

Terminologi på tværs af Danmark – og EU

Om **eTranslation TermBank** -
indsamling af terminologi til
maskinoversættelse

Sussi Olsen

**Lina Henriksen, Bolette S. Pedersen,
Claus Povlsen**

Center for Sprogteknologi,
Nordiske Studier og Sprogvidenskab

UNIVERSITY OF COPENHAGEN



Indhold

- Om projektet
- Motivation - terminologi og maskinoversættelse
- Forvaltning af terminologi i Danmark
- Resultater af indsamling
- Fremtidigt arbejde



eTranslation TermBank

- EU-projekt, løber sept. 2017- feb. 2019
- Formål: Anspore til indsamling af terminologiske ressourcer inden for specifikke fagdomæner for at forbedre kvaliteten og dækningsgraden af EU's maskinoversættelsessystem CEF eTranslation.
- Ressourcer for alle officielle EU-sprog + norsk (bokmål og nynorsk) og islandsk
- Fokus på domæner knyttet til følgende online-tjenester:
 - Sundhed – eHealth
 - Retsinformation – e-Justice
 - Forbrugerbeskyttelse – ODR Open Dispute Resolution



Projektet udsprunget af ELRC: European Language Resource Coordination

- Netværk der indsamler sprogressourcer for alle EU-sprog til CEF eTranslation.
- Fokus på parallelle ressourcer, oversatte tekster.
- Man kan selv uploade sprogressourcer til websitet: <http://lr-coordination.eu/resources>
- Indsamler også terminologi, tung proces -> eTTB-projektet

European Language
Resource Coordination



Før graviditeten

- Du kan forberede dig på en sund graviditet, allerede inden du bliver gravid. Begynd at tage et kosttilskud af folat. Og hold dig til alkohol og rygning.

Folat – et kosttilskud
Tag 400 mikrogram folat dagligt. Du kan overveje at blive gravid. Folat er en del af gruppen af B-vitaminer. Folat er nødvendigt for at barnet bliver født med mindstetallet af rygmarvs (neurorøds)fejl.

Ingen alkohol
Alkohol kan skade barnet. På graviditeten udsæder og frem. Da det kan være svært at holde det nøjagtigt tidspunkt for undfangelsen, bør du ikke være med at drikke alkohol, hvis du planlægger graviditet.
Hvis graviditeten er uoplagt, og du har drukket alkohol, er der dog kun sjældent grund til bekymring. Tal evt. med læger/pejdemoderen, om det ved første graviditetstestning.



Op med at ryge, allerede inden du planlægger.
Tjek medicinen
Er du i behandling med medicin i den periode, hvor du planlægger at blive gravid, bør du tale med

Before getting pregnant

- You can prepare for a healthy pregnancy even before you become pregnant. Start taking folic acid. And stay clear of alcohol and smoking.

Folic acid – a dietary supplement
You should begin taking 400 mcg folic acid daily as soon as you start planning your pregnancy. Folic acid is part of the vitamin B group. Folic acid reduces the risk of your baby being born with neural tube defects.

No alcohol
Alcohol can harm your baby from the start of pregnancy and onwards. As it is often difficult to know the exact time of conception, you should not drink alcohol if you are planning to have a baby.

However, if you find that you are pregnant without having planned it and you have been drinking alcohol, it is rarely a cause for concern. You may like to talk to your doctor/midwife about this when going for your first consultation.



Check medicines
If you are being treated with any medicines when you begin planning

Partnere i projektet

- Tilde, Letland (koordinator)
- Københavns Universitet, Danmark
- Árni Magnússon Institute for Icelandic Studies, Island
- Institute of the Estonian Language, Estland
- Jožef Stefan Institute, Slovenien
- International Network for Terminology, Østrig
- Institute for Language and Folklore, Sverige
- Swedish Centre for Terminology, Sverige
- Institute of the Lithuanian language, Litauen

Hver partner indsamler terminologi for flere sprog.
Danmark er også ansvarlig for norsk og ungarsk.



Hvem er eTTB i DK?



Medarbejdere ved Center for Sprogteknologi

- Lina Henriksen linah@hum.ku.dk
- Sussi Olsen saolsen@hum.ku.dk
- Bolette Sandford Pedersen bspedersen@hum.ku.dk
- Claus Povlsen cpovlsen@hum.ku.dk





Motivation – terminologi
og maskinoversættelse

Fejloversættelse af termer – ex: sundhedsdomænet

Eksemplerne er fra professor emerita Bodil Nistrup Madsen, CBS

Dansk (fra begrebsbasen http://sundhedsdata.iterm.dk)	Engelsk (Google)	Definition
engangsdosering	prefilled <i>Alternativ:</i> once a day regimen	"enkelstående <u>dosering</u> af én <u>lægemiddeldosis</u> "
pausering	break ring	"midlertidig afbrydelse af <u>lægemiddelordination</u> efter anvisning af en dertil autoriseret person"

Maskinoversættelse og terminologi



Regelbaseret maskinoversættelse – start i 1940'erne, populære i 80'erne

- Ordbøger og håndskrevne grammatiske regler, terminologi var let integrerbar

Statistisk maskinoversættelse – også kaldet phrase based - overtog i løbet af 90'erne

- Arbejder ud fra store mængder af parallelle tekster på de to sprog, danner sætninger ud fra lign. sætninger i træningsmaterialet. Mere flydende oversættelser. Terminologi – skal håndteres særskilt

Maskinoversættelse og terminologi

Neural maskinoversættelse – siden 2014

- Kontekstbaseret, betjener sig af 'deep learning', arbejder med vektorer, kræver mindre computerhukommelse
- Danner mere flydende sætninger (posteditering nedsat med 25%) og dermed bedre output
- Terminologi skal 'indlæres'

Anvendelse af terminologi i maskinoversættelse

- Monolingual: Præprocessering – domænebestemmelse
- Bilingual: en term skal 'overrule' det ord som ellers typisk optræder i en tilsvarende kontekst

Behov for terminologi inden for alle domæner

Hvad er CEF eTranslation?

CEF står for Connecting Europe Facility – et EU-program

Investerer i og støtter opbygningen af digitale infrastrukturer som kan forbedre dagliglivet for europæiske borgere, samt byggeblokke hertil.

Eksempler på digitale infrastrukturer (DSIs):

- ODR - online dispute resolution portal
- e-Justice
- eHealth

Eksempler på byggeblokke:

- eDelivery
- eSignature
- eInvoicing
- **eTranslation**

eTranslation er en udbygning af EU's tidligere maskinoversættelsesystem.



Terminologi i Danmark

Forvaltning af terminologi i Danmark

- Ingen central organisation for terminologiforvaltning, modsat fx Sverige og Norge
- 1998-2015: DanTermCentret arbejdede med at informere om og oplære i systematisk behandling af terminologi (få midler, lukket 2015).
- CBS – havde forskning i terminologi, afd. lukket
- CST, KU deltog i et EU-terminologiprojekt i 1992 hvor termbankerne IATE og EuroTermBank blev udviklet.
- Sprogteknologisk Udvalg – nedsat af regeringen 2018 - skal afklare behovet og perspektiverne for en national termbank.

Terminologinetværk i Danmark

FORVIR – Forum for Videnmodellering i Offentligt Regi

- Et netværk for ansatte i den offentlige administration. Formålet er at styrke terminologisk begrebsarbejde kvantitativt og kvalitativt

Terminologigruppen

- Et netværk for både den private og offentlige sektor. Formålet er at udveksle erfaringer og deltage i nordisk og internationalt terminologiarbejde

Terminologiarbejde i danske offentlige institutioner

Visse offentlige institutioner gør meget ud af terminologi- og begrebsarbejde, fx

- Sundhedsdatastyrelsen
- Socialstyrelsen
- Domstolsstyrelsen



Terminologiarbejde i danske offentlige institutioner

MEN

- Ca. 75 % af de offentlige institutioner som vi har henvendt os til, udfører ikke terminologiarbejde og har ingen termsamlinger.
- Kun 25% udfører terminologiarbejde, de fleste kun på dansk
- Resten 'outsourcer' deres oversættelser til oversættelsesbureauer eller klarer sig med google!



Indsamling

Hvem kan levere terminologi til eTTB?

I princippet alle.

Vi har primært henvendt os til offentlige institutioner:

- De er forpligtede til at levere deres data if. PSI-direktivet (Lov om videreanvendelse af den offentlige sektors informationer)
- Adgang til CEF eTranslation er lige nu forbeholdt offentlige institutioner (det forlyder at systemet bliver åbent for alm. borgere og private virksomheder fra 2020)

Vi har også henvendt os til halvoffentlige og private institutioner som vi har hørt skulle råde over terminologi.

Manglende interesse for og vilje til at levere terminologi

Det kræver meget tid og gentagne henvendelser at få leveret terminologi.

- Største problem: manglende kendskab til og interesse for CEF eTranslation


Mange spørgsmål og megen usikkerhed

- Hvem vil bruge terminologien og hvordan?
- Hvor lagres den?
- Kan data blive misbrugt?
- Hvad nu hvis vi opdaterer vores data og denne opdatering ikke sker i det vi afleverer?
- Kan vores data ende med at gøre mere skade end gavn?

Foreløbigt resultat af indsamlingen

Vi har henvendt os til godt 30 institutioner

Vi har modtaget terminologi fra nedenstående:

- Socialstyrelsen
- Beskæftigelsesministeriet
- Københavns Universitet, CIP 
- Forsknings- og Uddannelsesministeriet
- Sundhedsdatastyrelsen
- LO
- Privat terminolog



Konklusion - og
de sidste måneder
i projektet

Erfaringer fra indsamlingen – udfordringer

- Uvidenhed om hvor værdifulde tekster og termlister er for forskning og udvikling i dansk sprogteknologi
- Mange offentlige institutioner anvender hverken terminologiske resurser el.lign. som kan sikre ensartet tekstkvalitet
- sprogarbejde er ikke prioriteret - udliciteres
- Derfor ingen stor interesse for CEF eTranslation
- Usikkerhed om IPR, opdatering af versioner etc.
- Vi har ikke en klar fælles sprogpolitik og fælles fokus på sproglige problemstillinger i det offentlige (som andre nordiske lande)

Fremtidige opgaver i eTTB

- Flere ressourcer meget velkomne, så hvis I har kendskab til terminologi – så kontakt os endelig
- Det samme gælder hvis I kender til større parallelsprogede publikationer
- Alle domæner er relevante

Ellers vil vi nu koncentrere os om norske og ungarske ressourcer – for

- indsamling af dansk terminologi har vist sig at være temmelig tungt, men det er en leg sammenlignet med indsamling af ungarsk terminologi.



Tak for jeres
opmærksomhed!
Spørgsmål?