

Kako zaista radi automatizirano prevođenje?

Marko Tadić
Sveučilište u Zagrebu, Filozofski fakultet

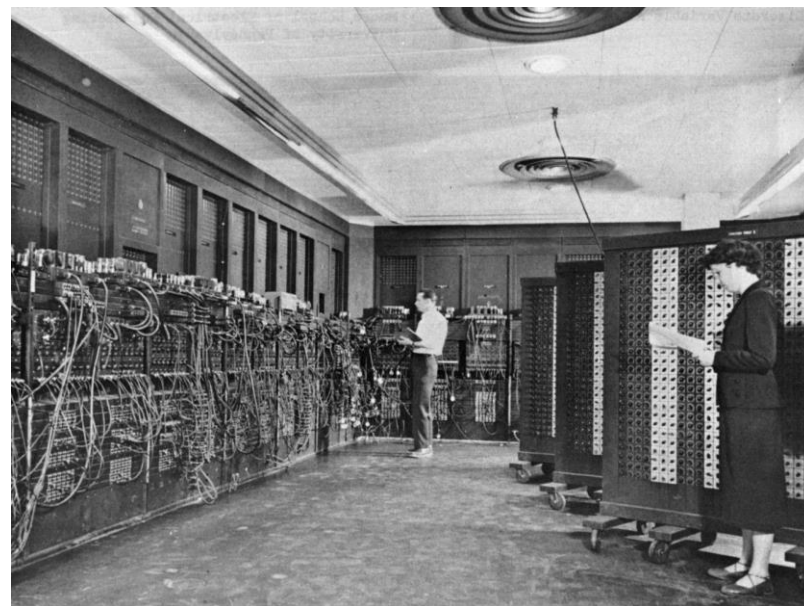
Zahvala za prilagođenu građu koja se pojavljivala na
MT Maratonima, radionicama MT-a itd.



Plan izlaganja

- Zašto strojno prevođenje: obujam, kakvoća, cijena?
- Zašto je MT teško?
- MT + ljudski prevoditelji = visoka kakvoća prijevoda
- Kako radi suvremeni statistički MT?
- Sve je u podacima!
- I to u pravoj vrsti podataka!

- Europa je višejezična!
- 24 službena jezika,
24+2 jezika u CEF-u
- Toliko toga za prevesti!
- Cijena prevođenja?
- Može li MT pomoći?
- A što je s kakvoćom prijevoda?



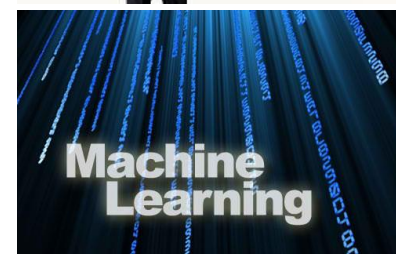
Slika: <https://en.wikipedia.org/wiki/ENIAC#/media/File:Eniac.jpg>
Licencija: javna domena

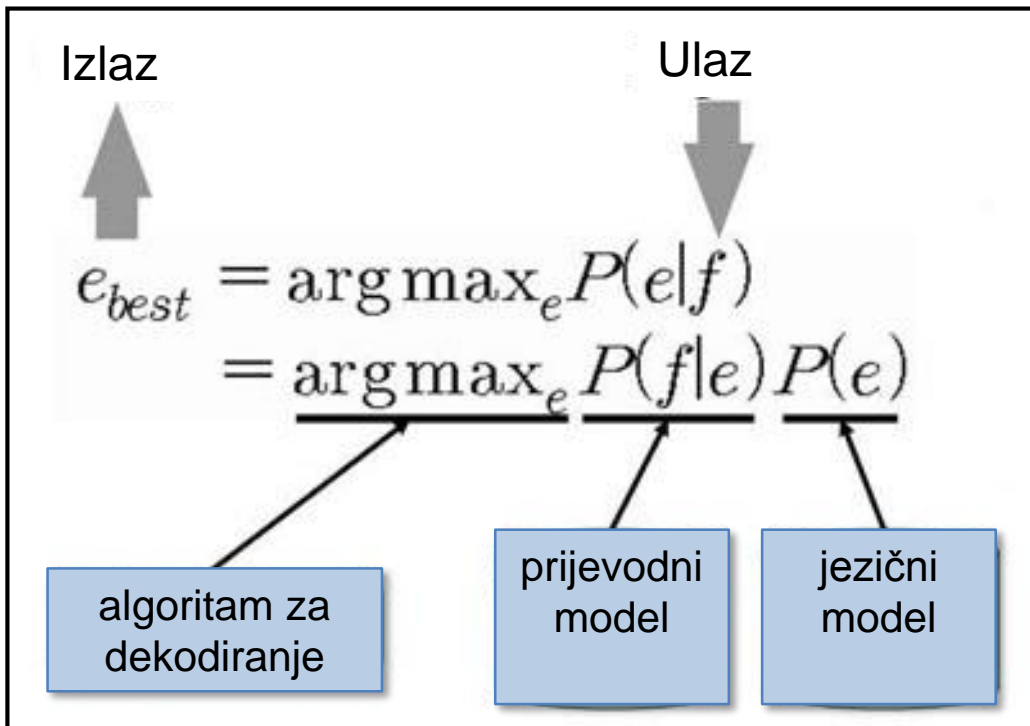
- prirodni su jezici:
 - elegantni
 - učinkoviti
 - prilagodljivi
 - složeni
- jedna riječ/rečenica može imati više značenja
- isto se značenje može prenijeti na razne načine
- značenje riječi ovisi o okolini
- doslovna i prenesena značenja (metafore)
- jezik je blisko vezan s kulturom (različiti načini poimanja iste pojave)
- strukturne osobine prirodnih jezika
 - red riječi u rečenici
 - morfologija
 - ...



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

- prevođenje s jezika na jezik je složeno
- ne možemo ga egzaktno izračunati
- pokušali smo: MT utemeljeno na pravilima + jezične tehnologije ...
- Što sad pokušavamo?
- strojno učenje (*Machine Learning*, ML)
 - stroj uči iz **podataka** \Rightarrow podatci su važni
 - približna rješenja \Rightarrow nisu savršena, potrebna je pomoć
 - profesionalnih prevoditelja
 - naknadnoga uređivanja (*post-editing*)
 - automatiziranoga prevođenja (\neq automatsko prevođenje)





- Ne bojte se!
Niti riječi danas o
matematici.
- umjesto toga:
- priča o statističkome
strojnom prevođenju u
slikama ...
- sve je u **podacima** ...



Statističko strojno prevođenje uči iz dvije vrste podataka:

- ljudski prijevodi
- tekstovi na ciljnome jeziku
- Što više podataka, to bolje!
- Ali: i prave vrste podataka!

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 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.



- koja je rečenica prevedena kojom rečenicom: **sravnjivanje rečenica**
- koja je riječ prevedena kojom riječju: **sravnjivanje riječi + prijevodna vjerojatnost**
- kako izgleda tekst na dobrom ciljnem jeziku: **jezični model**

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

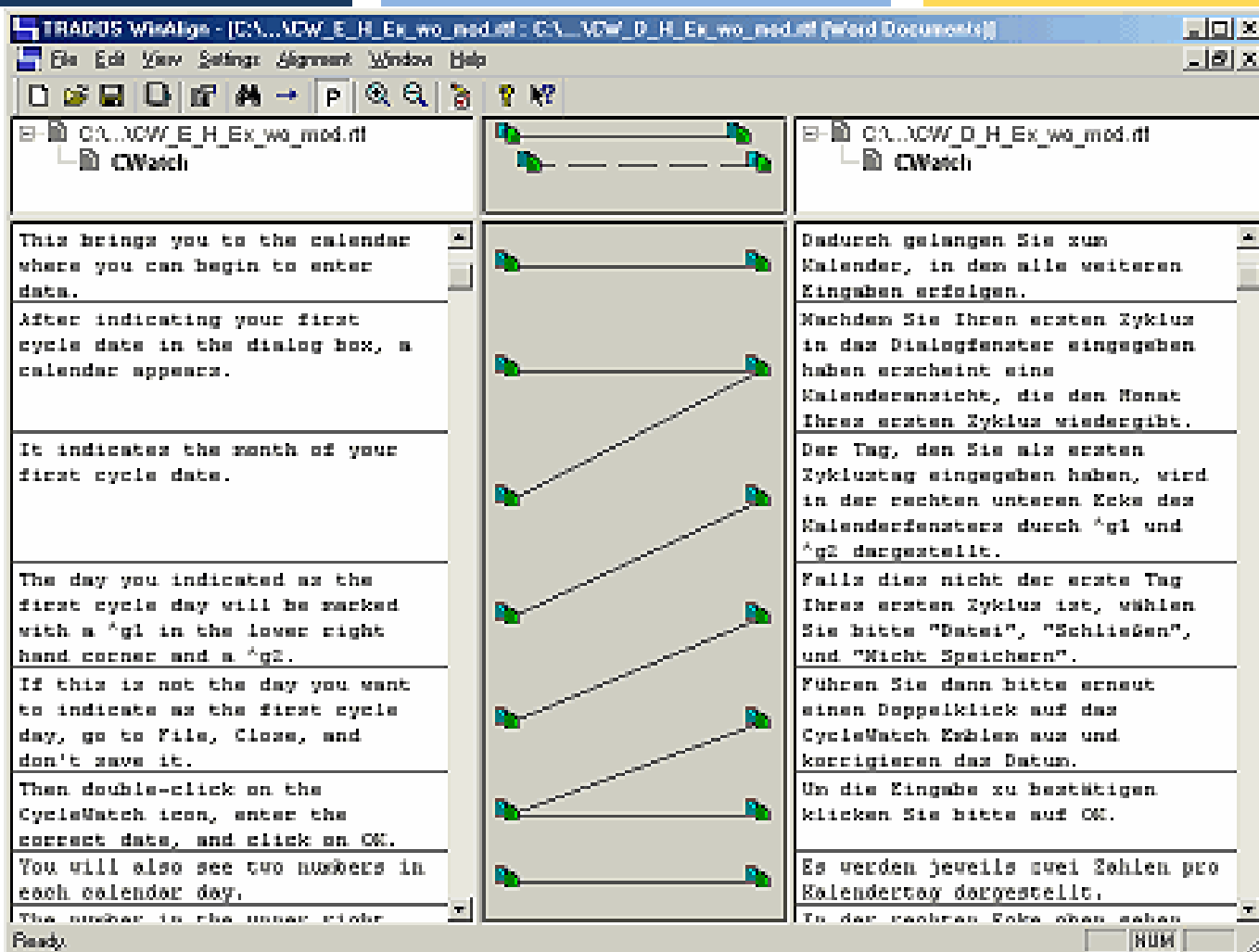
That all our knowledge begins with experience there can be no doubt. For how is it possible that the faculty of cognition should be awakened into exercise otherwise than by means of objects which affect our senses, and partly of themselves produce representations, partly rouse our powers of understanding into activity, to compare 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.



TRADOS WinAlign - [C:\...NOW_E_H_Ex_wo_mod.tif : C:\...NOW_D_H_Ex_wo_mod.tif (InWord Documents)]

File Edit View Settings Segment Window Help

C:\...NOW_E_H_Ex_wo_mod.tif
 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:\...NOW_D_H_Ex_wo_mod.tif
 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 Icon 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. NUM

		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	61.	Tomato Clear Egg Drop Soup	1.65	2.95		
雲吞	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	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		



- sustav za sravnjivanje riječi zna mnogo o nazivima kineskih juha
- ali ne zna ništa drugo
- zna samo ono što je već vidio u podacima za učenje
- kao i ljudi
- Ukoliko imamo skup sravnjenih riječi dvaju jezika, jesmo li time dobili prijevodni rječnik?
- Jesmo. I to izgleda stvarno jednostavno ...

I love the boy.
J'aime le garçon.
I love the dog.
J'aime le chien.
They love the dog.
Ils aiment le chien.
They talk to the girl.
Ils parlent à la fille.
They talk to the dog.
Ils parlent au chien.
I talk to the mother.
Je parle à la mère.

Aligned Data

I love the boy.
 J'aime le garçon.

I love the dog.
 J'aime le chien.

They love the dog.
 Ils aiment le chien.

They talk to the girl.
 Ils parlent à la fille.

They talk to the dog.
 Ils parlent au chien.

I talk to the mother.
 Je parle à la mère.

Aligned Data



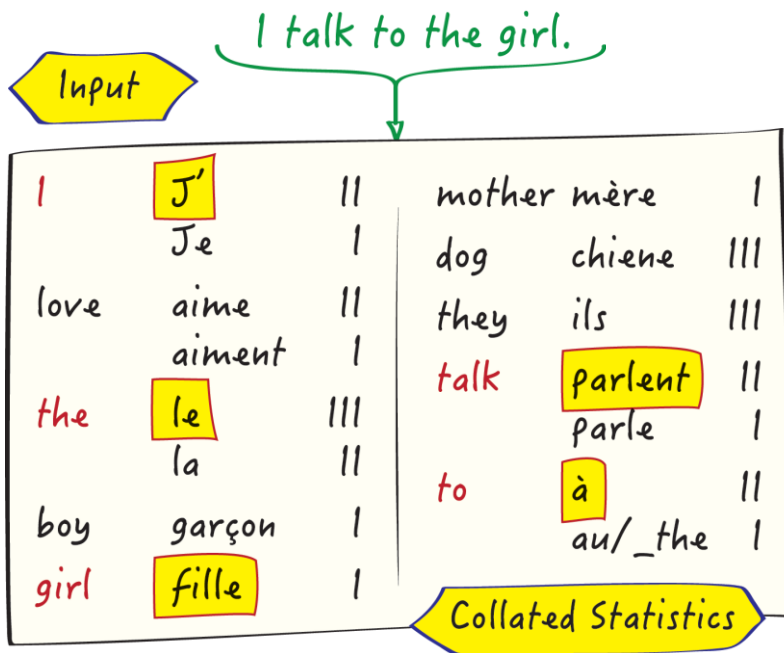
I	J'		mother	mère	
	Je		dog	chiene	
love	aime		they	ils	
	aiment		talk	parlent	
the	le			parle	
	la		to	à	
boy	garçon			au/_the	
girl	fille				

Collated Statistics



I love the boy.
J'aime le garçon.
I love the dog.
J'aime le chien.
They love the dog.
Ils aiment le chien.
They talk to the girl.
Ils parlent à la fille.
They talk to the dog.
Ils parlent au chien.
I talk to the mother.
Je parle à la mère.

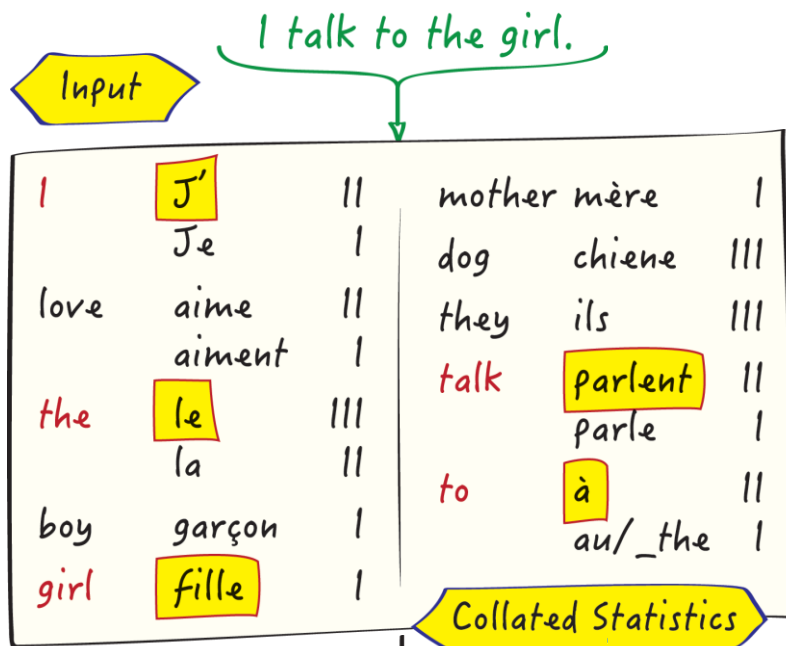
Aligned Data





I love the boy.
 J'aime le garçon.
 I love the dog.
 J'aime le chien.
 They love the dog.
 Ils aiment le chien.
 They talk to the girl.
 Ils parlent à la fille.
 They talk to the dog.
 Ils parlent au chien.
 I talk to the mother.
 Je parle à la mère.

Aligned Data



J'parlent à le fille.

Output



I love the boy.
 J'aime le garçon.
 I love the dog.
 J'aime le chien.
 They love the dog.
 Ils aiment le chien.
 They talk to the girl.
 Ils parlent à la fille.
 They talk to the dog.
 Ils parlent au chien.
 I talk to the mother.
 Je parle à la mère.



Aligned Data

I	talk	to	the	girl
J'	parlent	au	le	fille
2/3	2/3	2/3	3/5	1/1
Je	parle	à	la	fille
1/3	1/3	1/3	2/5	1/1

Kako odabrati prijevod?



I love the boy.
J'aime le garçon.
I love the dog.
J'aime le chien.
They love the dog.
Ils aiment le chien.
They talk to the girl.
Ils parlent à la fille.
They talk to the dog.
Ils parlent au chien.
I talk to the mother.
Je parle à la mère.



Aligned Data

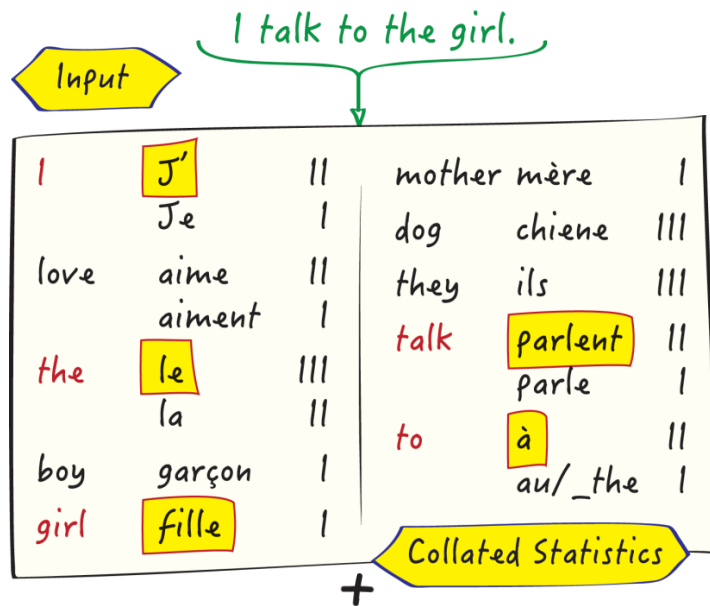
Jezični model:

- Što je ispravno na ciljnome jeziku?
- Koje riječi mogu slijediti nakon kojih riječi, a koje ne mogu... gramatika?
- naučeno iz podataka...
- *Je parle* ispravno ...
- *J' parlent* neispravno ...
- *la fille* ispravno ...
- *le fille* neispravno ...
- *Je parle à la fille* >> *J' parlent à le fille*



I love the boy.
J'aime le garçon.
I love the dog.
J'aime le chien.
They love the dog.
Ils aiment le chien.
They talk to the girl.
Ils parlent à la fille.
They talk to the dog.
Ils parlent au chien.
I talk to the mother.
Je parle à la mère.

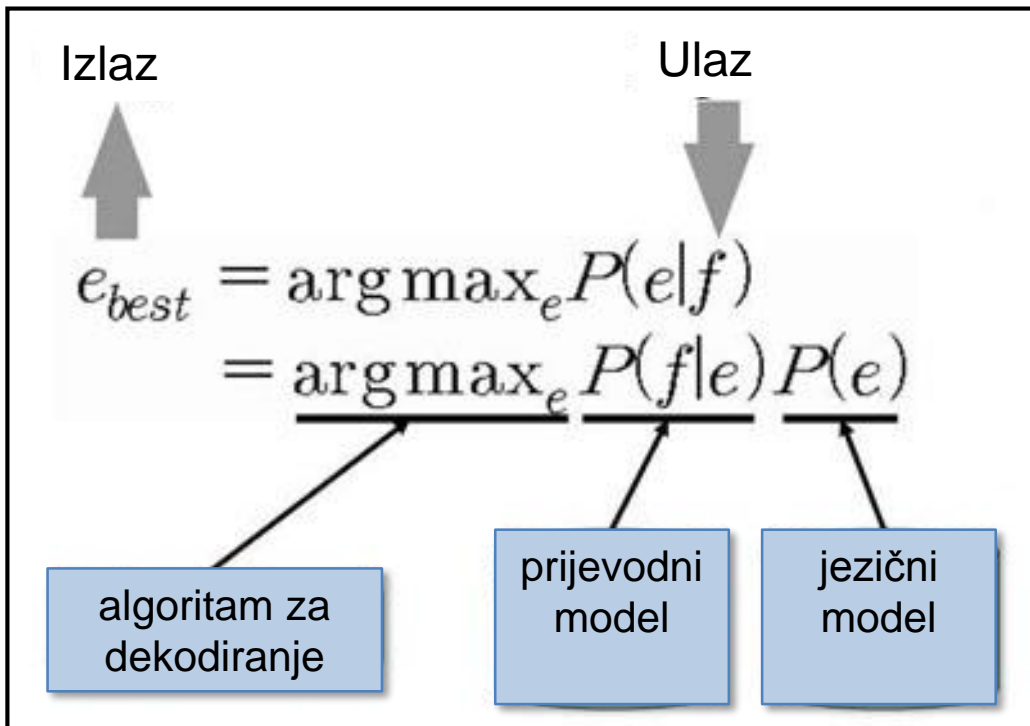
Aligned Data



Jezični model

J'parlent à le fille.

Output



- Ne bojte se!
Niti riječi danas o
matematici.
- umjesto toga:
- priča o statističkome
strojnom prevođenju u
slikama ...
- Sve je u **podacima** ...

- do sad smo prevodili samo pojedinačne riječi
- gubi se okolina: npr. sročnost prema rodu (*le fille* ...) itd.
- do neke to mjere može “popraviti” jezični model

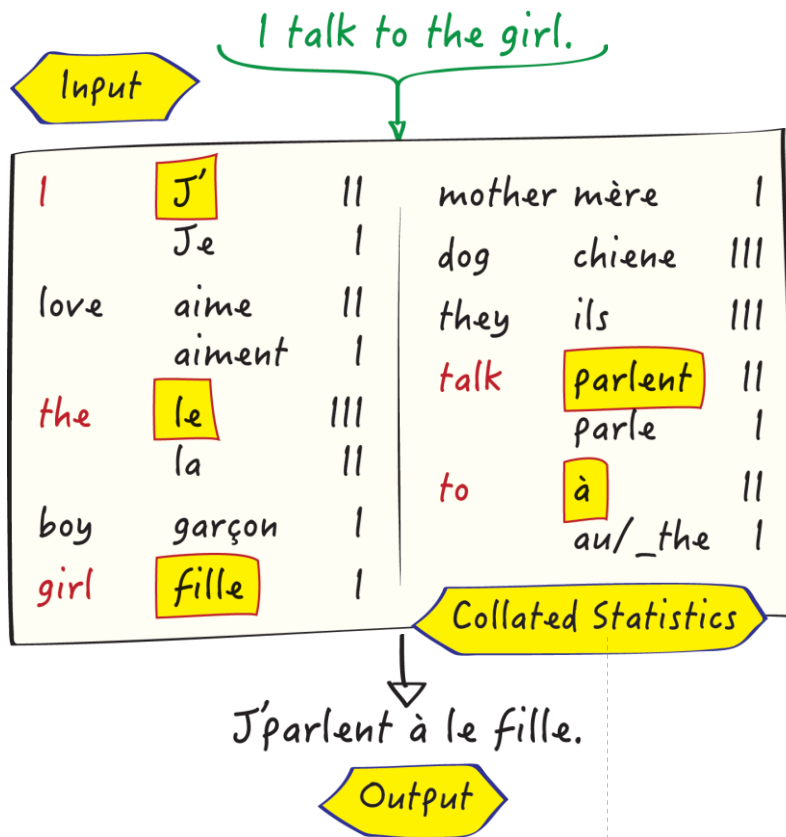
- bolji prijevodni modeli:
- nisu izrađeni samo na prijevodima pojedinačnih riječi
- nego čitavih fraza:

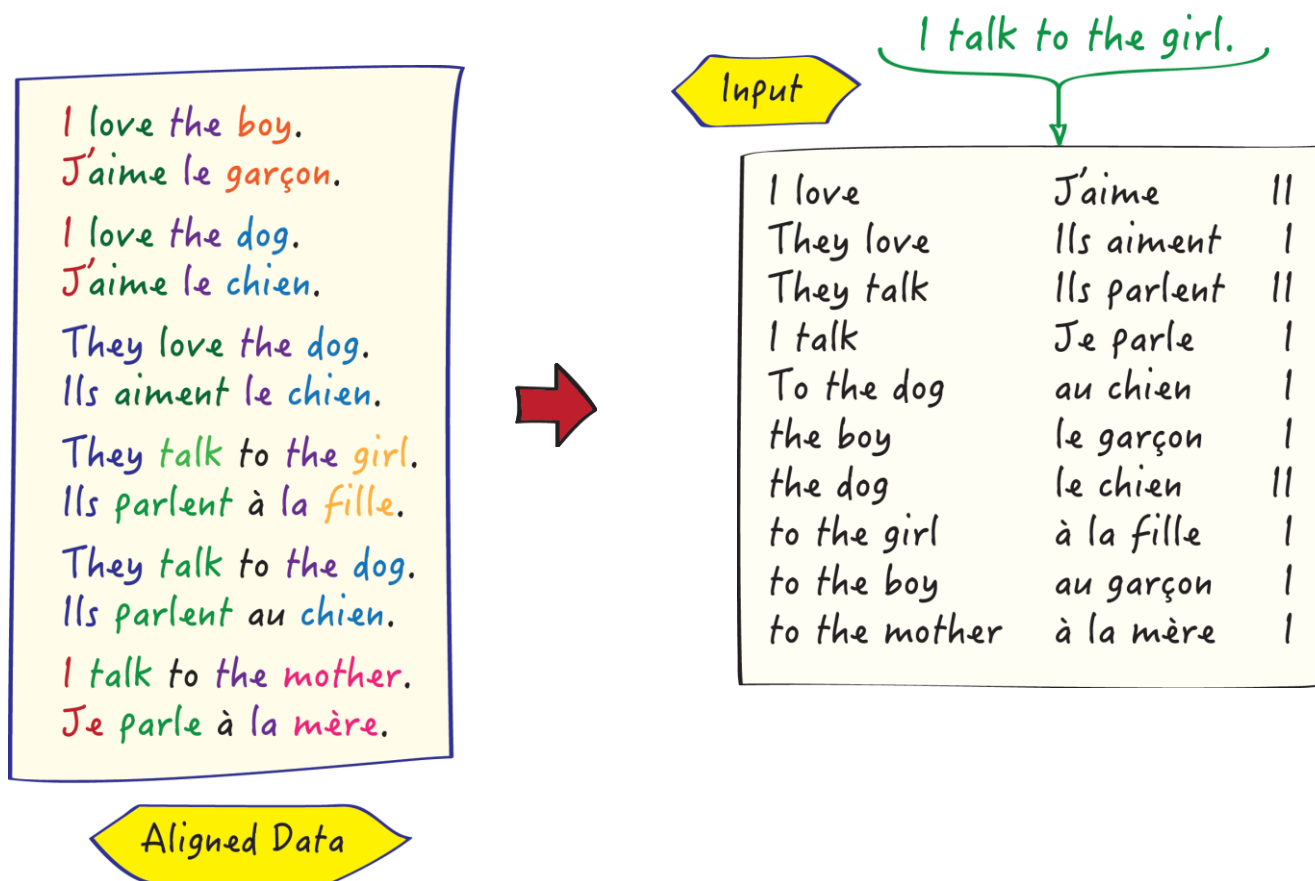
– *the girl : la fille*
– *to the girl : a la fille*
– *I talk : Je parle*

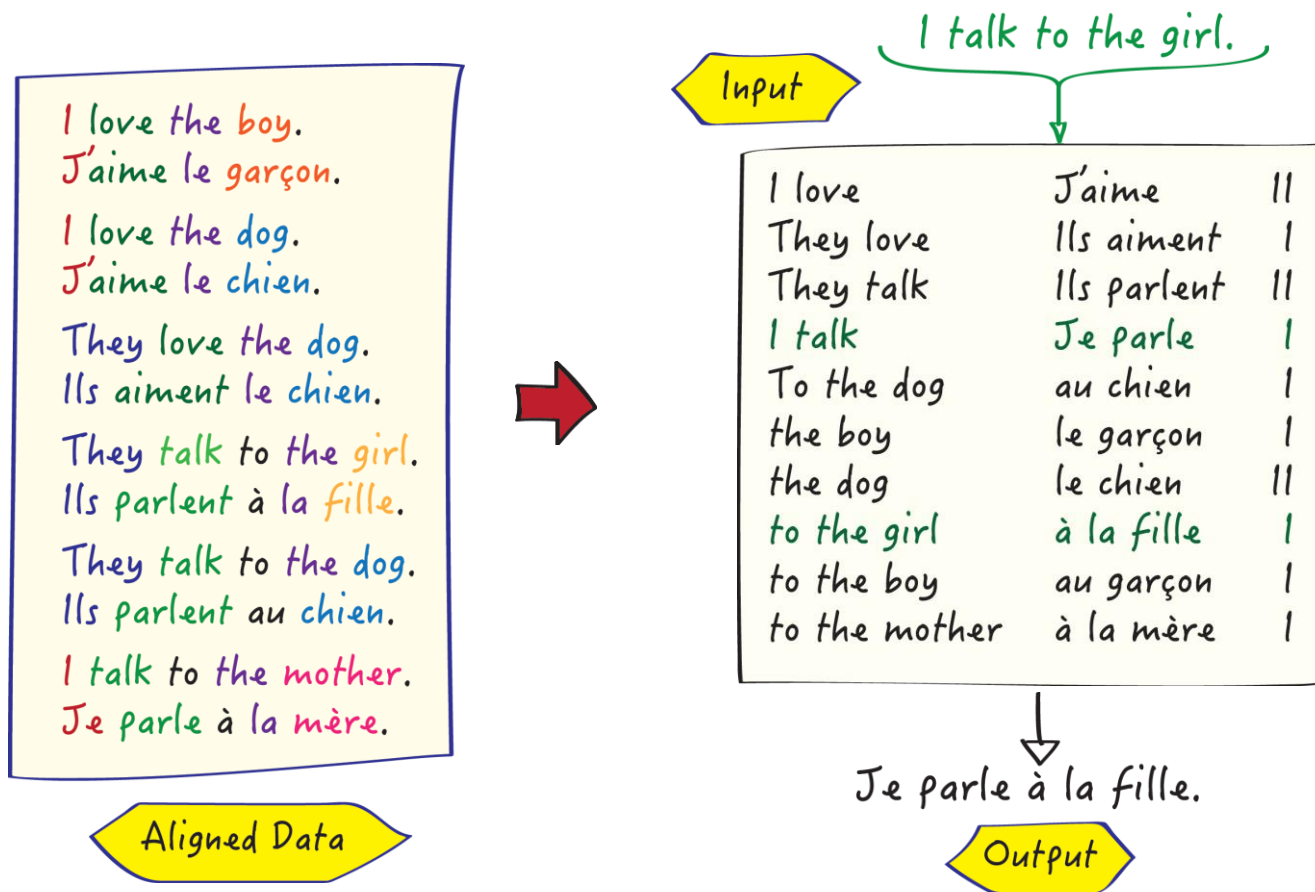


I love the boy.
J'aime le garçon.
I love the dog.
J'aime le chien.
They love the dog.
Ils aiment le chien.
They talk to the girl.
Ils parlent à la fille.
They talk to the dog.
Ils parlent au chien.
I talk to the mother.
Je parle à la mère.

Aligned Data





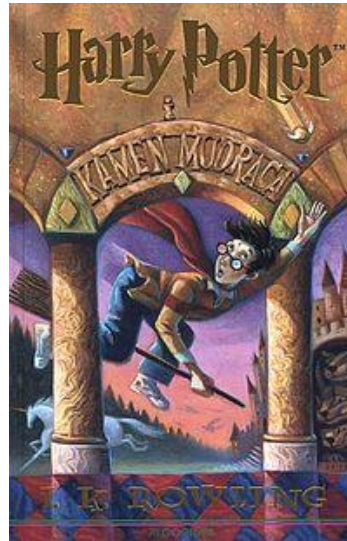
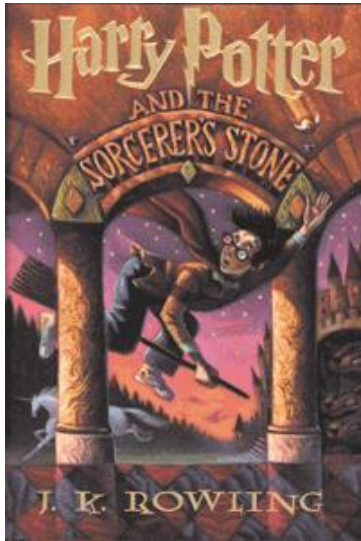


- znatno bolje nego statističko prevođenje pojedinih riječi
- standardna tehnologija: Google, Microsoft, Baidu, globalna lokalizacijska i prijevodna industrija



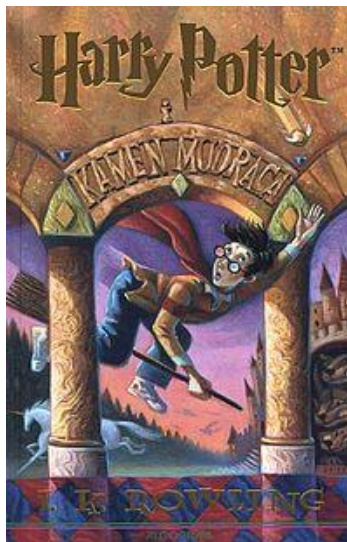
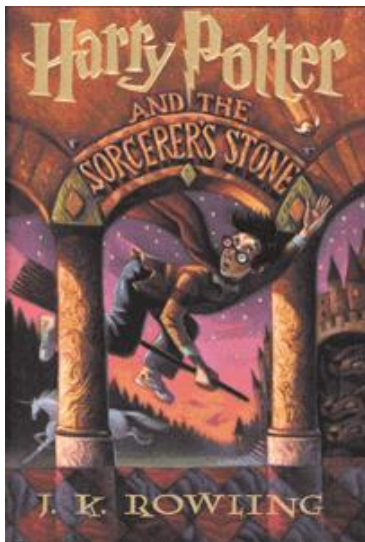
- Moses: sustav otvorenoga koda za SMT temeljen na frazama
- najrašireniji SMT sustav u svijetu
- nastao na temelju istraživanja koje je poduprla EK
- koristi ga EK u DGT-jevoj platformi MT@EC

- u statističkome strojnom prevođenju sve ovisi o podacima
- sustav za SMT iz podataka uči kako prevoditi
- podatci
 - prijevodi (dvojezični ili višejezični tekstovi, izvornici i prijevodi)
 - jednojezični podatci (tekstovi na ciljnome jeziku)
 - rječnici, zbirke termina, ontologije, popisi imena
- poput ljudi, SMT je dobar u onome što ga se naučilo



MOSES  CORE





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: __/__/____
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 disregarded 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

- CEF.AT treba pravu vrstu podataka
- nacionalne vlade, javna administracija, javne službe, NVO-i
- CEF nudi usluge koje će omogućuju višejezično sudjelovanje pojedinačnih državljana kao i tijela javne vlasti država-članica EU-a

- Pomozite nam pretvoriti CEF.AT u uspješnu
 - uslugu za državljane država-članica
 - uslugu za vas
 - potporu višejezičnosti
- Pomozite nam pronaći pravu vrstu podataka!
- Potpora hrvatskome jeziku potpora je ostalim europskim jezicima i obrnuto!

Zahvaljujem na pozornosti.