
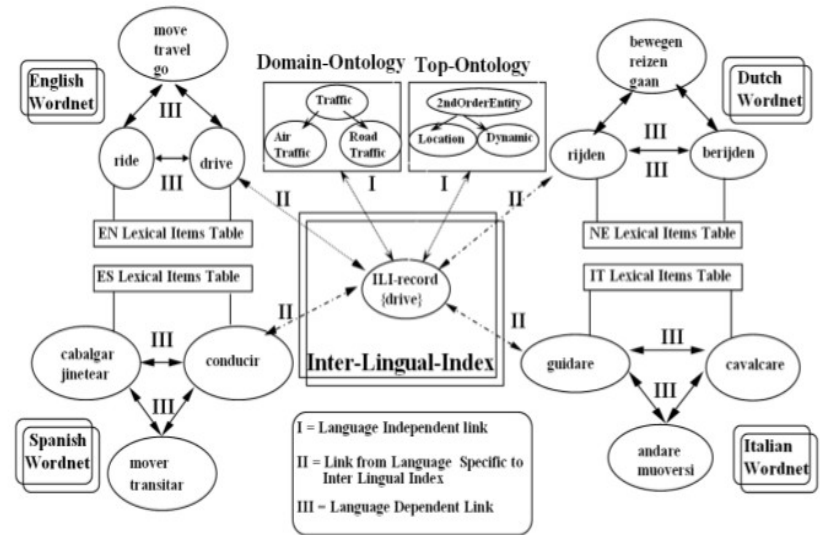
Topatutako sarrerak	Ordainak
<p>Begriff</p> <ul style="list-style-type: none"> ▪ Und konnte das sinnlose Wüten in jener Nacht etwas anderes gewesen sein als das Werk eines Kranken? Psychiatrische Begriffe wie „Psychose“ und „schizophrener Schub“ erschienen den D. s plötzlich wie Worte des Trostes. ▪ Ein romantisierender Begriff, ein Anachronismus in Zeiten der Globalisierung, wo jeder überall auf der Welt zu Hause sein kann und es scheinbar auch ist? Oder kann es sein, dass der Mensch ohne Verbundenheitsgefühl mit dem Ort, an dem er aufwuchs, in eine Sinnkrise stürzt, sich selbst verliert? Mit dem Zusammenspiel von Heimat und Identität setzt sich die 7. Internationale Foto-Triennale in Esslingen künstlerisch auseinander. ▪ Lenz: Es gibt anfangs ein paar Begriffe und Stichworte, um mich später daran zu erinnern. <p>Erakutsi adibide gehiago</p>	<p style="text-align: center;">1 itzulpen ziur eta 0 zalantzako aurkitu dira</p> <p>kontzeptu +1 -0 </p> <ul style="list-style-type: none"> ▪ Poesiarekin lotura handia duten kontzeptuak dira nolnahi ere, oso poetikoa baita zientifikoa ere badela deskubritu aurrekoa- txinatarren aspaldiko esaera: tximeleta baten hegoen dardara munduaren beste aldean sentitu liteke. ▪ Esan bezala, interesatzen zaizkion artisten lanak erreproduzitzen edo berinterpretatzen ditu, interes filosofiko edo kontzeptualagatik baino gehiago (autoretza edo originaltasuna zalantzan jartzea, esaterako), ezaugari formal eta estetikoek bultzatuta. ▪ Gizartearekin konpromisoa dugun GKEontzat, ordea, garapenerako laguntza-eredu ofiziala hermitasun globalaren kontzeptuan egon behar da oinarrituta, hau da, globalizazioaren ikuspegi sozial bezain demokratikoan eta giza eskubideen nahiz garapenerako eskubidearen erabateko indaraldian. <p>Erakutsi adibide gehiago</p>

4.4 TE extraction from Wikimedia

- ▼ First attempts: TE extraction from Wikimedia database dumps (source: <http://www.dicts.info>)
 - ▼ Wiktionary (by 2012):
 - ▼ DE-EN 2161 pairings. DE 3066 LU (905 syn.), EN 2161 LU (0 syn.)
 - ▼ DE-EU 2131 pairings. DE 2316 LU (185 syn.), EU 2338 LU (207 syn.)
 - ▼ *Eval. by hand: 80%+ nouns, very freq. words and toponyms, nearly no wrong translations*
 - ▼ *after MWE filtering and synonym cut-off: 2030 pairings. 1330 yes, 700 not on our lemmalist.*
 - ▼ Omegawiki interlanguage links (by 2012):
 - ▼ DE-EN 5603 pairings. DE 6836 LU (1233 syn.) EN 6393 (790 syn.)
 - ▼ DE-EU 4603 pairings. DE 5362 LU (759 syn.), EU 6771 LU (2168 syn.)
 - ▼ *Eval. by hand: mostly rare terms, some wrong translations*
 - ▼ Wikipedia interlanguage links (wikipedia article titles) (by 2012):
 - ▼ DE-EN 7433 pairings
 - ▼ DE-EU: ?
 - ▼ *Eval. by hand: only nouns, frequent and rare terminology, no wrong translations*
 - ▼ Erdmann 2007: Much better results!
 - no significant difference between the two language pairs
- Wiktionary & Omegawiki (still) not useful for automatic dictionary drafting
- Wikipedia bilingual term extraction: More research for DE-EU required

4.5 TE pairing via EuroWordNet

- ▼ Stops on the way:
 - ▼ German lemma-sign
 - ▼ GermaNet Lexical Unit
 - ▼ GermaNet synset(s)
 - ▼ EuroWordNet ILI record(s)
 - ▼ Princeton WordNet synset(s)
 - ▼ EusWN synset(s)
 - ▼ EusWN lemma-signs
- ▼ Advantages:
 - ▼ large databases
 - ▼ link on sense level
 - ▼ human-made WSD
- ▼ Known Problems:
 - ▼ Different approaches interling. WN
 - ▼ Different versions of PWN



5. Conclusions

- ▼ Importance of corpus methods in Lexicography increases
 - ▼ Indispensable in dictionary editing
 - ▼ Upcoming trend in dictionary publishing
 - ▼ Automatic Dictionary Drafting
- ▼ Not without large bilingual corpora
- ▼ Human Lexicographers (still) needed...

... The main resource remains the lexicographic knowledge of the project members in combination with large annotated text corpora which serve as a reference in all cases of doubt.

(Introduction to German Wordnet)

Thanks for your attention
Vielen Dank, Eskerrik asko

david.lindemann.soraluze@gmail.com