Towards Automatic Linking of Lexicographic Data

The case of *Den Danske Ordbog* and *Ordbog over det danske Sprog* for the Danish language

Sina Ahmadi, Sanni Nimb, John P. McCrae, Nicolai H. Sørensen

XIX EURALEX International Congress Alexandroupolis, Greece September 2021





This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731015.



Outline

- 1. Context
- 2. Danish Lexicographical Data
- 3. Word-Sense Alignment
- 4. Experiments with ODS and DDO
- 5. Next steps



1. Context

"Dictionaries are treasure houses of data on the uses of words. They are also our best starting point for all questions regarding word sense distinctions, in NLP, the humanities or lexicography. But to reveal the dictionary's treasures in a systematic way is no simple."

— Adam Kilgarriff (*Dictionary Word Sense Distinctions: An Enquiry into Their Nature*)



Some of the resources of Danish in the Society for Danish Language and Literature (DSL), Copenhagen

european lexicographic infrastructure

1. Context Lexical resources



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731015.

1. Context

ELEXIS – European Lexicographic Infrastructure

17 partners from 15 countries, 73 institutions from 37 countries (February 2018 - July 2022)



Find out more about ELEXIS at <u>https://elex.is/</u>

european lexicographic



1. Context Why linking resources?

- To improve word and concept coverage
 - e.g., named entities, new senses
- To improve domain coverage
- To improve multilingualism
- Creating resources for new languages
- To combine expert-made semantic relations
 - e.g., Hypernymy, meronymy, etc.



BabelNet (https://babelnet.org/)



Context A few applications

- Entity linking
- Word-sense disambiguation
- Semantic role labeling

"Paris is the capital of France"

wikipedia.org/wiki/Paris wikipedia.org/wiki/France

Source: https://en.wikipedia.org/wiki/Entity_linking





1. Context A few applications

- Entity linking
- Semantic role labeling
- Word-sense disambiguation



Source: https://web.stanford.edu/~jurafsky/slp3/slides/22_SRL.pdf

european lexicographic infrastructure

1. Context A few applications

- Entity linking
- Semantic role labeling
- Word-sense disambiguation

WordNet Search - 3.1 - <u>WordNet home page</u> - <u>Glossary</u> - <u>Help</u>

Word to search for: tie Search WordNet

Display Options: (Select option to change)

Key: "S:" = Show Synset (semantic) relations, "W:" = Show Word (lexical) relations Display options for sense: (gloss) "an example sentence"

Noun

- <u>S: (n) necktie</u>, **tie** (neckwear consisting of a long narrow piece of material worn (mostly by men) under a collar and tied in knot at the front) *"he stood in front of the mirror tightening his necktie"; "he wore a vest and tie"*
- <u>S:</u> (n) <u>affiliation</u>, <u>association</u>, <u>tie</u>, <u>tie-up</u> (a social or business relationship)</u> "a valuable financial affiliation"; "he was sorry he had to sever his ties with other members of the team"; "many close associations with England"
- <u>S:</u> (n) tie (equality of score in a contest)
- <u>S:</u> (n) tie, <u>tie beam</u> (a horizontal beam used to prevent two other structural members from spreading apart or separating) *"he nailed the rafters together with a tie beam"*
- <u>S: (n) link, linkup, tie, tie-in</u> (a fastener that serves to join or connect) "the walls are held together with metal links placed in the wet mortar during construction"
- <u>S:</u> (n) <u>draw</u>, <u>standoff</u>, **tie** (the finish of a contest in which the score is tied and the winner is undecided) *"the game ended in a draw"; "their record was 3 wins, 6 losses and a tie"*
- <u>S:</u> (n) tie ((music) a slur over two notes of the same pitch; indicates that the note is to be sustained for their combined time value)
- <u>S:</u> (n) tie, <u>railroad tie</u>, <u>crosstie</u>, <u>sleeper</u> (one of the cross braces that support the rails on a railway track) *"the British call a railroad tie a sleeper"*
- <u>S:</u> (n) tie (a cord (or string or ribbon or wire etc.) with which something is tied) *"he needed a tie for the packages"*



2. Danish lexicographical resources

- Lexicographic resources at DSL
 - DanNet the Danish WordNet
 - Ordbog over det danske Sprog (ODS) (open access: https://ordnet.dk/ods)
 - Den Danske Ordbog (DDO) (open access: https://ordnet.dk/ddo)
 - The Danish Thesaurus
 - Danish FrameNet (aligned with DDO)
- Our focus is on on linking ODS and DDO at sense level





2. Danish lexicographical resources ODS vs. DDO

	afstand _{substantiv} , fælleskøn Vis overblik	Vis fork	Afstand, en. ['au,sdan'] flte ['au-,sdanə] (efter ty. abstand (if. S lat. distantia); Moth(S740) har ordet som vbs. til afstaa: "fravigelse.
	BØJNING -en, -e, -ene UDTALE ['aw,sdan'] & @ OPRINDELSE efter tysk Abstand Betydninger		I) fjernhed; længden af mellemrummet (mat.: af en ret linie) mellem to punkter. udfinde Solens Afstand fra Jorden. Heitm. Physik.67. (jf. Steners. CritBet.29 og Marg. Klopstock.Breve.(1760).40). Søfarende have
1	 rumlig udstrækning der adskiller SYNONYMER distance sjældent fra ORD I NÆRHEDEN hul¹ gab² spri GRAMMATIK afstand mellem NOGET/I EKSEMPLER indbyrdes afstand e behørig afstand e i/på passende Vi passerede skibene i en afstand 	to punkter, linjer eller flader, målt som længden af en linje eller rute stand ng plads tom plads ledig pladsvis mere NOGEN og NOGET/NOGEN afstand af/på MÅL afstand fra/til NOGET den geografiske afstand bedømme/måle afstanden på lang afsta afstand i/på sikker afstand af ca. 40 meter HvidovA1989	mellem dem ofte stor Færdighed i at bedømme Afstandene. Heib.Pros.II.369. *Seer jeg en Hatfuld sydet Damp Tidens Maal, Rummets Afstande flytte. Ploug.VV.II.13. Afstanden fra Kærsholm til Bøstrup Præstegaard var fem- seks Kilometer. Pont.LP. VII.65. (sj.:) han vandt Afstand (dvs.: kom længere og længere bort) fra (de angribende vilde heste). Rist.FT.28. × spec. om regelmæssige mellemrum mellem (afdelinger af) soldater, som staar bag ved hinanden. Sal. IX.531. jf.: Der er Gæssene med Retning og Afstand
	Brevduer kan finde hjem over lar 1.a tidsmæssig udstrækning der ORD I NÆRHEDEN tidsrum tids han vidste at såret til trods f	ge afstande skoleb-fys.91 adskiller to begivenheder interval interval tidsafstand tidsmarginvis mere or de tyve års afstand endnu ikke var lægt SvÅMad89	som en Trup Soldater. Bogan.I.127. det er daarlig ridning; her er ikke spor af afstand (dvs.: der er ulige stor afstand mellem de enkelte ryttere) efter præp. i. *Alt, hvad Naturen I maalløs Afstand fra hinanden spreder. Bagges.L. I.156. Munken holder sig i en ærbødig Afstand. Oehl.IV.161. Medens vi talte, saa jeg i lang Afstand en Dame komme. Goldschm.VI.274. i
	1.b OVERFØRT mangel på fortrolighed, kontakt eller personligt engagement SYNONYM distance ORD I NÆRHEDEN modsætningsforhold delte meninger påstand mod påstand en strid om ord uforenelighedvis mere GRAMMATIK især i singularis EKSEMPLER en vis afstand 🖻		afstand (ell. T i en afstand. Gylb.III.215. IV.280.331. VIII.223), (sj.) S ikke tæt ved ell. paa nært hold; (temmelig) langt borte. (nu oftere paa afstand). *Jeg vendte om, og som en ydmyg Slave I Afstand troe jeg fulgte Deres Vei. Heib. Poet.VII.272. Vandet saae klart ud nær ved, men seet i Afstand, sort som Blæk. HCAnd.VI.300. Jeg elsker Franskmændene – i Afstand. Goldschm.I.354. De dræbte ham i Afstand med Pile. smst.III. 197.
	brugen af fagsprog skaber er 1.c OVERFØRT forskel eller modsæ ORD I NÆRHEDEN forskellighed	afstand mellem læge og patient fagb-psyk.92a tningsforhold mellem to parter eller størrelser divergens skisma kløft grøft gab²vis mere	Image: Præp. paa. *Ei een af Huys Mænd er under Vaaben Paa mange S Miles Afstand. Hauch Æ.70. Der skulde skydes (dvs.: ved en duel) paa 15 S Skridts Afstand. JakKnu A.211. paa afstand, d. s. s. i afstand. Hun havde ogsaa noget Godt i sine Øjne, naar hun var lidt paa Afstand. Schand.BS.118.
https://ordnet.dk/ddo		n og oppositionen er mindsket BT1991 og realiteter kalder jeg svindel BerlT1991	Aftenklokken lød saa smukt paa Afstand. <i>Pont.LP</i> . der har den mindste Katar bør holde sig paa Afstand <i>Sundhedstid.1916.266</i> .

DEN DANSKE ORDBOG



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731015.



2. Danish lexicographical resources ODS vs. DDO

ODS

- Historical Danish dictionary (1700-1955)
- Richly-described entries
- Retrodigitized → noisy senses (both, in printed and XML versions)
- Senses not linked to any external resource
- Obsolete senses





2. Danish lexicographical resources *afstand* ['aw,sdan[?]] (distance in Danish)*

1) fiernhed; længden af mellemrummet (mat.: af en ret linie) mellem to punkter. udfinde Solens Afstand fra Jorden. Heitm. Physik.67. (jf. Steners. CritBet.29 og Marg. Klopstock.Breve.(1760).40). Søfarende . . have ofte stor Færdighed i at bedømme Afstandene. Heib.Pros.II.369. *Seer jeg... en Hatfuld sydet Damp | Tidens Maal, Rummets Afstande flytte. Ploug.VV.II.13. Afstanden fra Kærsholm til Bøstrup Præstegaard var femseks Kilometer. Pont.LP. VII.65. (sj.:) han vandt Afstand (dvs.: kom længere og længere bort) fra (de angribende vilde heste). Rist.FT.28. || 🗙 spec. om regelmæssige mellemrum mellem (afdelinger af) soldater, som staar bag ved hinanden. Sal. IX.531. jf .: Der er Gæssene . . med Retning og Afstand som en Trup Soldater. Bogan. I.127. det er daarlig ridning; her er ikke spor af afstand (dvs.: der er ulige stor afstand mellem de enkelte ryttere) S efter præp. i. *Alt, hvad Naturen . . | I maalløs Afstand fra hinanden spreder. Bagges.L. I.156. Munken . . holder sig i en ærbødig Afstand. Oehl.IV.161. Medens vi talte, saa jeg i lang Afstand en Dame komme. Goldschm.VI.274. i S afstand (ell. [†] i en afstand. Gylb.III.215. IV.280.331. VIII.223), (sj.) ikke tæt ved ell. paa nært hold; (temmelig) langt borte. (nu oftere paa afstand). *Jeg vendte om, og som en ydmyg Slave | I Afstand troe jeg fulgte Deres Vei. Heib. Poet.VII.272. Vandet saae klart ud nær ved, men seet i Afstand, sort som Blæk. HCAnd.VI.300. Jeg elsker Franskmændene - i Afstand. Goldschm.I.354. De dræbte ham i Afstand med Pile. smst.III. 197. S | efter præp. paa. *Ei een af Tillys Mænd er under Vaaben | Paa mange S Miles Afstand. Hauch.Æ.70. Der skulde skydes (dvs.: ved en duel) paa 15 Skridts Afstand. JakKnu.A.211. paa afstand, d. s. s. i afstand. Hun havde ogsaa noget Godt i sine Øjne, naar hun var lidt paa Afstand. Schand.BS.118. Aftenklokken . . lød saa smukt paa Afstand. Pont.LP.VII.89. Enhver Voksen der har den mindste Katar bør holde sig paa Afstand fra Børnene. Sundhedstid.1916.266.

* https://ordnet.dk/ddo/ordbog?query=afstand

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731015.

SemID="29014340"><Semem BetNo="1" SememID="29401382" odsID="Afstand 1"><dotPlus id="2666"/><SemIndhold> o>1)</dotBetNo> <i><dotSpaced>fjernhed;</dotSpaced> længden af <dotSp ced>mellemrummet</d</pre> col="0293" lin="25" orig="mellemrum-met"/>(mat.: af en ret linie) mellem to punkter.</i> <dotLn col="0293" lin=" 26"/>udfinde Solens Afstand fra Jorden. <i>Heitm. <dotLn col="0293" lin="27"/>Physik.67. (jf. Steners. CritBet.29 og Marg. <dotLn col="0293" lin="28"/>Klopstock.Breve.(1760).40).</i> Søfarende . . <dotLn col="0293" lin="29"/>have ofte stor Færdighed i at bedømme <dotLn col="0293" lin="30"/>Afstandene. <i>Heib.Pros.II.369.</i> d>Seer jeg <dotLn col="0293" lin="31"/>. . en Hatfuld sydet Damp | Tidens Maal, <dotL d>*</do col="0293" lin="32"/>Rummets Afstande flytte. <i>Ploug.VV.II.13.</i> <dotLn col="0293" lin="33"/>Afstanden fra Kærsholm til Bøstrup Præstegaard<dotLn col="0293" lin="34" orig="Præste-gaard"/>var fem-seks Kilometer. <i> Pont.LP. <dotLn col="0293" lin="35"/>VII.65. (sj.:)</i> han vandt Afstand <i>(dvs.: kom <dotLn col="0293" lin=" 36"/>længere og længere bort)</i> fra <i>(de angribende <dotLn col="0293" lin="37"/>vilde heste). Rist.FT.28.</ > <Planke PlankeID="29900969"><dotEdBreak/> <piktogram src="swords" title="sværd"/> <i>spec. om regelmæssige< otLn col="0293" lin="38" orig="regel-mæssige"/>mellemrum mellem (afdelinger af) <dotLn col="0293" lin="39"/> soldater, som staar bag ved hinanden. Sal. <dotLn col="0293" lin="40"/>IX.531. jf.:</i> Der er Gæssene . . med Retning<dotLn col="0293" lin="41" orig="Ret-ning"/>og Afstand som en Trup Soldater. <dotLn col="0293" lin="42"/> <i>Bogan.I.127.</i> det er daarlig ridning; her <dotLn col="0293" lin="43"/>er ikke spor af afstand <i>(dvs.: der er ulige stor <dotLn col="0293" lin="44"/>afstand mellem de enkelte ryttere)</i> <piktogram src="dotpipe" title="dobbeltbrudt streg"/> </Planke><Planke PlankeID="29900970"><dotPlus id="2667"/><dotE reak/> <i>eft<u>er <</u> lotLn col="0293" lin="45"/>præp.</i> i. <dotRaised>*</dotRaised>Alt, hvad Naturen . . | I maalløs <dotLn col=" 0293" lin="46"/>Afstand fra hinanden spreder. <i>Bagges.L. <dotLn col="0293" lin="47"/>I.156.</i> Munken . . holder sig i en ærbødig <dotLn col="0293" lin="48"/>Afstand. <i>0ehl.IV.161.</i> Medens vi talte, saa <dotLn col ="0293" lin="49"/>jeg i lang Afstand en Dame komme. <i>Goldschm.VI.274.</i><dotLn col="0293" lin="50" orig=" Gold-schm.VI.274."/>i afstand <i>(ell.</i> <piktogram src="cross" title="kors"/> i en afstand. <dotLn col ="0293" lin="51"/><i>Gylb.III.215. IV.280.331. VIII.223), (sj.) <dotLn col="0293" lin="52"/><dotPlus col="293" id="2668" lin="50"/><do d>ikke tæt ved</dotSpaced> ell. paa nært hold; (temmelig) <dotLn col="0293" lin="53</pre> "/>langt borte. (nu oftere</i> paa afstand<i>).</i> <dotRaised>*</dotR d>Jea <dotLn col="0293" lin="54"/>

vendte om, og som en ydmyg Slave | I <dotLn col="0293" httn col="0293" lin="56"/>Poet.VII.272.</i> Vandet sa i Afstand, sort som Blæk. <dotLn col="0293" lin="58"/> 0293" lin="59" orig="Franskmæn-dene"/>- i Afstand. <i>> ham i Afstand med Pile. <i>smst.III. <dotLn col="0293 <dotPlus id="2669"/><dotPlus id="2670" kontrol="inger</pre> <i>efter præp.</i> paa. <dotRaised>*</d tRaised>Ei Vaaben | Paa mange <dotLn col="0293" lin="63"/>Miles lin="64"/>skydes <i>(dvs.: ved en duel)</i> paa 15 Sł <i>JakKnu.A.211.</i> paa afstand, <i>d. s. s.< noget Godt <dotLn col="0294" lin="01"/>i sine Øjne, na >Schand.BS.118.</i> Aftenklokken . . lød saa <dotl</pre> </i> Enhver<dotLn col="0294" lin="04" orig="En-hver"/ '/>bør holde sig paa Afstand fra Børnene. <dotLn col='</pre> 0294" lin="07"/></Planke></SemIndhold></Semem></Se





2. Danish lexicographical resources ODS vs. DDO

DDO

- modern Danish dictionary 1955-->
- DSL's direct follower to ODS --> many lemmas and senses in common
- Well-structured sense descriptions
- ~50% linked to DanNet, ~90% to The Danish Thesaurus at sense level





european lexicographic infrastructure

3. Word sense alignment (WSA)

→ linking lexical content at sense level, including glosses

lead Video English: lead¹ English: lead² American: lead¹ American: lead² Specialist English: lea ▶ lead 1 \rightarrow the lead 16. countable noun 2 [singular] the amount or distance by which one competitor is ahead of another A lead is a piece of information or an idea which may help people to discover the Intersection of the section of th facts in a situation where many facts are not known, for example in the lead over investigation of a crime or in a scientific experiment. The Socialists now have a commanding lead over their opponents. The inquiry team is also following up possible leads after receiving 400 calls from the public. 3 [singular] if someone follows someone else's lead, they do the same as the other person has done ● Other countries are likely to follow the U.S.'s lead. Synonyms: clue, tip, suggestion, trace More Synonyms of lead ■ The Government should give industry a lead in tackling racism (=show what other people 17. countable noun should do). ■ The black population in the 1960s looked to Ali for a lead (=looked to him to show them what The lead in a play, film, or show is the most important part in it. The person who they should do). plays this part can also be called the lead. 4 \rightarrow take the lead (in doing something) Performers from the Bolshoi Ballet dance the leads. Both the leads in the play are impressive. 5 [countable] a piece of information that may help you to solve a crime or mystery SYN clue ■ The police have checked out dozens of leads, but have yet to find the killer. Synonyms: leading role, principal, protagonist, title role More Synonyms of lead 6 [countable] the main acting part in a play, film etc, or the main actor 18. countable noun play the lead/the lead role A dog's lead is a long, thin chain or piece of leather which you attach to the dog's ▲ He will play the lead role in 'Hamlet'. collar so that you can control the dog. Dowers was cast in the lead role (=he was chosen to play it). the male/female lead [mainly British] ▲ They were having trouble casting the female lead. An older man came out with a little dog on a lead. • the film's romantic lead **REGIONAL NOTE:** Longman in AM, use leash 7 \rightarrow lead singer/guitarist etc DICTIONAR

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731015.



Sense alignment is challenging

- Differences in structure
 - sense vs. sub-sense vs. sub-sub-sense etc.
 - formalizations, such as WordNet [8], FrameNet [9] and generative lexicon [10]
- Differences in content
 - Lexical choice:
 - Alcohol: <u>vandklar</u> vædske (water-clear liquid) vs. <u>farveløs</u> (colorless) in Danish
 - Definition paradigms [11], as in:
 - footnote: A footnote is a note of text placed at the bottom of a page (analytical)
 - méchant : qui est dangereux, nuisible, néfaste (synonymous)
 - good: the opposite of bad (relational)

3. Word sense alignment Monolingual WSA datasets

To address some of the current main limitations in WSA:

- Multilingualism
- Monolingual resources
- Gold-standard datasets
- Semantic relationship annotation



A Multilingual Evaluation Dataset for Monolingual Word Sense Alignment

Sina Ahmadi^{*}, John P. M^cCrae^{*},

Sami Nimb¹, Fahar Khan², Monica Monachin², Holette S. Pedersen², "Therry Declerck^{2,12}, Tanja Wisak², Andrea Belland¹, Ircene Fissi⁴, Thomas Torolgard², Siao Oren⁴, Simon Krek², Veronika Lipp⁶, Tamás Várrad⁶, László Simon⁶, András Cyöft⁵, Cardet Therias⁶, Tamake Schoonheim⁹, Yifat Ben Mohe¹⁰, Maya Ruich¹⁰, Raya Aba Manad¹¹, Doriele Lonke⁶, Kirs Kovarkach¹⁰, Margit Langemets¹¹, Jehen Kalla¹¹ Oksana Dereza¹, Theodorus Fransen⁷, David Cillessen⁷, David Linderman¹⁴, Mikel Alonso¹⁴, Ana Salgado¹⁴ José Luis Sanch¹⁶, Rathel¹, Urein-Ruick¹⁴, Ordf Porta Zameran¹⁰, Kirs I Simo¹⁰, Petry Osenora¹⁷, Zara Kancheva¹¹, Iraylo Rade¹¹, Ranka Stankovi¹⁴, Andrej Pertih¹⁰, Dejan Gahrovšet¹⁹ ¹Insight Catte for Data Analytics, National University of Tecland, Galvay ^(sin admatdi, john mecrea) ^{(sin sight-center org (other affiliations in Appendix A)}

Abstract

Algring energies across resources and languages is a challenging task with beneficial applications in the field of natural language processing and electronic lectocography. In this paper, we describe our efforts in manually algring monolingual dicionaries. The languages is carried out at semi-level for various resources in 15 languages. Morrower, senses are anotated with possible semantic relationships such of languages and resources and focuses on the nove challenging task of linking general-purpose language. We believe that or data will pave the way for further advances in alignment and evaluation of word senses by creating new solutions, particularly those nototionsby requiring data such as neural networks. Our resources are probletly vaniable bet the part / part. https://withs.

Keywords: lexical semantic resources, sense alignment, lexicography, language resource

Mihalcea, 2005; Ponzetto and Navigli, 2010; Gurevych et al., 2012).

1. Introduction Lexical semantic resources (LSRs) are knowledge repositories that provide the vocabulary of a language in a descriptive and structured way. One of the famous examples of LSRs are dictionaries. Dictionaries form an important foundation of numerous natural language processing (NLP) tasks, including word sense disambiguation, machine translation, question answering and automatic summarization. However, the task of combining dictionaries from different sources is difficult, especially for the case of mapping the senses of entries, which often differ significantly in granularity and coverage. Approaches so far have mostly only been evaluated on named entities and quite specific domain language. In order to support a shared task at the GLOB-ALEX workshop1, we have developed a new baseline that covers 15 languages and will provide a new baseline for the task of monolingual word sense alignment.

Different dictionaries and related resources such as worthnets and encyclopedia have significant differences in structure and heterogeneity in content, which makes aligning information across resources and languages a challenging task. Word sense alignment (WSA) is a more specific task of linking dictionary content at sense level which has been proved to be beneficial in various NLP tasks, such as wordsense disambiguiton (Navigil and Ponzetto, 2012), semantic role labeling (Palmer, 2009) and information extraction (More et al., 2013). Moreover, combining LSRs can enhance domain coverage in terms of the number of lexicalters and types of lexical-semantic information (Shi and

Given the current progress of artificial intelligence and the usage of data to train neural networks, annotated data with specific features play a crucial role to tackle data-driven challenges, particularly in NLP. In recent years, a few efforts have been made to create gold-standard dataset, i.e., a dataset of instances used for learning and fitting parameters for aligning senses across monolingual resources including collaboratively-curated ones such as Wikipedia2, and expert-made ones such as WordNet. However, the previous work is limited to a handful of languages and much of it is not on the core vocabulary of the language, but instead on named entities and specialist terminology. Moreover, de spite the huge endeavour of lexicographers to compile dictionaries, proper lexicographic data are rarely openly accessible to researchers. In addition many of the resources ar quite small and the extent to which the mapping is reliable is unclear.

In this paper, we present a set of datasets for the task of WSA containing manully-monotated monolingual regources in 15 languages. The annotation is carried out at sense level where four semantic relationships, namely, nelatedness, equivalence, broadness, and narrowness, are selected for each pair of senses in the two resources by mative lexicographics. Given the lexicographic context of this study, we have tried to provide lexicographic data from expert-made dictionairs. We believe that our datasets will pave the way for further developments in exploring statistical and neural methods, as well as for evaluation purposes. The rest of the paper is organized as follows: we first describe the previous work in Section 2. After having de-

²https://www.wikipedia.org

* Contact Authors https://globalex2020.globalex.link/

→ 17 manually-annotated monolingual resources for the task of WSA covering 15 languages, including **Danish**



3. Word sense alignment Monolingual WSA datasets

(i montgolfier (noun) F	(used of style of speaking) overly embellished of or relating to elocution French inventor who (with his brother Josef Michel French inventor who (with his brother Jacques Etie	NONE - exact - NONE -	▼ 0 -pertaining to elc ▼	0 -pertaining to elocution.	"lemma": "splenetic",
o montgolfier (noun) F	of or relating to elocution French inventor who (with his brother Josef Michel French inventor who (with his brother Jacques Etie	exact •	0 -pertaining to elc 🔻		renna . sprenetic ,
montgolfier (noun) F F	French inventor who (with his brother Josef Michel French inventor who (with his brother Jacques Etie	NONE -			"POS tag": "adjective"
F F	French inventor who (with his brother Josef Michel French inventor who (with his brother Jacques Etie	NONE -			"gender": ""
F	French inventor who (with his brother Jacques Etie		0 -a balloon which -	0 -a balloon which ascends by the buoyancy	"meta ID": "".
diag (work)		NONE -	0 -a balloon which 👻		"resource 1 senses":
aice (verb)					{
c	cut into cubes	exact -	1. -to cut into smal 💌	1. -to cut into small cubes; .	"#text": "of or relating to the spleen"
p	play dice	NONE -	~	2. -to ornament with squares, diamonds, or	<pre>"external_ID": "splenic.a.01"},</pre>
ebb (verb)					{
fl	flow back or recede	exact -	0 -to cause to flow 👻	0 -to cause to flow back.	"#text": "very irritable",
fr	fall away or decline	NONE -			<pre>"external_ID": "bristly.s.01"}</pre>
h	hem in fish with stakes and nets so as to prevent th	NONE -],
educated (adjective)					"resource_2_senses":
C	characterized by full comprehension of the problem	Į		0 -formed or developed by education; .	(
p	possessing an education (especially having more t	exact			" #Cext ": "allected with spieen; mailcion
quaver (verb)		parrower			→ spiterui, peevisi, itetiui. ,
g	give off unsteady sounds, alternating in amplitude o	narrower		0 -to utter with quavers.	
S	sing or play with trills, alternating with the half note	broader	~		"alignment": [
rangy (adjective)		related			
а	adapted to wandering or roaming	NONE	.	0 -inclined or able to range, or rove about, fo	"sense_source": "very irritable",
а	allowing ample room for ranging	T	-		"sense_target": "affected with spleen;
tr	all and thin and having long slender limbs	Ψ	•		→ malicious; spiteful; peevish;
smiler (noun)					
а	a person who smiles	v	•	0 -one who smiles.	<pre>"semantic_relationship": "exact"}</pre>
tł	the human face (`kisser' and `smiler' and `mug' are	*			
chair (noun)					}



Semantic relations (based on SKOS*)

- **exact**: The sense are the same, for example the definitions are simply paraphrases
- **broader**: The sense in the first dictionary completely covers the meaning of the sense in the second dictionary and is applicable to further meanings
- **narrower**: The sense in the first dictionary is entirely covered by the sense of the second dictionary, which is applicable to further meanings
- **related**: There are cases when the senses may be equal but the definitions in both dictionaries differ in key aspects
- **none** : There is no match for this sense

* https://www.w3.org/2004/02/skos/



Semantic relations: an example



Source: https://www.w3.org/2004/02/skos/core/guide/2005-10-06/

european lexicographic infrastructure

Manual alignment of ODS-DDO senses, examples

- **ODS exact DDO**: noun *passager* ('passenger'): ODS 'person traveling with mail coach etc.', DDO: 'person traveling with private or public means of transportation'.
- ODS broader than DDO: noun værge ('guardian'): ODS 'a guardian of anything or anybody', DDO 'a guardian in legal context'.
- ODS narrower than DDO: adjective spids ('sharp'): ODS two narrower senses, one about a sound and another one about a smell → DDO 'pungent in an unpleasant way (about smell, taste or sound)'.
- ODS related to DDO: noun søvn ('sleep'): ODS 'being able to sleep', DDO 'the state of sleeping'



Manual alignment of ODS-DDO senses

- 380 nouns, 93 verbs, 63 adjectives.
- ODS 3,595 senses, DDO 1,667 senses
 → 1,000 exact matches
- Polysemy: average DDO 3.11 senses per lemma (ODS 6.7)
- Difference in dictionary structure:
 - ODS senses include MWU
 - DDO senses: no MWU
 - ODS: heading main senses
 - DDO: no heading main senses

Matches ODS-DDO



Language	Resource	Nouns	Verbs	Adjectives	Adverbs	Other	All
Pasqua	Basque Wordnet	929 (6836)	0 (0)	0 (0)	0 (0)	0 (0)	929 (6836)
Basque	Euskal Hiztegia	971 (7754)	0 (0)	0 (0)	0 (0)	0 (0)	971 (7754)
Dulasia	BTB-WN	1394 (15649)	175 (1698)	305 (3187)	50 (338)	0 (0)	1924 (20872)
Bulgarian	Bulgarian Wik- tionary	1273 (12883)	164 (1107)	194 (1418)	39 (306)	0 (0)	1670 (15714)
Danish	Ordbog over det danske Sprog	2176 (282040)	983 (119163)	436 (60599)	0 (0)	0 (0)	3595 (461802)
	Den Danske Ordbog	1036 (12326)	383 (4045)	248 (2228)	0 (0)	0 (0)	1667 (18599)
Dutch	Woordenboek der Nederlandsche Taal	1459 (28979)	405 (5185)	527 (7878)	106 (2662)	0 (0)	2497 (44704)
	Algemeen Neder- lands Woordenboek	497 (8443)	140 (1542)	109 (1393)	13 (172)	0 (0)	759 (11550)
Enalish (KD)	Global	92 (532)	107 (617)	80 (457)	57 (257)	61 (283)	397 (2146)
English (KD)	Password	66 (536)	72 (417)	62 (324)	33 (177)	46 (188)	279 (1642)
	Webster	1131 (11606)	741 (4622)	373 (2585)	45 (269)	0 (0)	2290 (19082)
English (NUIG)	Princeton WordNet	730 (12166)	496 (6980)	249 (2892)	24 (207)	0 (0)	1499 (22245)
Estonian	Dictionary of Esto- nian (EKS)	543 (4012)	273 (1598)	151 (747)	98 (451)	78 (370)	1143 (7178)
	Estonian Basic Dic- tionary (PSV)	543 (4492)	273 (1983)	151 (1097)	98 (596)	79 (468)	1144 (8636)
German	German Wiktionary	2026 (15160)	0 (0)	0(0)	0(0)	0 (0)	2026 (15160)
	German OmegaWiki	1266 (14354)	0 (0)	0(0)	0(0)	0 (0)	1266 (14354)
TT 1. (2011) (2010)	Comprehensive				a contrato		1355 (14654)
Hungarian	Explanatory						1038 (10934)
Irish	An Foclóir Beag	891 (8053)	11 (95)	55 (267)	10 (56)	36 (171)	1003 (8642)
	Irish Wiktionary	1209 (6696)	8 (45)	61 (181)	10 (41)	36 (109)	1324 (7072)
Italian	ItalWordNet	408 (3128)	352 (2411)	0 (0)	0(0)	0 (0)	760 (5539)
	SIMPLE	290 (1990)	218 (1240)	0 (0)	0 (0)	0 (0)	508 (3230)
0.11	Serbian WordNet	691 (5864)	985 (6522)	92 (713)	0(0)	0 (0)	1768 (13099)
Serbian	Dictionary of Serbo- Croatian Literary Language	289 (2360)	281 (1527)	29 (215)	0 (0)	0 (0)	599 (4102)
1010 D. 1922-020	Slovene WordNet	409 (1106)	303 (901)	237 (733)	44 (133)	0 (0)	993 (2873)
Slovenian (JSI)	Slovene Lexical Database	284 (2237)	191 (1047)	220 (1486)	29 (102)	0 (0)	724 (4872)
Slovenian (ISJFR)	Standard Slovenian Dictionary (eSSKJ)	229 (2060)	109 (911)	76 (620)	0 (0)	60 (588)	474 (4179)
	Kostelski slovar	151 (1050)	61 (308)	45 (257)	0(0)	38 (263)	295 (1878)
Spanish	Diccionario de la lengua española	617 (7986)	225 (2426)	305 (3269)	26 (161)	24 (250)	1197 (14092)
	Spanish Wiktionary	602 (6421)	227 (2045)	294 (2825)	25 (129)	22 (123)	1170 (11543)
Deter	Dicionário da Língua Portuguesa Contemporânea	285 (4060)	58 (686)	110 (1287)	9 (143)	1 (9)	463 (6185)
Portuguese	Dicionário Aberto	199 (1521)	53 (203)	67 (372)	3 (15)	1 (5)	323 (2116)
Duration	Ozhegov-Shvedova	258 (2038)	109 (615)	101 (533)	15 (77)	44 (368)	527 (3631)
Kussian	Dictionary of the Russian Language (MAS)	310 (2811)	173 (1338)	190 (1219)	20 (114)	71 (1010)	764 (6492)

european lexicographic infrastructure

MWSA Datasets **Statistics** Number of senses and number of tokens (in parentheses)

Data openly available at: https://github.com/elexis-eu/mwsa

23

This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 731015.



4. Experiments with ODS and DDO

→ linking lexical content at sense level, including glosses

In order to predict the semantic similarity between two given senses, we trained various models based on the following similarity features:

1. String metrics

- Longest common substring
- Length ratio
- Average word length ratio
- Jaccard, Dice, and Containment
- 2. Word Embeddings (based on cosine similarity)
- 3. Automatic feature extraction

4. Experiments with ODS and DDO **Naisc**

- Naisc is a system for aligning RDF datasets
- It takes as input 2 RDF documents
- It outputs an alignment (set of RDF triples) between these two documents





4. Experiments with ODS and DDO Results

Performance of our similarity detection models for automatic alignment of DDO and ODS within a specific limit of space-separated tokens (15, 20, 25 and all tokens)

ODS sense size	Model	Precision	Recall	F-measure
15	String metrics	65.3%	48.1%	55.4%
	Word Embeddings	66.7%	48.0%	55.8%
	Auto	64.0%	46.6%	54.0%
20	String metrics	61.5%	44.3%	51.5%
	Word Embeddings	64.7%	46.7%	54.3%
	Auto	63.3%	45.8%	53.2%
25	String metrics	57.5%	21.9%	31.7%
	Word Embeddings	55.9%	21.2%	30.8%
	Auto	58.5%	22.2%	32.1%
All	String metrics	54.7%	9.8%	16.7%
	Word Embeddings	50.7%	9.7%	16.3%
	Auto	50.3%	9.4%	15.8%



4. Experiments with ODS and DDO **Results**



The correlation of sense sizes in ODS with F-measure using various methods



Future steps

- Text paraphrasing: converting ODS senses into more representative ones by
- Removing excessive descriptions, e.g. historical citations
- Lexical and orthographic normalization
- Paraphrasing techniques