

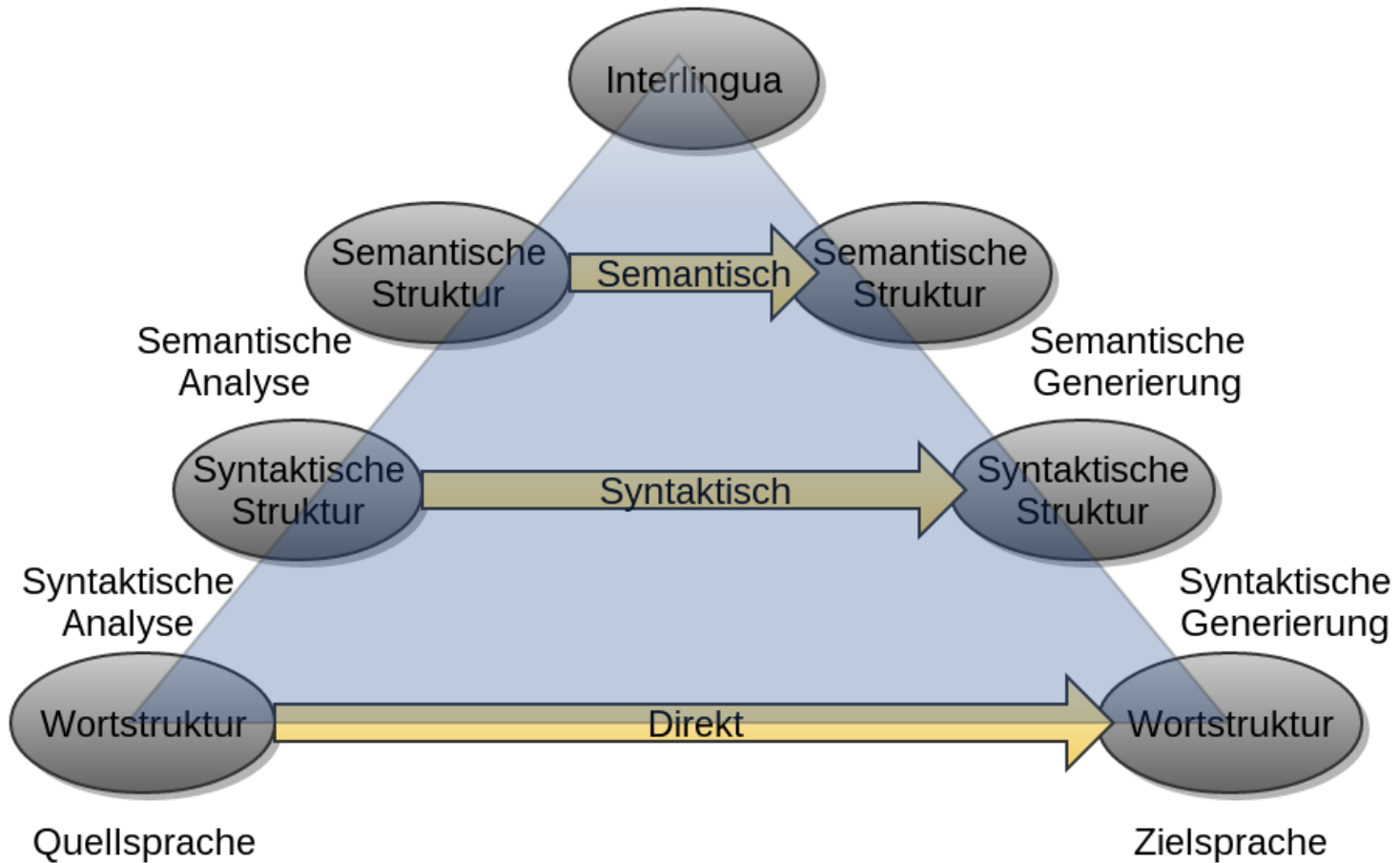
“Wie funktioniert maschinelle Übersetzung?”

Bartholomäus Wloka (Universität Wien)

Basierend auf der Präsentation von Prof. Josef van Genabith
(Deutsches Forschungszentrum für Künstliche Intelligenz)

Überblick:

- Was macht maschinelle Übersetzung so schwierig?
- FAHQMT
- Warum MÜ: Datenmenge, Qualität und Kosten?
- MÜ + **professionelle** Übersetzer = Qualität
 - CAT Tools
 - MAHT
 - HAMT





Überblick:

- Wie funktioniert die moderne statistisch-basierte MÜ?
- Es geht vor allem um Daten
- Und um die richtige Art von Daten
- Vor allem: Parallele Korpora und Sprachmodelle
- Am besten Domänenspezifisch

- 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.

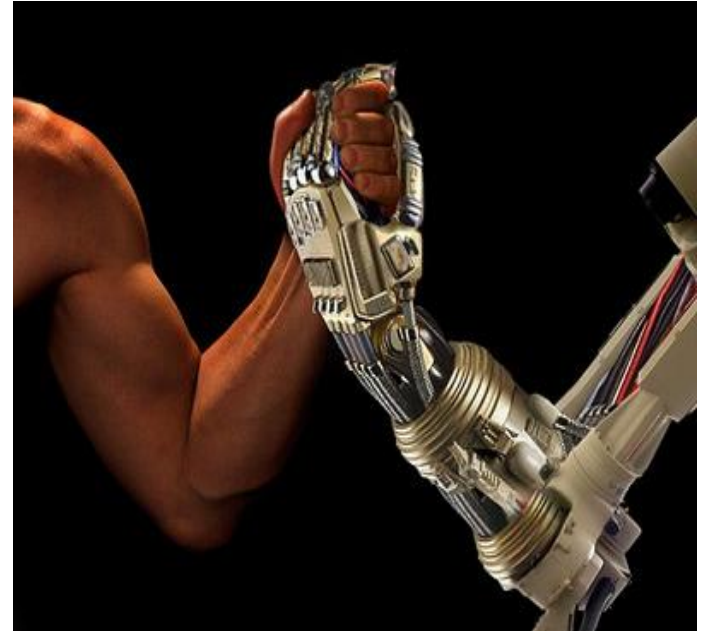


Image: <http://workingtropes.lmc.gatech.edu/wiki/index.php/File:Man-vs-machine.jpg>
License: CC BY-NC-SA 3.0



- Sprache ist komplex, Übersetzung noch komplexer
- Wir können sie nicht genau berechnen
- Verschiedene MÜ Methoden erforscht
- Hybride Methoden
- Maschinelles Lernen
 - Aus **Daten** lernen \Rightarrow zentrale Rolle von Daten
 - Grobe Lösung
 - schafft ersten Überblick, MAHT
 - *Post-editing* durch professionelle Übersetzer





Schlagzeilen:

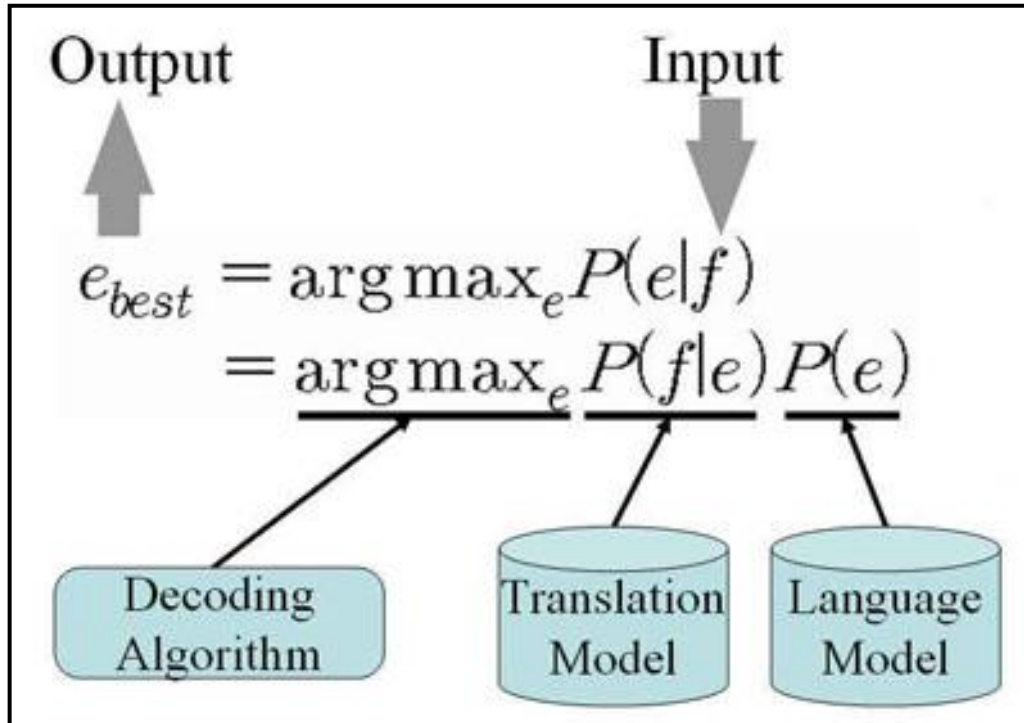
1. Minister Accused Of Having 8 Wives In Jail
2. Teacher Strikes Idle Kids
3. Miners refuse to work after death
4. Local High School Dropouts Cut in Half
5. Red Tape Holds Up New Bridges
6. Clinton Wins on Budget, but More Lies Ahead
7. Police: Crack Found in Man's Buttocks

Lecture notes:
Chris Manning

Was macht die MÜ so schwierig?



- Sprichwörter:
 - The early bird catches the worm
 - Morgenstund hat Gold im Mund
- Satzstellung
 - The **German chancellor** Angela Merkel **will make** an announcement on Thursday.
 - Angela Merkel **wird** am Donnerstag eine Ankündigung **machen**
- Polysemie
 - Der Angeklagte hat gestanden, jetzt muss er sitzen.



- 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 **adäquate** Daten wie möglich

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

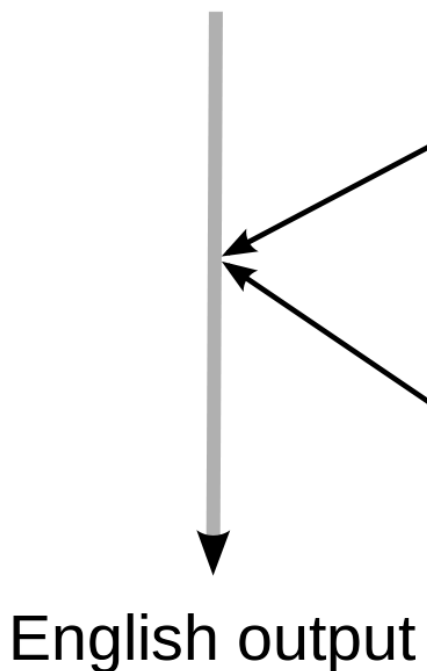
Introduction

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.



似乎格式有問題



**translation
model**

**language
model**

parallel corpus

网站资讯分析网数
据显示的主域名为
全世界访问量最高
的站点除此之外搜
索在其他国家或地
区域名下的多个站
点等等及旗下的等

The corporation has been estim
to run more than one million pag
in data centers around the world
to process over one billion search
requests and about twenty-four i
of user-generated data each dat
December 2012 Alexa listed as

monolingual corpus

started functioning in 1928 and established the tradition of
large exhibitions and trade fairs held in Brno, and nowadays
also ranks among the sights of the city. Brno is also
known for hosting big motorbike and other races on the
Masaryk Circuit, a tradition established in 1930 in which
the Road Racing World Championship Grand Prix is
one of the most prestigious races. Another notable cultural
tradition is an international fireworks competition.

- Welche Sätze wurden wie übersetzt: **Satz-Alignierung**
- Welche Wörter wurden wie übersetzt: **WSD + 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.

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.

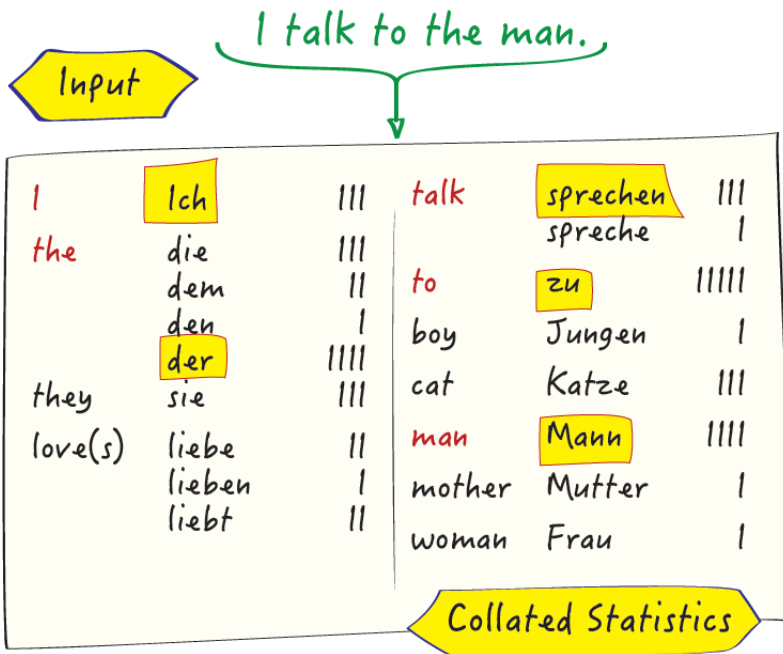
They talk to the cat.
Sie sprechen zu der Katze.

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.

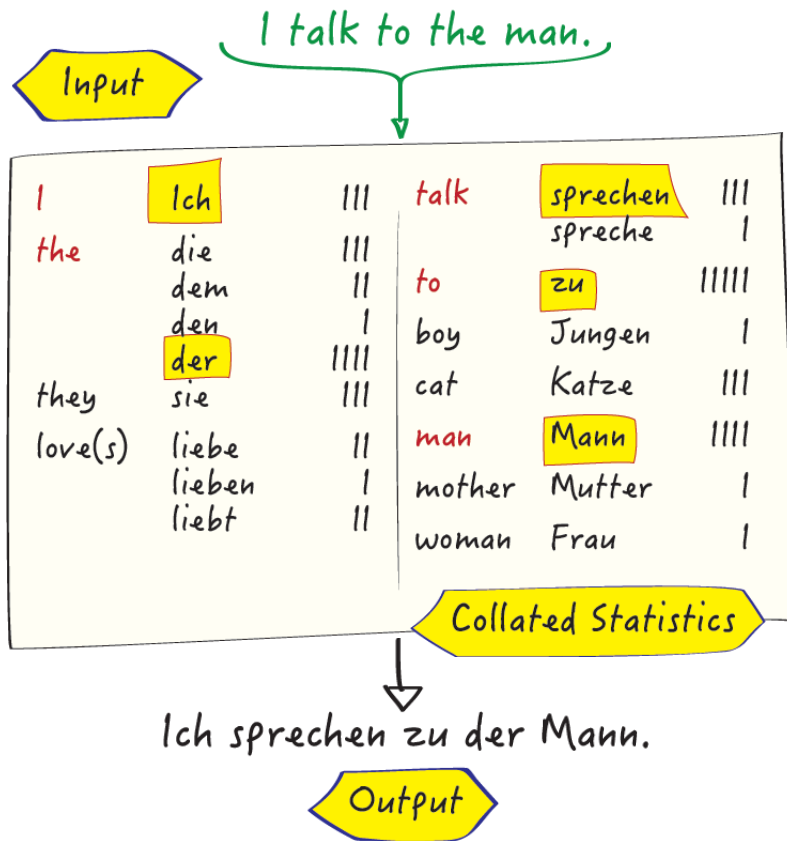
They talk to the cat.
Sie sprechen zu der Katze.

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



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

Auswahlkriterien?

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

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

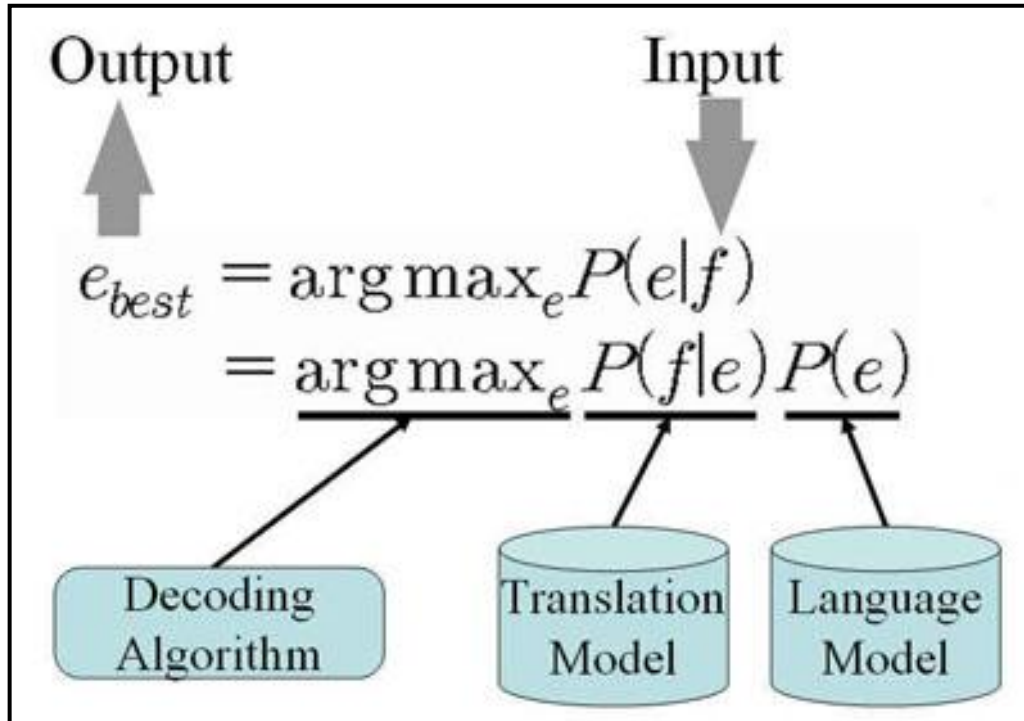
Collated Statistics

+

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.

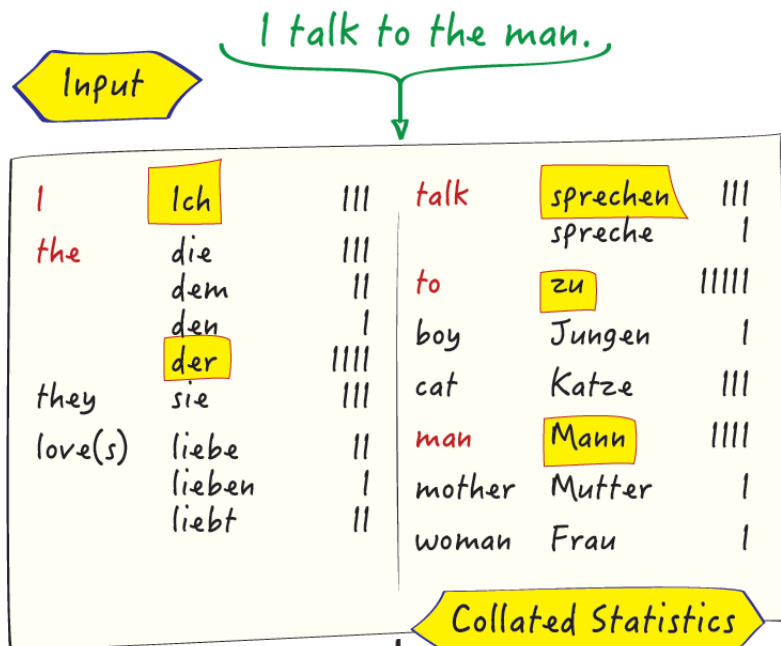
They talk to the cat.
Sie sprechen zu der Katze.

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



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



- 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

- Ziel: Der “Multilingual Digital Single Market” (mSDM):
 - Keine sprachliche Barriere
 - Freier Verkehr von Leuten, Information, Dienste, Kultur, und Kommerz
- Ziel: CEF.AT:
 - Unterstützung von öffentlichen Diensten, Regierungen, Verwaltung, **NGOs** in ganz Europa



- 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)
- Die Qualität der SMÜ hängt vom „Gelernten“ ab
- Nachbearbeitung möglich mit
 - Lexikalische Ressourcen, Terminologie, Ontologien, Eigennamen



- 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 mit die richtigen Daten zu finden bzw. zur Verfügung zu stellen