

“Wie funktioniert automatisierte Übersetzung?”

Prof. Josef van Genabith

(Deutsches Forschungszentrum für Künstliche Intelligenz)

Überblick:

- Warum MÜ: Datenmenge, Qualität und Kosten?
- Was maschinelle Übersetzung so schwierig?
- MÜ + **Professionelle** Übersetzer = Qualität
- Wie funktioniert die moderne statistisch-basierte MÜ?
- Es geht vor allem um Daten!
- Und um die richtige Art von Daten!

- Europa = Mehrsprachigkeit
- 24 Amtssprachen
24+2 CEF Sprachen
- Vieles zu übersetzen!
- Übersetzungskosten!?
- Kann MÜ Hilfe leisten?
- Und in welcher Qualität?

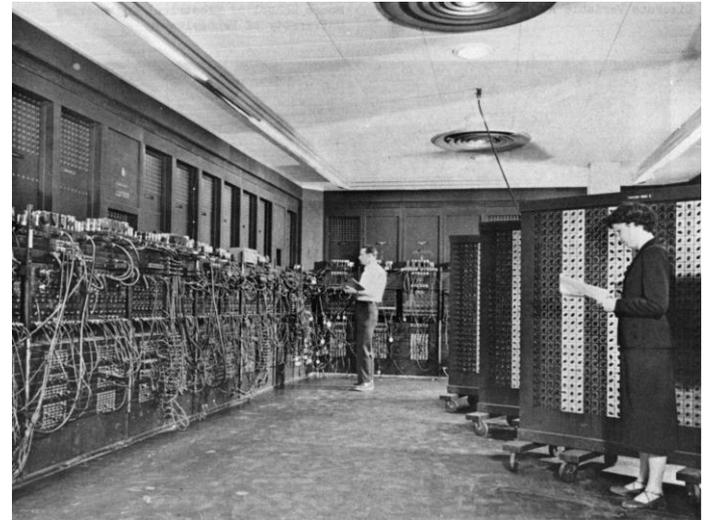


Image: <https://en.wikipedia.org/wiki/ENIAC#/media/File:Eniac.jpg>
License: public domain

- Natürliche Sprachen sind:
 - Elegant
 - Effizient
 - Flexibel
 - Komplex
- Ein Wort/Satz kann verschiedene bedeuten
- Mehrere Möglichkeiten, das Gleiche zu sagen
- Bedeutung hängt von Kontext ab
- Übertragener Sinn (Metapher)
- Sprache und Kultur (unterschiedliche Konzeptualisierungen des gleichen Sachverhalts)
- Wortstellung
- Morphologie u.v.m.

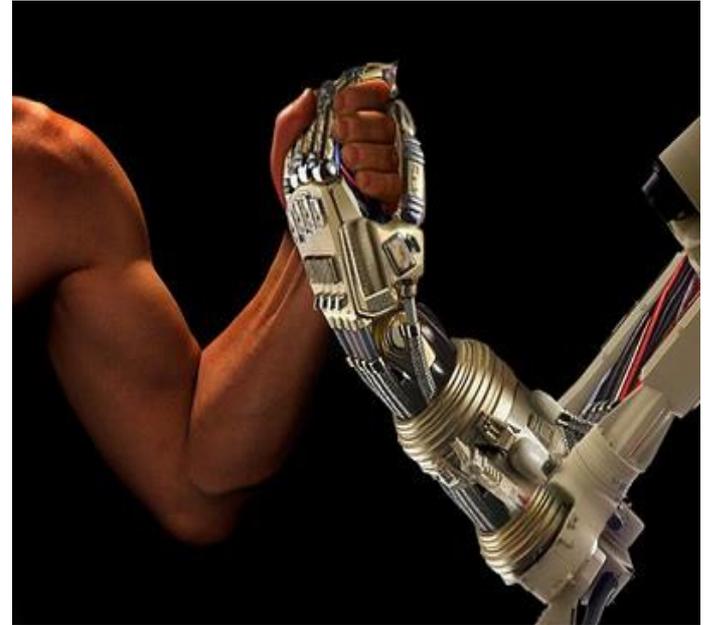
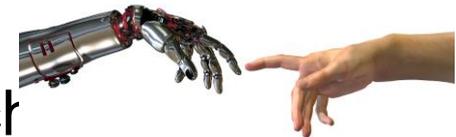
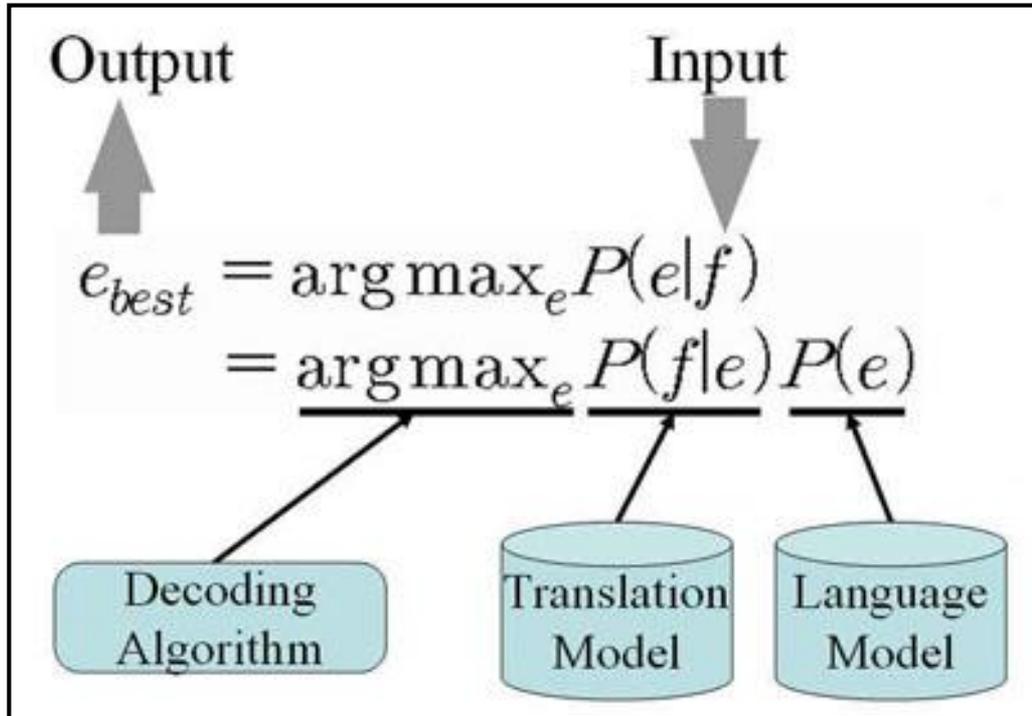


Image: <http://workingtropes.lmc.gatech.edu/wiki/index.php/File:Man-vs-machine.jpg>
License: CC BY-NC-SA 3.0



- Sprache und Übersetzung sind komplex
- Wir können sie nicht genau berechnen
- Wir haben bereits regel-basierte MÜ und Sprachtechnologie erforscht und eingesetzt
- Was nun?
- Maschinelles Lernen
 - Aus **Daten** lernen \Rightarrow zentrale Rolle von Daten
 - Grobe Lösung \Rightarrow nicht perfekt, Hilfe
 - Professionelle Übersetzer
 - Nach-editieren
 - Automatisierte Übersetzung \neq automatisch





- Heute kein Mathematikunterricht!
- Sondern:
- Die Geschichte der statistischen MÜ in Bildern ...
- Es dreht sich einzig und allein um **Daten**
...

Die statistische MÜ lernt aus zwei Typen von Daten:

- Übersetzungen von Menschen
- Text in der Zielsprache
- So viele Daten wie möglich!
- Aber: die richtige Art von Daten!

GERMAN

Einleitung

I. Von dem Unterschiede der reinen und empirischen Erkenntnis

Daß alle unsere Erkenntnis mit der Erfahrung anfangt, daran ist gar kein Zweifel; denn wodurch sollte das Erkenntnisvermögen sonst zur Ausübung erweckt werden, geschähe es nicht durch Gegenstände, die unsere Sinne rühren und teils von selbst Vorstellungen bewirken, teils unsere Verstandstätigkeit in Bewegung bringen, diese zu vergleichen, sie zu verknüpfen oder zu trennen, und so den rohen Stoff sinnlicher Eindrücke zu einer Erkenntnis der Gegenstände zu verarbeiten, die Erfahrung heißt? Der Zeit nach geht also keine Erkenntnis in uns vor der Erfahrung vorher, und mit dieser fängt alle an.

ENGLISH

Introduction

I. Of the difference between Pure and Empirical Knowledge

That all our knowledge begins with experience there can be no doubt. For how is it possible that the faculty of cognition should be awakened into exercise otherwise than by means of objects which affect our senses, and partly of themselves produce representations, partly rouse our powers of understanding into activity, to compare to connect, or to separate these, and so to convert the raw material of our sensuous impressions into a knowledge of objects, which is called experience? In respect of time, therefore, no knowledge of ours is antecedent to experience, but begins with it.

FRENCH

Introduction

I. De la différence de la connaissance pure et de la connaissance empirique.

Que toute notre connaissance commence avec l'expérience, cela ne soulève aucun doute. En effet, par quoi notre pouvoir de connaître pourrait-il être éveillé et mis en action, si ce n'est par des objets qui frappent nos sens et qui, d'une part, produisent par eux-mêmes des représentations et, d'autre part, mettent en mouvement notre faculté intellectuelle, afin qu'elle compare, lie ou sépare ces représentations, et travaille ainsi la matière brute des impressions sensibles pour en tirer une connaissance des objets, celle qu'on nomme l'expérience? Ainsi, chronologiquement, aucune connaissance ne précède en nous l'expérience et c'est avec elle que toutes commencent.

- Welche Sätze wurden wie übersetzt: **Satz-Alignierung**
- Welche Wörter wurden wie übersetzt: **Wort-Alignierung**
+ **Übersetzungswahrscheinlichkeiten**
- Wie sieht eine gute Zielsprache aus: **Sprachmodell**

GERMAN

Einleitung

I. Von dem Unterschiede der reinen und empirischen Erkenntnis

Daß alle unsere Erkenntnis mit der Erfahrung anfangt, daran ist gar kein Zweifel; denn wodurch sollte das Erkenntnisvermögen sonst zur Ausübung erweckt werden, geschähe es nicht durch Gegenstände, die unsere Sinne rühren und teils von selbst Vorstellungen bewirken, teils unsere Verstandstätigkeit in Bewegung bringen, diese zu vergleichen, sie zu verknüpfen oder zu trennen, und so den rohen Stoff sinnlicher Eindrücke zu einer Erkenntnis der Gegenstände zu verarbeiten, die Erfahrung heißt? Der Zeit nach geht also keine Erkenntnis in uns vor der Erfahrung vorher, und mit dieser fängt alle an.

ENGLISH

Introduction

I. Of the difference between Pure and Empirical Knowledge

That all our knowledge begins with experience there can be no doubt. For how is it possible that the faculty of cognition should be awakened into exercise otherwise than by means of objects which affect our senses, and partly of themselves produce representations, partly rouse our powers of understanding into activity, to compare to connect, or to separate these, and so to convert the raw material of our sensuous impressions into a knowledge of objects, which is called experience? In respect of time, therefore, no knowledge of ours is antecedent to experience, but begins with it.

FRENCH

Introduction

I. De la différence de la connaissance pure et de la connaissance empirique.

Que toute notre connaissance commence avec l'expérience, cela ne soulève aucun doute. En effet, par quoi notre pouvoir de connaître pourrait-il être éveillé et mis en action, si ce n'est par des objets qui frappent nos sens et qui, d'une part, produisent par eux-mêmes des représentations et, d'autre part, mettent en mouvement notre faculté intellectuelle, afin qu'elle compare, lie ou sépare ces représentations, et travaille ainsi la matière brute des impressions sensibles pour en tirer une connaissance des objets, celle qu'on nomme l'expérience? Ainsi, chronologiquement, aucune connaissance ne précède en nous l'expérience et c'est avec elle que toutes commencent.

Satz-Alignierung



TRADOS WinAlign - [C:\..._CW_E_H_Ex_wo_mod.rtf : C:\..._CW_D_H_Ex_wo_mod.rtf]

File Edit View Settings Segment Window Help

C:\..._CW_E_H_Ex_wo_mod.rtf
 CycleMatch

This brings you to the calendar where you can begin to enter data.
 After indicating your first cycle date in the dialog box, a calendar appears.
 It indicates the month of your first cycle date.
 The day you indicated as the first cycle day will be marked with a ^g1 in the lower right hand corner and a ^g2.
 If this is not the day you want to indicate as the first cycle day, go to File, Close, and don't save it.
 Then double-click on the CycleMatch icon, enter the correct date, and click on OK.
 You will also see two numbers in each calendar day.
 The number in the upper right

C:\..._CW_D_H_Ex_wo_mod.rtf
 CycleMatch

Dadurch gelangen Sie zum Kalender, in dem alle weiteren Eingaben erfolgen.
 Nachdem Sie Ihren ersten Zyklus in das Dialogfenster eingegeben haben erscheint eine Kalenderansicht, die den Monat Ihres ersten Zyklus wiedergibt.
 Der Tag, den Sie als ersten Zyklustag eingegeben haben, wird in der rechten unteren Ecke des Kalenderfensters durch ^g1 und ^g2 dargestellt.
 Falls dies nicht der erste Tag Ihres ersten Zyklus ist, wählen Sie bitte "Datei", "Schließen", und "Nicht Speichern".
 Führen Sie dann bitte erneut einen Doppelklick auf das CycleMatch Erblen aus und korrigieren das Datum.
 Um die Eingabe zu bestätigen klicken Sie bitte auf OK.
 Es werden jeweils zwei Zahlen pro Kalendertag dargestellt.
 In der rechten Ecke oben sehen

Ready.

		CLASSIC SOUPS		Sm.	Lg.
清 燉 雞	57.	House Chicken Soup (Chicken, Celery, Potato, Onion, Carrot)	1.50	2.75	
雞 飯	58.	Chicken Rice Soup	1.85	3.25	
雞 麵	59.	Chicken Noodle Soup	1.85	3.25	
廣 東 雲 吞	60.	Cantonese Wonton Soup.....	1.50	2.75	
蕃 茄 蛋	61.	Tomato Clear Egg Drop Soup	1.65	2.95	
雲 吞	62.	Regular Wonton Soup	1.10	2.10	
酸 辣	63.	Hot & Sour Soup	1.10	2.10	
蛋	64.	Egg Drop Soup.....	1.10	2.10	
雲 吞	65.	Egg Drop Wonton Mix.....	1.10	2.10	
豆 腐 菜	66.	Tofu Vegetable Soup	NA	3.50	
雞 玉 米	67.	Chicken Corn Cream Soup	NA	3.50	
蟹 肉 玉 米	68.	Crab Meat Corn Cream Soup.....	NA	3.50	
海 鮮	69.	Seafood Soup.....	NA	3.50	

		CLASSIC SOUPS		Sm.	Lg.			
清	燉	雞	湯	57.	House Chicken Soup (Chicken, Celery, Potato, Onion, Carrot)	1.50	2.75	
雞	飯	湯	58.	Chicken Rice Soup	1.85	3.25		
雞	麵	湯	59.	Chicken Noodle Soup	1.85	3.25		
廣	東	雲	吞	60.	Cantonese Wonton Soup.....	1.50	2.75	
蕃	茄	蛋	湯	61.	Tomato Clear Egg Drop Soup	1.65	2.95	
雲	吞	湯	62.	Regular Wonton Soup	1.10	2.10		
酸	辣	湯	63.	Hot & Sour Soup	1.10	2.10		
蛋	花	湯	64.	Egg Drop Soup.....	1.10	2.10		
雲	吞	湯	65.	Egg Drop Wonton Mix	1.10	2.10		
豆	腐	菜	湯	66.	Tofu Vegetable Soup	NA	3.50	
雞	玉	米	湯	67.	Chicken Corn Cream Soup	NA	3.50	
蟹	肉	玉	米	湯	68.	Crab Meat Corn Cream Soup.....	NA	3.50
海	鮮	湯	69.	Seafood Soup.....	NA	3.50		

- Im Wort-Alignierung Modus vieles ist über chinesische Suppe bekannt
- Aber sonst recht wenig ...
- Es wird nur das erkannt, was in den Trainingsdaten gesehen wurde.
- Eigentlich, wie bei Menschen ...
- Ein gemeinsames Thema denn ...
- Kann auf der Basis von aus wort-alignierten Übersetzungen bestehenden Daten ein Übersetzungslexikon gelernt werden?
- Ja, und zwar recht einfach ...

I love the woman.
Ich liebe die Frau.

The man loves the cat.
Der Mann liebt die Katze.

The man loves the woman.
Der Mann liebt die Frau.

I love the man.
Ich liebe den Mann.

They talk to the cat.
Sie sprechen zu der Katze.

They talk to the boy.
Sie sprechen zu dem Jungen.

They talk to the man.
Sie sprechen zu dem Mann.

I talk with the mother.
Ich spreche mit der Mutter.

Aligned Data



I love the woman.
Ich liebe die Frau.

The man loves the cat.
Der Mann liebt die Katze.

The man loves the woman.
Der Mann liebt die Frau.

I love the man.
Ich liebe den Mann.

They talk to the cat.
Sie sprechen zu der Katze.

They talk to the boy.
Sie sprechen zu dem Jungen.

They talk to the man.
Sie sprechen zu dem Mann.

I talk with the mother.
Ich spreche mit der Mutter.

Aligned Data



I	Ich		talk	sprechen	
the	die			spreche	
	dem		to	zu	
	den		boy	Jungen	
	der		cat	Katze	
they	sie		man	Mann	
love(s)	liebe		mother	Mutter	
	lieben		woman	Frau	
	liebt				

Collated Statistics

*I love the woman.
Ich liebe die Frau.*

*The man loves the cat.
Der Mann liebt die Katze.*

*The man loves the woman.
Der Mann liebt die Frau.*

*I love the man.
Ich liebe den Mann.*

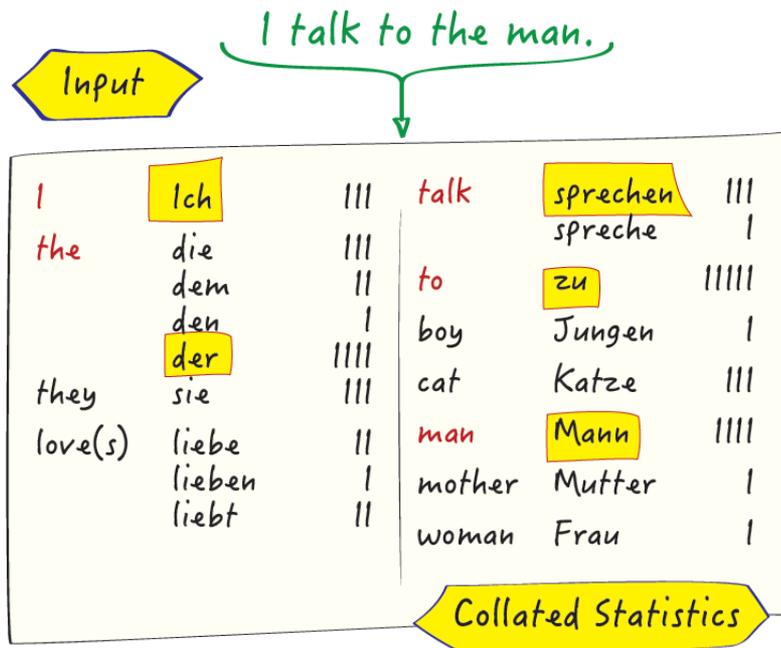
*They talk to the cat.
Sie sprechen zu der Katze.*

*They talk to the boy.
Sie sprechen zu dem Jungen.*

*They talk to the man.
Sie sprechen zu dem Mann.*

*I talk with the mother.
Ich spreche mit der Mutter.*

Aligned Data



*I love the woman.
Ich liebe die Frau.*

*The man loves the cat.
Der Mann liebt die Katze.*

*The man loves the woman.
Der Mann liebt die Frau.*

*I love the man.
Ich liebe den Mann.*

*They talk to the cat.
Sie sprechen zu der Katze.*

*They talk to the boy.
Sie sprechen zu dem Jungen.*

*They talk to the man.
Sie sprechen zu dem Mann.*

*I talk with the mother.
Ich spreche mit der Mutter.*

Aligned Data



Input I talk to the man.

I	Ich		talk	sprechen	
the	die			spreche	
	dem		to	zu	
	den		boy	Jungen	
	der		cat	Katze	
they	sie		man	Mann	
love(s)	liebe		mother	Mutter	
	lieben		woman	Frau	
	liebt				

Collated Statistics

Ich spreche zu der Mann.

Output

I love the woman.
 Ich liebe die Frau.
 The man loves the cat.
 Der Mann liebt die Katze.
 The man loves the woman.
 Der Mann liebt die Frau.
 I love the man.
 Ich liebe den Mann.
 They talk to the cat.
 Sie sprechen zu der Katze.
 They talk to the boy.
 Sie sprechen zu dem Jungen.
 They talk to the man.
 Sie sprechen zu dem Mann.
 I talk with the mother.
 Ich spreche mit der Mutter.



I	talk	to	the	man
Ich	sprechen zu	der	Mann	
3/3	3/4	5/5	4/10	4/4
Ich	spreche	zu	dem	Mann
3/3	1/4	5/5	2/10	4/4

Aligned Data

Auswahlkriterien?



I love the woman.
Ich liebe die Frau.
The man loves the cat.
Der Mann liebt die Katze.
The man loves the woman.
Der Mann liebt die Frau.
I love the man.
Ich liebe den Mann.
They talk to the cat.
Sie sprechen zu der Katze.
They talk to the boy.
Sie sprechen zu dem Jungen.
They talk to the man.
Sie sprechen zu dem Mann.
I talk with the mother.
Ich spreche mit der Mutter.



Aligned Data

Sprachmodell:

- Was ist eine gute Zielsprache?
- Welche Wörter können aufeinander folgen, und welche nicht...? Die Grammatik
- Aus den Daten lernen ...
 - *Ich spreche* is good ...
 - *Ich sprechen* is bad ...
 - *zu dem Mann* is good ...
 - *zu der Mann* is bad ...
- *Ich spreche zu dem Mann* >>
Ich sprechen zu der Mann



I love the woman.
 Ich liebe die Frau.
 The man loves the cat.
 Der Mann liebt die Katze.
 The man loves the woman.
 Der Mann liebt die Frau.
 I love the man.
 Ich liebe den Mann.
 They talk to the cat.
 Sie sprechen zu der Katze.
 They talk to the boy.
 Sie sprechen zu dem Jungen.
 They talk to the man.
 Sie sprechen zu dem Mann.
 I talk with the mother.
 Ich spreche mit der Mutter.

Aligned Data



Input I talk to the man.

I	Ich		talk	sprechen	
the	die			spreche	
	dem		to	zu	
	den		boy	Jungen	
	der		cat	Katze	
they	sie		man	Mann	
love(s)	liebe		mother	Mutter	
	lieben		woman	Frau	
	liebt				

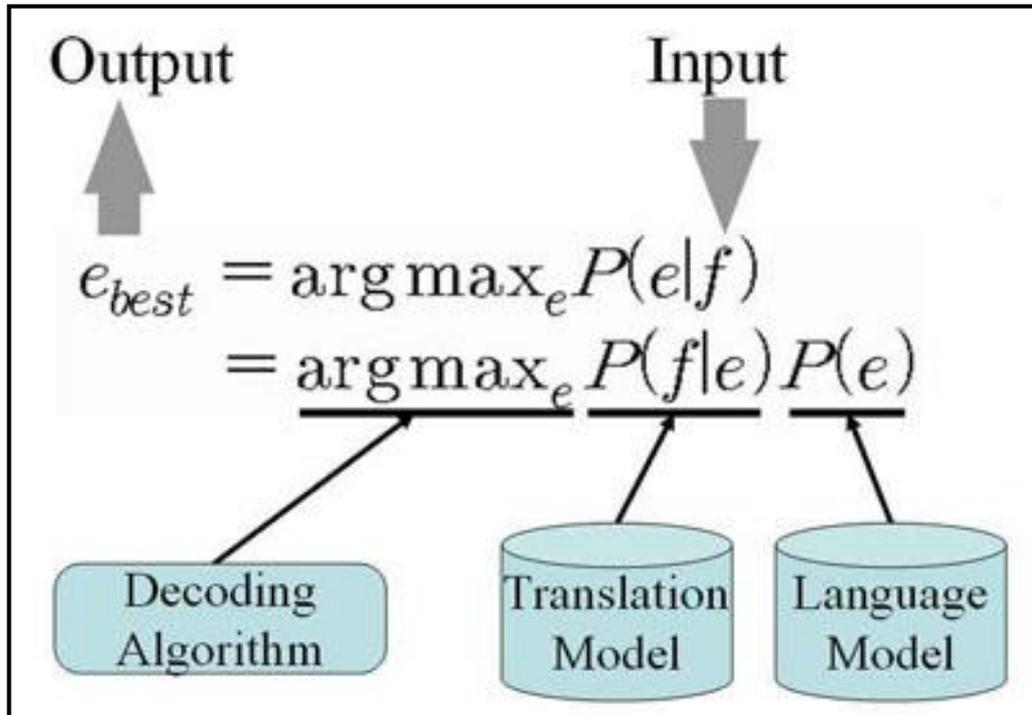
Collated Statistics

+

Language Model

Ich spreche zu dem Mann.

Output



- Heute kein Mathematikunterricht!
- Sondern:
- Die Geschichte der statistischen MÜ in Bildern ...
- Es dreht sich einzig und allein um **Daten**
...

- Bis jetzt: nur einzelne Wörter übersetzt
- Kontext, wie Kongruenz, fehlt (*zu dem Mann ...*) usw.
- Bis zu einem gewissen Grad “repariert” mit Hilfe des Sprachmodells
- Ein besserer Ansatz:
- Nicht nur einzelne Wörter, sondern auch Phrasen übersetzen:

*the man : der Mann
to the man : zu dem Mann
I talk : Ich spreche*

I love the woman.
Ich liebe die Frau.

The man loves the cat.
Der Mann liebt die Katze.

The man loves the woman.
Der Mann liebt die Frau.

I love the man.
Ich liebe den Mann.

They talk to the cat.
Sie sprechen zu der Katze.

They talk to the boy.
Sie sprechen zu dem Jungen.

They talk to the man.
Sie sprechen zu dem Mann.

I talk with the mother.
Ich spreche mit der Mutter.

Aligned Data



Input

I talk to the man.

I	Ich		talk	sprechen	
the	die			spreche	
	dem		to	zu	
	den		boy	Jungen	
	der		cat	Katze	
they	sie		man	Mann	
love(s)	liebe		mother	Mutter	
	lieben		woman	Frau	
	liebt				

Collated Statistics

Ich spreche zu der Mann.

Output

I love the woman.
 Ich liebe die Frau.
 The man loves the cat.
 Der Mann liebt die Katze.
 The man loves the woman.
 Der Mann liebt die Frau.
 I love the man.
 Ich liebe den Mann.
 They talk to the cat.
 Sie sprechen zu der Katze.
 They talk to the boy.
 Sie sprechen zu dem Jungen.
 They talk to the man.
 Sie sprechen zu dem Mann.
 I talk with the mother.
 Ich spreche mit der Mutter.

Aligned Data



Input

I talk to the man.

I love	Ich liebe	11
__ loves	__ liebt	11
they talk	sie sprechen	111
I talk	ich spreche	1
the man	der man	11
the woman	die Frau	1
the cat	die Katze	1
to the cat	zu der Katze	1
to the boy	zu dem Jungen	1
to the man	zu dem Mann	1
with the mother	mit der Mutter	1

I love the woman.
Ich liebe die Frau.
The man loves the cat.
Der Mann liebt die Katze.
The man loves the woman.
Der Mann liebt die Frau.
I love the man.
Ich liebe den Mann.
They talk to the cat.
Sie sprechen zu der Katze.
They talk to the boy.
Sie sprechen zu dem Jungen.
They talk to the man.
Sie sprechen zu dem Mann.
I talk with the mother.
Ich spreche mit der Mutter.

Aligned Data



Input

I talk to the man.

I love	Ich liebe	11
__ loves	__ liebt	11
they talk	sie sprechen	111
I talk	ich spreche	1
the man	der man	11
the woman	die Frau	1
the cat	die Katze	1
to the cat	zu der Katze	1
to the boy	zu dem Jungen	1
to the man	zu dem Mann	1
with the mother	mit der Mutter	1

Ich spreche zu dem Mann.

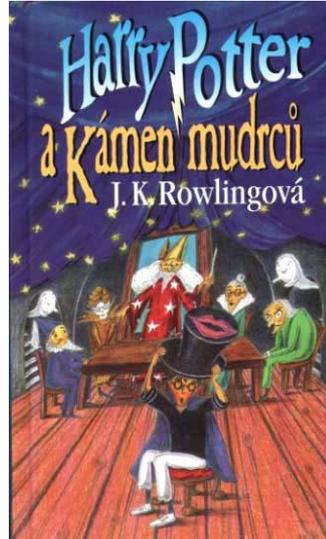
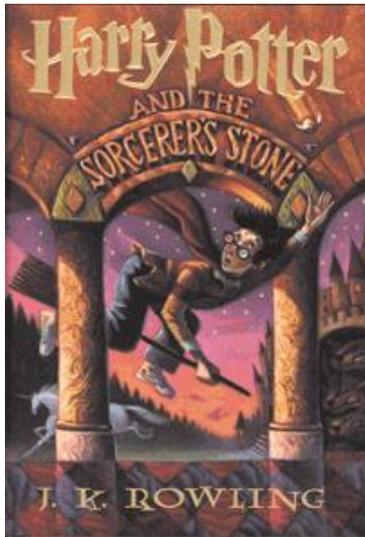
Output

- Viel besser als wortbasierte SMÜ!
- Standard Technologie: Google, Microsoft, Baidu, globale Lokalisierungs- und Übersetzungsindustrie
- Moses Open Source PB-SMÜ
- Am meisten verwendetes System für SMÜ
- Forschung auch von der EC finanziert
- Eingesetzt bei dem Direktorat EC DGT's MT@EC

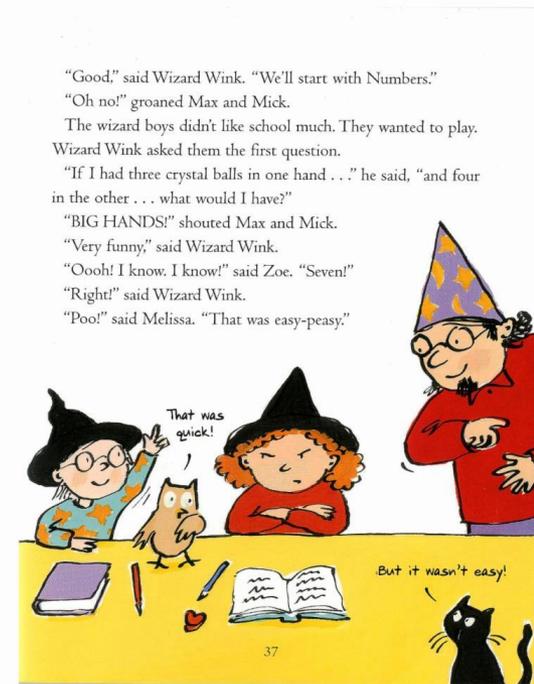
MOSES  CORE

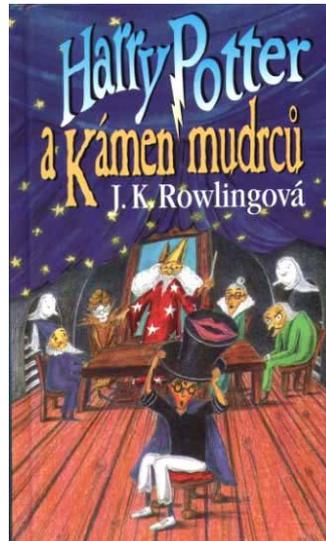
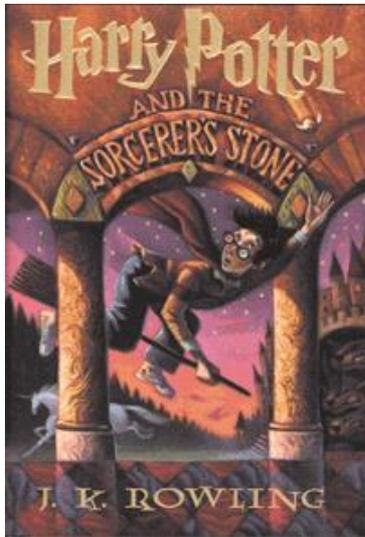


- Bei der Statistischen Maschinellen Übersetzung dreht sich alles um Daten
- SMÜ lernt das Übersetzen aus den Daten
- Daten
 - Übersetzungen (zweisprachige Daten)
 - Einzelsprachliche Daten (Text in der Zielsprache)
 - Wörterbücher, Terminologie, Ontologien, Eigennamen
- Genauso wie bei den Menschen, die SMÜ ist gut bei dem was sie gelernt hat.



MOSES  CORE





Protect - Personal Information CIVMEANS7
Legal Aid Agency
Financial Assessment for Family Mediation
Provider reference/case code: MED12/1GBHST/1/451
This form must be completed in ink.

Applicant's Details
Surname: Mr/Ms/Miss/Ms _____ First name(s) _____
Surname at birth if different: _____ Date of birth: JJ ____
Address: _____ Postcode: _____
National Insurance number: L _____
Job: _____

Financial Eligibility

- The client has a partner whose means are to be aggregated:
 - Yes Please provide details of both client's and partner's means.
 - No Please provide details of both client's means only.
- The case is about ownership or possession of assets and / or financial provision:
 - Yes Go to question 3.
 - No Go directly to Part B Capital.
- The client's assets (held in sole name or jointly held) have been claimed by the opponent:
 - Yes Please complete Part A Capital - Subject matter of dispute.
 - No Go directly to Part B Capital.

The subject matter of dispute disregard only applies to assets that are specially claimed by the opponent. All assets that have not been specifically claimed by the opponent must be included in Part B Capital.

CIVMEANS7 Page 1 Version 8 April 2013 © Crown Copyright

MOSES  CORE



- CEF.AT braucht die richtigen Daten
- Nationale Regierungen, öffentliche Verwaltungen, öffentliche Dienste, NRO/NGOs
- CEF bietet Diensten für multilinguale Interaktion mit den nationalen Bürgern, EU Bürgern und anderen Nutzern von öffentlichen Verwaltungen.

- Helfen Sie uns, CEF.AT zum Erfolg zu führen
 - Dienste für Europäische Bürger
 - Dienste für Sie
 - Unterstützung von Mehrsprachigkeit
- Helfen Sie, die richtigen Daten zu finden
- Unsere Sprachen zu unterstützen heißt Europa zu unterstützen, und umgekehrt