



UPPSALA
UNIVERSITET



Maskinöversättning

Vad och hur?

Anna Sågvall Hein

2016-03-10

Översikt

Vad är maskinöversättning?

Kort tillbakablick

Varför är det så svårt?

Språkvetenskaplig
– regelbaserad maskinöversättning

Statistisk maskinöversättning

Maskinöversättning i praktiken

Forskning i fokus



UPPSALA
UNIVERSITET



Kort tillbakablick

De första försöken gjordes redan på 50-talet i USA och Sovjet med en primitiv teknologi.

De gällde översättning mellan ryska och engelska för politiska och militära ändamål och aktörerna var tekniker.

De högt ställda förhoppningarna grusades och efter några års stiltje inleddes en ny period med språkvetare/datorlingvister och datavetare i spetsen.

De språkvetenskapligt baserade metoder som började utvecklas på 60-talet har i modifierad form överlevt till i dag, då de fått stark konkurrens av statistiska metoder.



UPPSALA
UNIVERSITET



Varför är det så svårt?

Flertydighet

- ett ord kan representera olika ordklasser, former och betydelser

Variation

- olika ord för samma företeelse – synonymi

Språkskillnader

- ordböjning, ordföljd



UPPSALA
UNIVERSITET



Exempel

Boken *ligger under* bordet.
The book *is under* the table.

Studenten *somnade under* föreläsningen.
The student *fell asleep during* the lecture.

Bara ett *under* kan rädda tidningen.
Only a *miracle/wonder* can save the newspaper.

Under tiden stannade bilen.
Meanwhile, the car stopped.

Var snäll och skriv under här.
Please, sign here.



UPPSALA
UNIVERSITET



Regelbaserad översättning

ANALYS

skapar en grammatisk struktur, där översättningsalternativ känns igen och väljs

TRANSFER

översätter den grammatiska strukturen till en grammatisk struktur på målspråket

GENERERING

bygger upp en översättning på den grammatiska strukturen

Alla komponenterna använder sig av omfattande lexikon, grammatiker och transferregler



Analysstruktur

Under tiden stannade bilen

[phr.cat	cl]		
	cl.type	main			
	mode	decl			
	verb	[]	
		word.cat			VERB
		diat			act
		inff			fin
		verb.type			main
	tense	past			
	head	[]	
	lex	stanna.vb.1			
	trglex	stop.vb.1			
adv.in.fund	[]			
	word.cat		ADV		
	lex		under_tiden.ab.1		
	trglex	meanwhile.ab.1			
subj	[]			
	word.cat		NOUN		
	gender		utr		
	numb		sing		
	def		def		
	case		basic		
	person		3		
head	[]			
	lex		bil.nn.1		
	trglex	car.nn.1			

”Under tiden stannade bilen”



Transferstruktur

[phr.cat	cl]								
	cl.type	main										
	mode	decl										
	adv.in.fund	[word.cat	ADV]							
		lex	meanwhile.ab.1									
	pred	[subj	[word.cat	NOUN]					
				numb	sing							
				def	def							
				case	basic							
				person	3							
				head	[lex car.nn.1]							
				verb	[word.cat	VERB]				
									diat	act		
									inff	fin		
									verb.type	main		
	tense	past										
	head	[lex stop.vb.1]										



Genereringsstruktur

Meanwhile, the car stopped

"Meanwhile"	word.cat	ADV
ABP	degree	pos
	lex	meanwhile.ab.1
" , "	word.cat	SEP
SRI	type	minor
	lex	comma.sr.1
"the"	word.cat	ART
AL	det.type	det
	lex	the.al.1
"car"	word.cat	NOUN
NNSB	numb	sing
	case	basic
	person	3
	lex	car.nn.1
"stopped"	word.cat	VERB
VBRM	inff	fin
	verb.type	main
	tense	past
	pv	-
	numb	sing
	person	3
	lex	stop.vb.1

"Meanwhile, the car stopped"

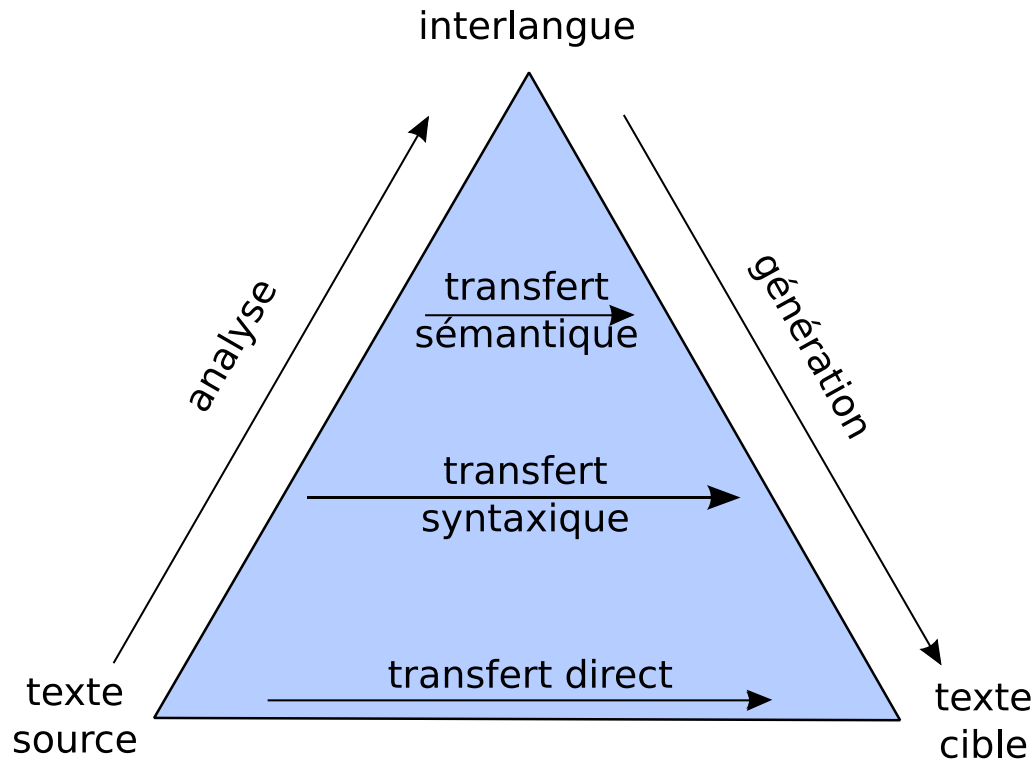


UPPSALA
UNIVERSITET



Transferöversättning på olika nivåer

Vauquois triangel





UPPSALA
UNIVERSITET



Transferbaserad MT

Källspråkstroagna översättningar

Kräver stora lexikon och grammatiker

Tar tid att bygga

Brister som uppdagas kan åtgärdas genom ändringar i lexikon och grammatik



UPPSALA
UNIVERSITET



Statistisk maskinöversättning

Ett statistiskt översättningssystem, en SMT, använder sig inte av någon språkkunskap utan bygger helt på tidigare översättningar

En SMT består av två huvudkomponenter

- en sannolikhetsbaserad översättningsmodell

- en modell över målspråket, s.k. språkmodell

Kvaliteten på de översättningar som en SMT producerar beror helt på kvaliteten på de båda huvudkomponenterna



UPPSALA
UNIVERSITET



Översättningsmodellen

Modellen består av sannolikhetsbaserade översättningar av ord och längre uttryck.

För varje delöversättning anges ett sannolikhetsvärde.

Översättningar av hela meningar byggs upp genom att delöversättningar kombineras på ett sätt som ger den totalt sett högsta sannolikheten.

Ett stort antal översättningar produceras.

Problem kan uppstå i gränsen mellan olika delöversättningar.



UPPSALA
UNIVERSITET



Språkmodellen

Språkmodellen innehåller stora mängder målspråksdata.

De är organiserade som enordningar och flerordningar – s.k. ngram.

De högst rankade översättningarna jämförs mot språkmodellen och den översättning som stämmer bäst väljs.

Träning av översättningsmodellen

Insamling och rensning av träningsdata i form av tidigare översättningar
av rätt slag

språkpar, texttyp, ämnesområde

Parallellställning – länkning – av träningsdata på olika nivåer
meningar, ord, fraser

Meningslänkade träningsdata

Regeringsförklaringen 1988

s1.1	REGERINGSFÖRKLARING	Statement of Government Policy by the Prime Minister , Mr Ingvar Carlsson , at the Opening of the Swedish Parliament on Tuesday , 4 October , 1988 .	s1.1
s2.1	Eders Majestäter , Eders Kungliga Högheter , herr talman , ledamöter av Sveriges riksdag	Your Majesties , Your Royal Highnesses , Mr Speaker , Members of the Swedish Parliament .	s2.1
s3.1	Sveriges neutralitetspolitik är av avgörande betydelse för vårt lands fred och oberoende	Sweden' s policy of neutrality is of decisive importance for our peace and independence .	s3.1
s3.2	Den bidrar också till stabilitet och avspänning i vår del av världen	It also contributes to stability and détente in our part of the world .	s3.2
s3.3	Kring denna politik finns en bred folklig uppslutning	There is wide popular support for this policy .	s3.3
s3.4	Den kommer att fullföljas med kraft och konsekvens	It will be pursued with firmness and consistency .	s3.4
s4.1	Neutralitetspolitiken stöds av ett starkt försvar till värn för vårt oberoende	Our policy of neutrality is underpinned by a strong defence . That safeguards our independence .	s4.1 s4.2
s4.2	Kränkningar av svenskt territorium kommer aldrig att accepteras	Violations of Swedish territory will never be accepted .	s4.3
s4.3	Armén kommer att reformeras och effektiviseras	The army will be reorganized with the aim of making it more effective .	s4.4
s4.4	Det är regeringens föresats att söka breda lösningar i frågor som är av betydelse för vår nationella säkerhet	It is the Government' s intention to seek broad solutions in issues that are of importance for our national security .	s4.5
s5.1	Regeringen har välkomnat överenskommelsen mellan Förenta staterna och Sovjetunionen om att avskaffa de landbaserade medeldistanskärnvapnen	The Government welcomed the agreement between the United States and the Soviet Union on the elimination of land- based intermediate- range nuclear weapons .	s5.1



UPPSALA
UNIVERSITET



Ordlänkning - samförekomst

Sveriges **neutralitetspolitik** är av avgörande betydelse för vårt lands fred och **oberoende**.

Sweden's **policy of neutrality** is of decisive importance for our peace and **independence**.

Neutralitetspolitiken stöds av ett starkt försvar till värn för vårt **oberoende**.

Our **policy of neutrality** is underpinned by a strong defence. That safeguards our **independence**.



UPPSALA
UNIVERSITET



Fraslänkning

Kring **denna politik finns** en bred politisk uppslutning.

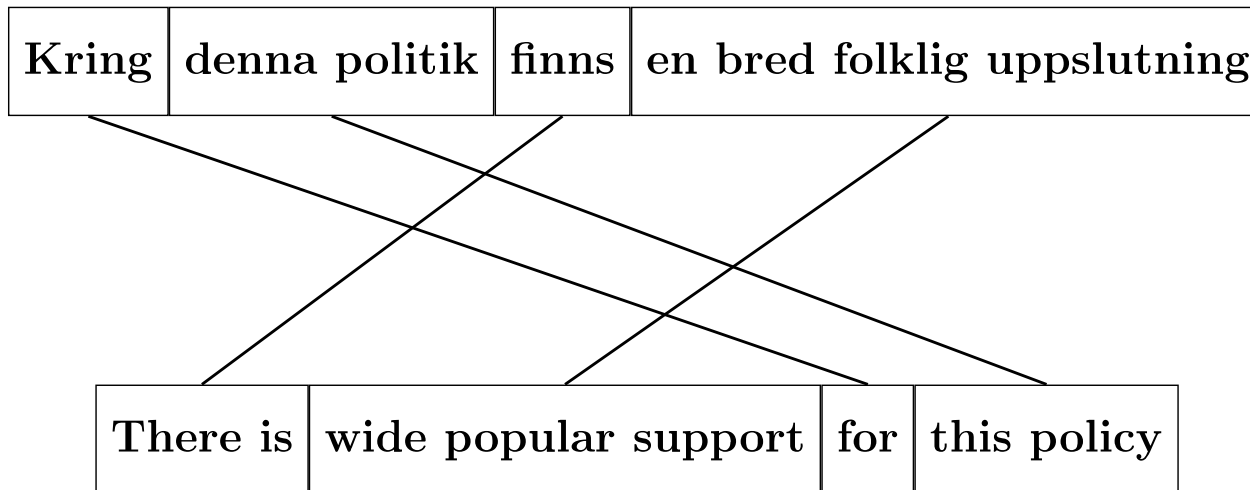
There is wide popular support for **this policy**.



UPPSALA
UNIVERSITET



Ord- och fraslänkning





UPPSALA
UNIVERSITET



Resurser

För ett nytt språkpar behövs en översättningskorpus på minst 1 miljon ord och enspråkiga korpusar på över 1 miljard ord vardera (Franz Josef Och, Google)

Det finns open-source-mjukvara för att träna upp en SMT (t.ex. Moses)



UPPSALA
UNIVERSITET



SMT:s beroende av träningsdata

Under tiden somnade studenten.
Meanwhile, the student fell asleep.

Under tiden vaknade studenten.
Meanwhile woke student.

(Google 2016-03-09)



Utvärdering

Manuell

Automatisk



UPPSALA
UNIVERSITET



Automatisk utvärdering

Ett facit upprättas i form av en översättningskorpus

Översättningskorpusen hålls utanför träningsdata

En automatisk jämförelse görs mellan översättningen i korpusen och den maskingenererade översättningen

Ett likhetsmått beräknas, t.ex. BLEU som går mellan 0 och 1

Det har sagts att gränsen för att det ska löna sig att redigera en maskinöversatt text ligger vid ett BLEU-värde runt 0.4



UPPSALA
UNIVERSITET






Maskinöversättning i praktiken

Maskinöversättning + postredigering för publicering

Maskinöversättning som stöd för mänskliga översättare

Snabb råöversättning - gisting

Post-redigering Convertus-systemet

[View original](#) | [View translation](#) |  100 |  17 |  0

 Smooth scroll off |  Spotlight off |  

Close

Mark as done

Source text

Translation

KURSPLAN

Course syllabus

M0051H Diagnostiska modaliteter, 7,5 Högskolepoäng

M0051H Diagnostic modalities, 7.5 Credits


Diagnostic modalities

Diagnostic modalities

Gäller för perioden

Applies for the period

Vald version visar för vilken termin och läsperiod som denna kursplanen gäller för.

 Copy source

 Clear

 Reset



The selected version shows for which semester and study period this course syllabus applies.



Senaste version visas först.

Latest version is shown first.

Kursplanen fastställd

Course syllabus established

av Prefekt vid Institutionen för hälsovetenskap 2008-12-12

by the Head of the Department of Health Sciences 12/12/2008

Forskning i fokus

Hybridsystem

- kombination av regelbaserade och statistiska system

Beroenden över meningsgränserna

- t.ex. syftningsproblem

Språkanpassning

- Inför träningen förbehandlar man källspråket för att göra det mer likt målspråket