

Как работи машинният превод?



Проф. д-р Светла Коева
Институт за български език
„Проф. Любомир Андрейчин“
Българска академия на науките

- През 2020 година:
- 44 зетабайта информация в интернет на различни езици;
- по 6 компютърни устройства на жител на планетата.





- Базира се на граматични и семантични правила за анализ и синтез на словосъчетания и изречения.
- Необходими са изчерпателни речници и правила, включващи подробна морфологична, синтактична и семантична информация за двойката езици, за които се осъществява преводът.
- Подобни ресурси обикновено се изготвят от квалифицирани специалисти и изработването им е сложна и продължителна задача.



- Статистическият машинен превод се базира на вероятността дадена последователност от символи да е превод на друга последователност от символи.
- При статистическия машинен превод се използват статистически модели, които се създават на базата на тренировъчни корпуси.
- Тренировъчните корпуси са паралелни преводни текстове и текстове от езика, на който се превежда.

- Езиковият знак е арбитрарен (случаен).

<i>zaldi</i>	<i>igel</i>	<i>txoria</i>	<i>behi</i>	<i>sagua</i>

<http://grammar.about.com/od/ab/g/Arbitrariness.htm>

- Езиковият знак е арбитрарен (случаен).

zaldi

igel

txoria

behi

sagua

horse

frog

bird

cow

mouse

- Граматичността не винаги означава приемливост.

Colorless green ideas sleep furiously.

**Furiously sleep ideas green colorless.*

Chomsky. N. Syntactic Structures. The Hague/Paris: Mouton. 1957.

Бе сгладне и честлинните комбурси
тарляха се и сврецоваха във плите;
съвсем окласни бяха тук щурпите
и отма равапсатваха прасурси.



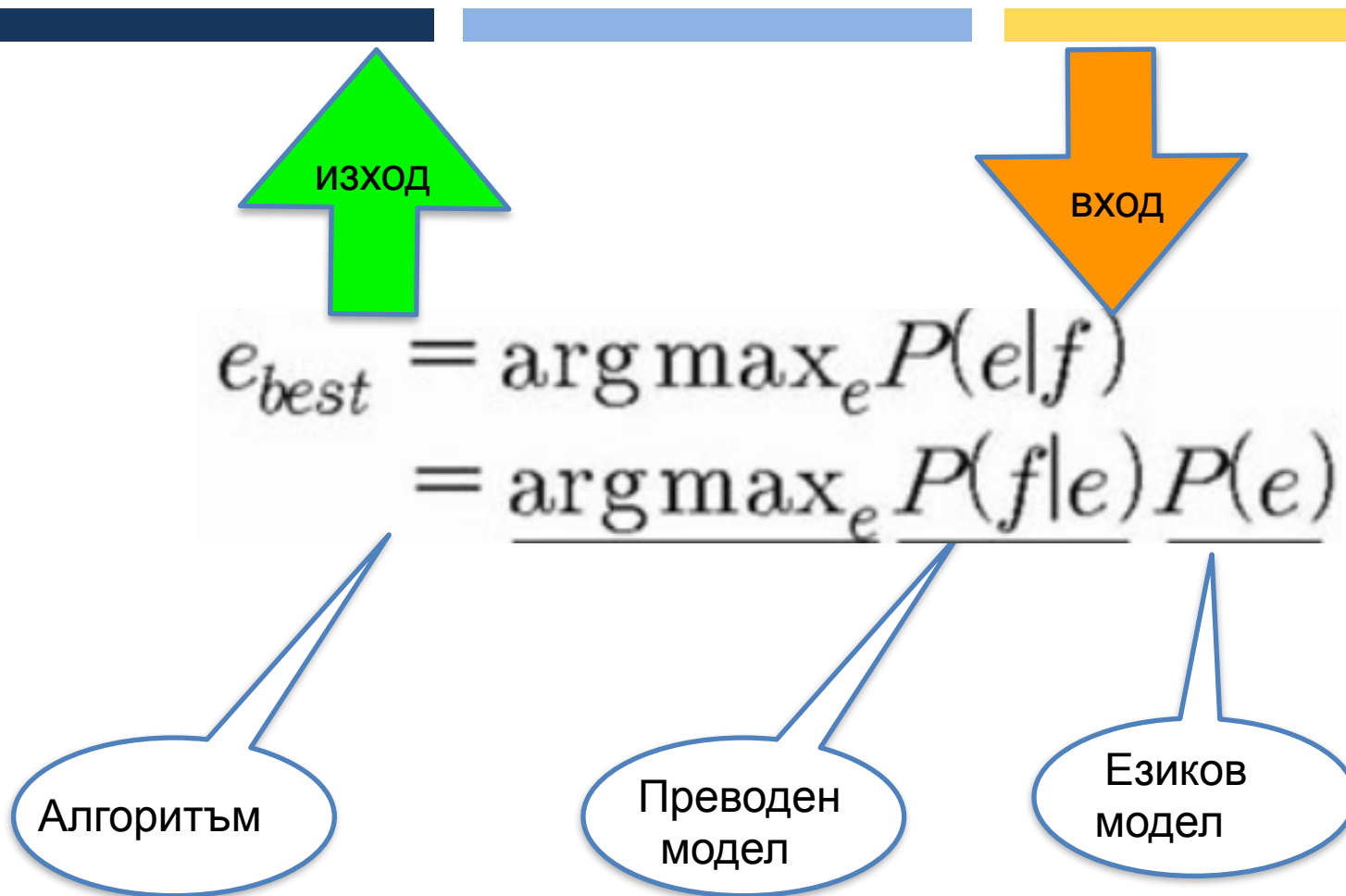
Луис Карол. Алиса в Страната на чудесата и в Огледалния свят. софия, 1996.

- Многозначност

Старите приятели и познати на Ана

си спомниха нейната младост в Пловдив.

- Междуетзикова асиметрия
 - The German chancellor Angela Merkel will make an announcement on Thursday.
 - Angela Merkel wird am Donnerstag eine Ankündigung machen.
 - Германският канцлер Ангела Меркел ще направи изявление в четвъртък.



- Статистическият МТ се обучава от данни:
 - преводни текстове;
 - текстове от езика, на който се превежда.
- Колкото повече данни, толкова по-добре.
- Колкото повече подходящи данни, толкова по-добре.



- Кое изречение в кое изречение се превежда: **съотнасяне по изречения.**
- Коя дума с коя дума се превежда: **съотнасяне по думи и преводна вероятност.**
- Дали полученният превод е лексикално, граматически и стилистично коректен: **езиков модел.**

	Английски	Български
ID01	INFORMATION FOR NATIONAL PARLIAMENTS	ИНФОРМАЦИЯ ЗА НАЦИОНАЛНИТЕ ПАРЛАМЕНТИ
ID02	Article 1	Член 1
ID03	Commission consultation documents (green and white papers and communications) shall be forwarded directly by the Commission to national Parliaments upon publication.	Консултативните документи на Комисията (зелени книги, бели книги и съобщения) се изпращат директно от Комисията на националните парламенти на държавите членки при публикуването им.
ID04	The Commission shall also forward the annual legislative programme as well as any other instrument of legislative planning or policy to national Parliaments, at the same time as to the European Parliament and the Council.	Комисията изпраща на националните парламенти и годишната законодателна програма, както и всеки друг инструмент за законодателно планиране или за политическа стратегия, като едновременно с това ги изпраща и на Европейския парламент и на Съвета.



I had dropped — бях изпуснал
the lantern — фенера
when **I seized** — когато се хванах
the bars — [за] решетката
but it **still** — но той все още
burned upon the floor — гореше на пода

<http://www.bglibrary.net/en-bg-title.htm>



- The
 - oficial
 - Journal
 - of
 - the
 - Community
- Официалният
вестник
на
общността



酸辣汤

Люто кисела супа

鲜蘑牛肉汤 Супа от телешко и печурки

鲜蘑鸡肉汤 Супа от пиле с печурки

鲜蘑鸡 Пиле с печурки



I love this boy

Аз обичам това момче

I love this dog

Аз обичам това куче

They love this dog

Те обичат това куче

I talk to this dog

Аз говоря на кучето

Съотнесени думи и
изречения



I love this boy

Аз обичам това момче

I love this dog

Аз обичам това куче

They love this dog

Те обичат това куче

I talk to this dog

Аз говоря на кучето

I	Аз	3
love	обичам	2
love	обичат	1
this	това	3
this		1
boy	момче	1
dog	куче	2
dog	кучето	1
They	Те	1
talk	говоря	1
talk	говорят	1

I love this boy

Аз обичам това момче

I love this dog

Аз обичам това куче

They love this dog

Те обичат това куче

I talk to this dog

Аз говоря на кучето

I	Аз	3/3
love	обичам	2/3
love	обичат	1/3
this	това	3/4
this		1/4
boy	момче	1/1
dog	куче	2/3
dog	кучето	1/3
They	Те	1/1
talk	говоря	1/1
talk	говорят	1/1



I love this boy
Аз обичам това момче
I love this dog
Аз обичам това куче
They love this dog
Те обичат това куче
I talk to this dog
Аз говоря на кучето

They love this boy.

They	Те	1/1
love	обичам	2/3
love	обичат	1/3
this	това	3/4
this		1/4
boy	момче	1/1

Езиков модел



I love this boy
Аз обичам това момче
I love this dog
Аз обичам това куче
They love this dog
Те обичат това куче
I talk to this dog
Аз говоря на кучето

They love this boy.

They	Те	1/1
love	обичат	1/3
this	това	3/4
boy	момче	1/1

ЕЗИКОВ МОДЕЛ

**Какво е правилно в
езика, на който се
превежда.**



- По-добър модел:
 - Не се превежда дума по дума;
 - Съотнасянето по фрази позволява по-коректно да се представят редица езикови черти, например съгласуване.

- this boy : това момче
- to this dog : на кучето
- I talk : Аз говоря



I love this boy
Аз обичам това момче
I love this dog
Аз обичам това куче
They love this dog
Те обичат това куче
I talk to this dog
Аз говоря на кучето

I love	Аз обичам	2
They love	Те обичат	1
this boy	това момче	1
this dog	това куче	2
I talk	Аз говоря	1
to	на	1
this dog	кучето	1

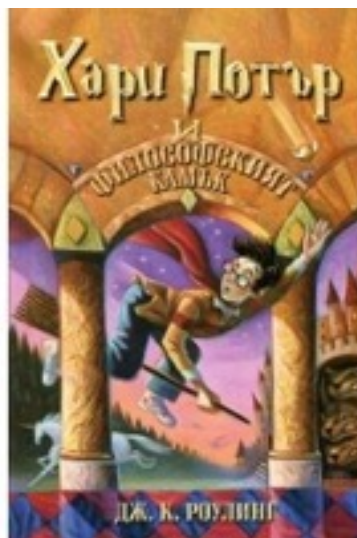
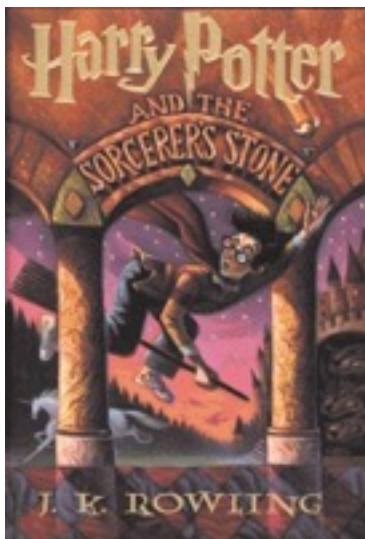


I love this boy
Аз обичам това момче
I love this dog
Аз обичам това куче
They love this dog
Те обичат това куче
I talk to this dog
Аз говоря на кучето

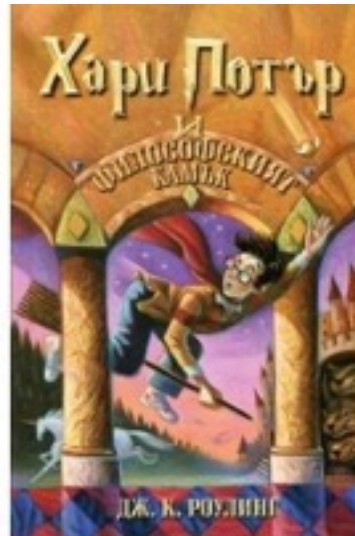
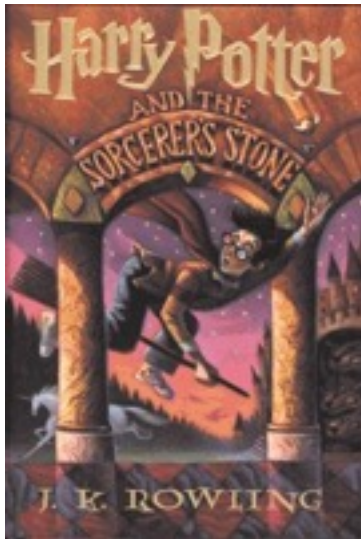
I love	Аз обичам	2
They love	Те обичат	1
this boy	това момче	1
this dog	това куче	2
I talk	Аз говоря	1
to	на	1
this dog	кучето	1
They love this dog.		
Те обичат това куче.		
Те обичат кучето.		



- Статистическият машинен превод зависи от данните.
- Статистическият машинен превод се учи как да превежда от данните.
- Данните са:
 - преводни паралелни текстове;
 - едоезикови текстове (за езика, на който се превежда) ;
 - речници, терминологични бази, онтологии, имена на лица, организации, географски имена.
- Също като при хората статистическият машинен превод е толкова по-добър, колкото повече се е учил.



MOSES  CORE



Protected - Personal Information CIVMEANS7
Legal Aid Agency
Financial Assessment for Family Mediation
Provider reference case code: MED121688273 LAB1
This form must be completed in ink.

Applicant's Details
Surname Mr/Ms/Ms/Ms _____ First name(s) _____
Surname at birth if different _____ Date of birth _____
Address _____ Postcode _____
National insurance number _____
Job _____

Financial Eligibility

1. The client has a partner whose means are to be aggregated.
 Yes Please provide details of both client's partner's means.
 No Please provide details of both client's means only.
2. 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.
3. 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 designated only applies to assets that are specifically 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 5 April 2013 © Crown Copyright

MOSES  CORE

- Механизмът за свързване на Европа осигурява услуги за гражданите на Европейския съюз.
- *Платформата за автоматичен превод CEF.AT* се нуждае от подходящи данни, за да работи по-добре.
- Свободният софтуер Мозес е научна разработка, финансирана от Европейската комисия, и се използва от MT@EC – част от *CEF.AT*.

MOSES  CORE

Координация на езиковите ресурси в Европа

*Да се подкрепят нашите езици
означава да се подкрепя обединена
Европа и обратно.*