# Automated Translation: How Does It Work?

**Stelios Piperidis**
**ELRC, ILSP/Athena RC**

**Simon Krek**
**Jožef Stefan Institute**

Acknowledgements: with adaptations of materials that appeared at MT Marathons, WMTs, etc.

# Machine Translation

Agenda:

- Why MT: Volume, Quality and Cost?
- Why is MT hard?
- MT + Human Translators = Quality
- How does modern statistical MT work?
- Its all about Data!
- And the right kind of Data!

# Machine Translation, Quality and Cost?

- Europe = Multilinguality
- 24 official languages, 24+2 CEF languages
- So much to translate!
- Translation costs!?
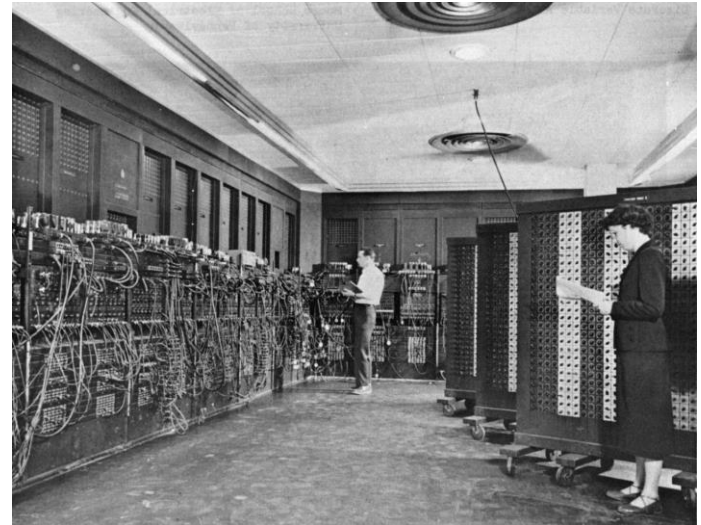- Can MT help?
- What about the Quality?



Image: https://en.wikipedia.org/wiki/ENIAC#/media/File:Eniac.jpg
License: public domain
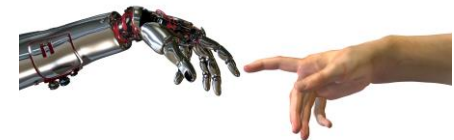
# Why is MT Hard?

- Human languages are:

  – Elegant
  – Efficient
  – Flexible
  – Complex

- One word/sentence may mean many things
- Many ways of saying the same thing
- Meaning depends on context
- Literal and figurative language (metaphor)
- Language and culture (different ways of conceptualising the same thing)
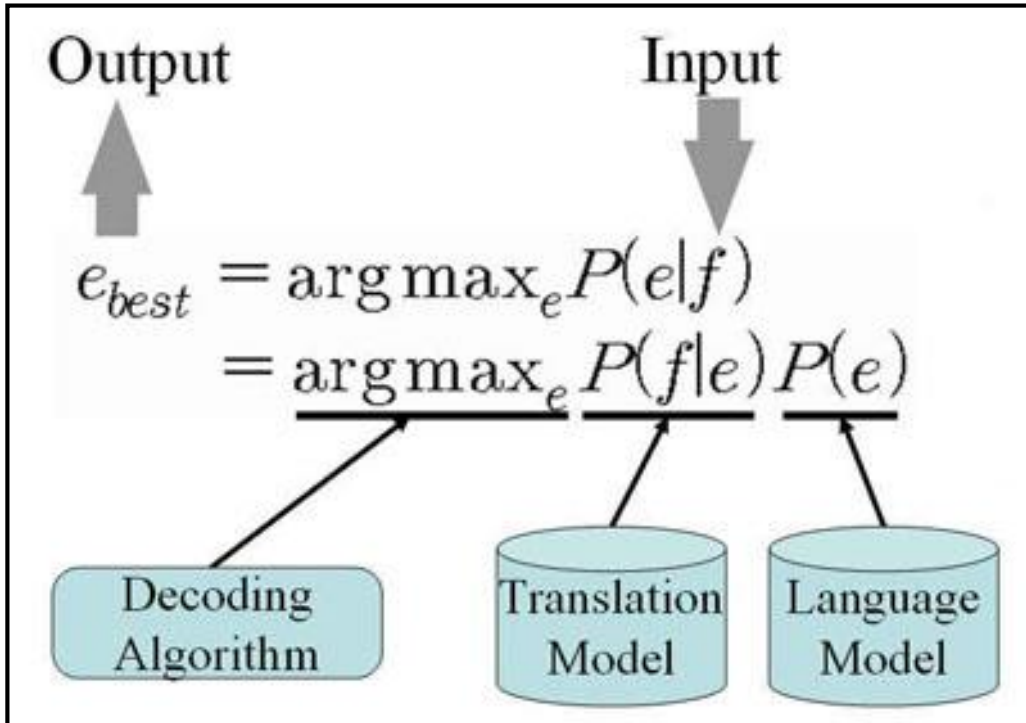
- Word order
- Morphology
- …



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

# Language and Translation is Complex

- Language/translation is complex
- We cannot compute it exactly
- We tried: rule-based MT and LT …
- What do we do?
- Machine Learning
  - Learns from data $\Rightarrow$ data is all important
  - Approximate solution $\Rightarrow$ not perfect, needs help
    - human professional translators
    - Post-editing
    - Automated Translation $\neq$ Automatic

# How does Modern MT Work?

$$e_{best} = \arg\max_e P(e|f)$$
$$= \arg\max_e P(f|e)P(e)$$

Output

Input

Decoding Algorithm

Translation Model

Language Model

- No maths today

- Instead:

- The story of Statistical MT in pictures …

- Its all about **Data** …

# How does Modern MT Work?

Statistical MT learns from data

Two kinds of data:

- Human translations
- Text in the target language

- The more data the better!
- Also: the right kind of data!

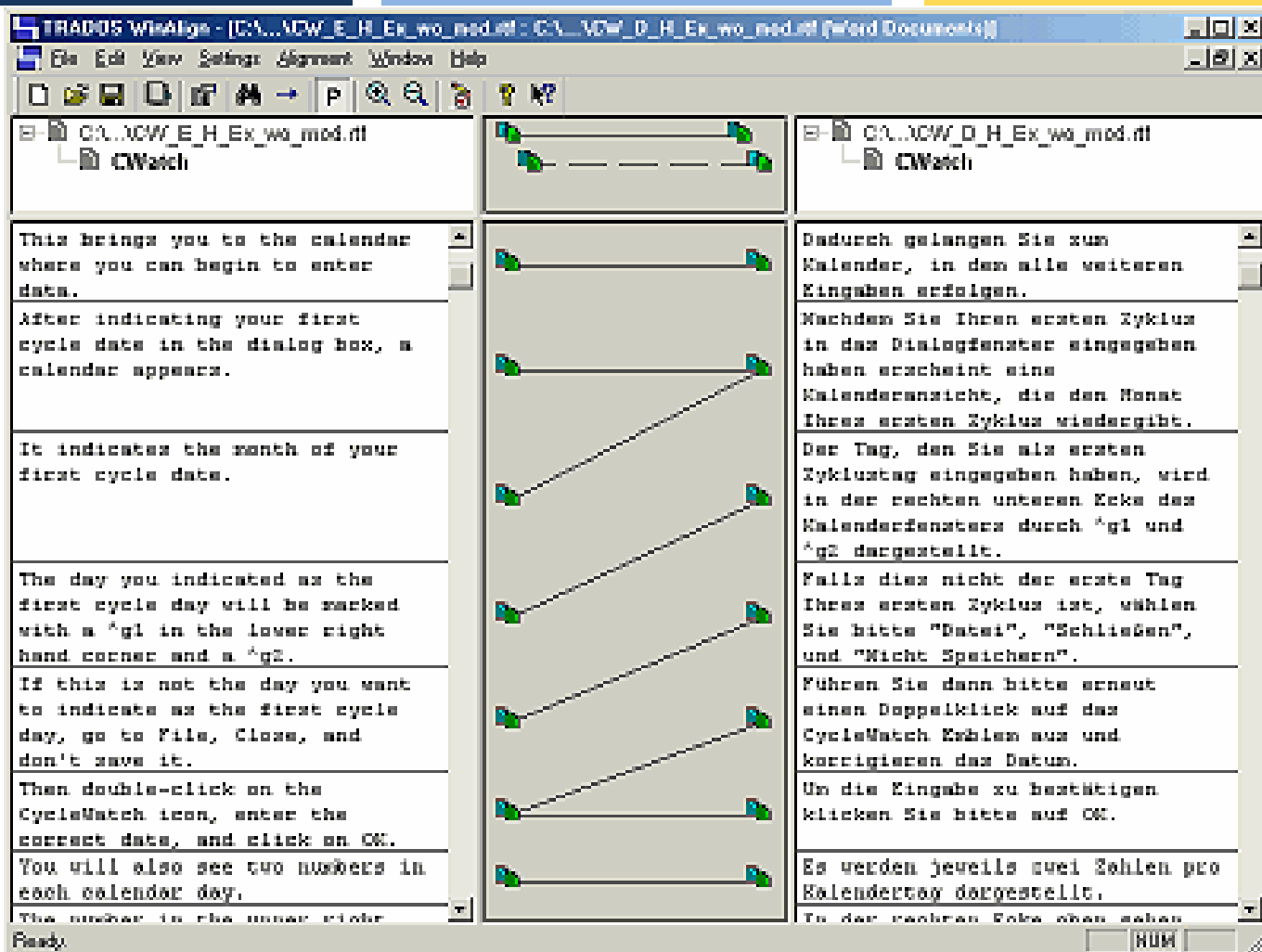| GERMAN | ENGLISH | FRENCH |
| --- | --- | --- |
| Einleitung | Introduction | Introduction |
| *I. Von dem Unterschiede der reinen und empirischen Erkenntnis* | *I. Of the difference between Pure and Empirical Knowledge* | *I. De la différence de la connaissance pure et de la connaissance empirique.* |
| Daß alle unsere Erkenntnis mit der Erfahrung anfange, 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 Verstandestätigkeit 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. | 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. | 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. |

# What can/do we Learn from Data?

- Which sentences translate as which: sentence alignment
- Which words translate as which: word alignment + translation probabilities
- What is good target language like: language model

| GERMAN | ENGLISH | FRENCH |
|---|---|---|
| Einleitung | Introduction | Introduction |
| I. Von dem Unterschiede der reinen und empirischen Erkenntnis | I. Of the difference between Pure and Empirical Knowledge | I. De la différence de la connaissance pure et de la connaissance empirique. |
| Daß alle unsere Erkenntnis mit der Erfahrung anfange, 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 Verstandestätigkeit 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. | 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. | 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. |

# Sentence Alignment

# Word Alignment:

## 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 |
| 雲蛋湯 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 |

# Learning to Translate Words:

- Word alignment mode knows a lot about Chinese soups
- Doesn't know much else …

- Only knows what it has seen in the training data
- Like people …
- A common theme …

- Given word aligned translation data, can we learn a translation dictionary?
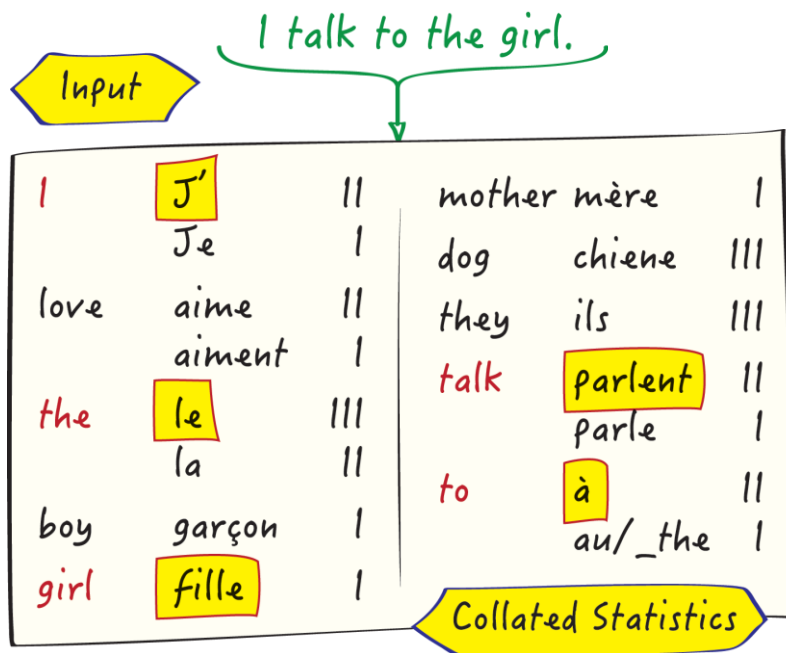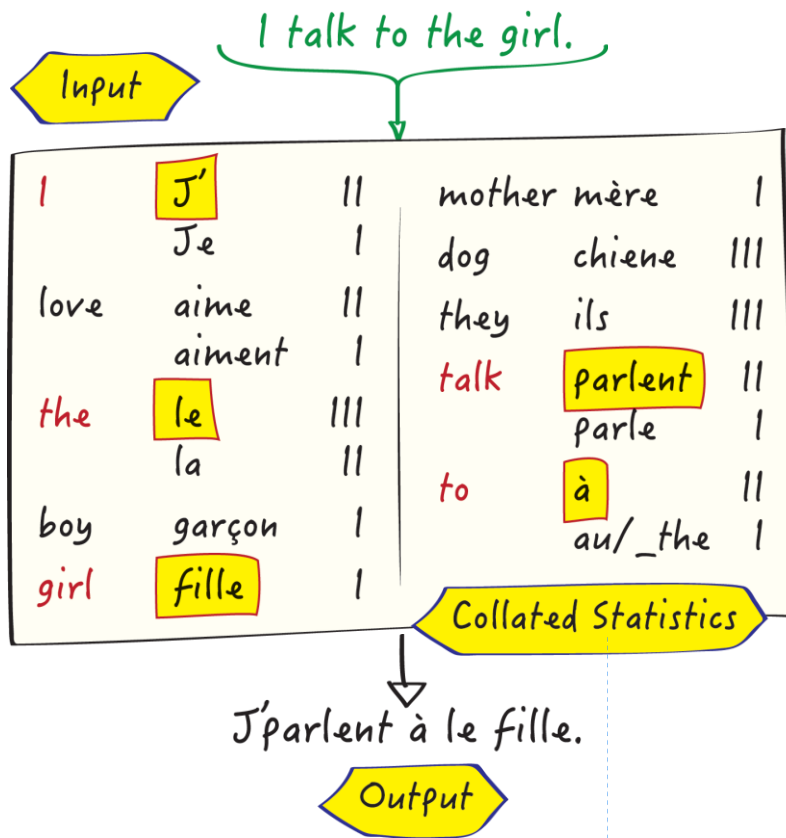- Yes, really easy …

# Statistical Machine Translation



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.
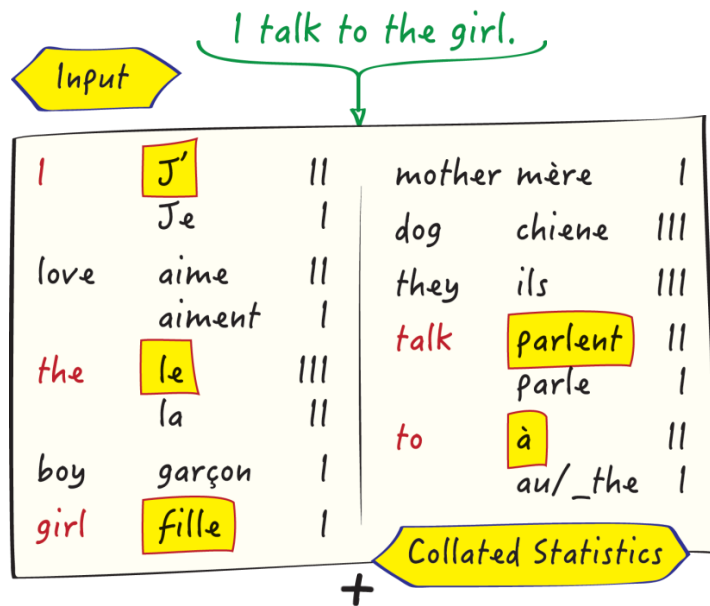Je parle à la mère.

Aligned Data

# Statistical Machine Translation



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.
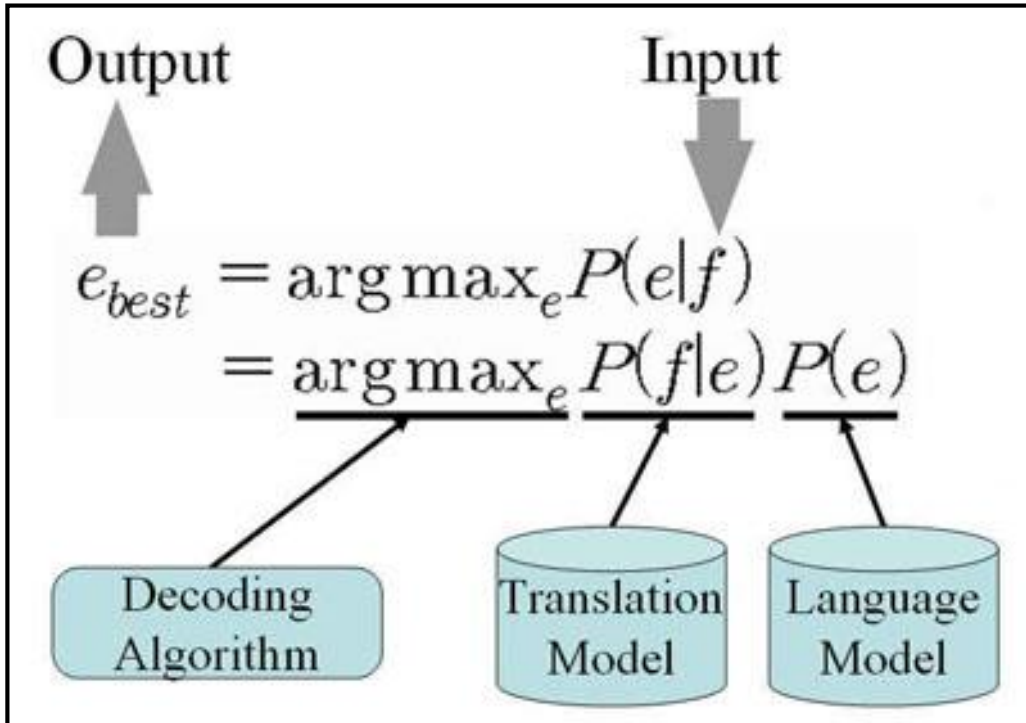Je parle à la mère.

Collated Statistics

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

# Statistical Machine Translation

# Statistical Machine Translation

# Statistical Machine Translation



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 |

## How to choose?
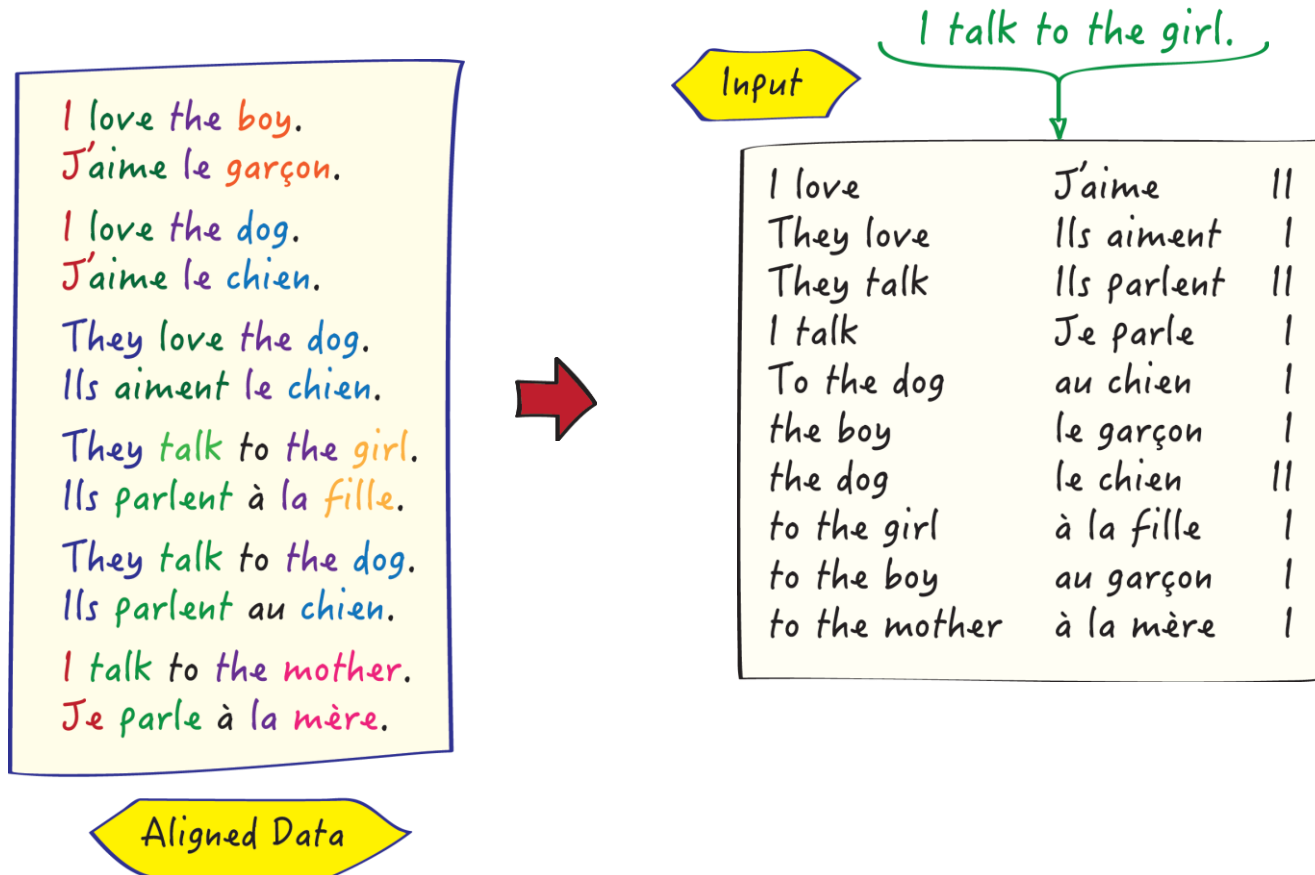
# Statistical Machine Translation

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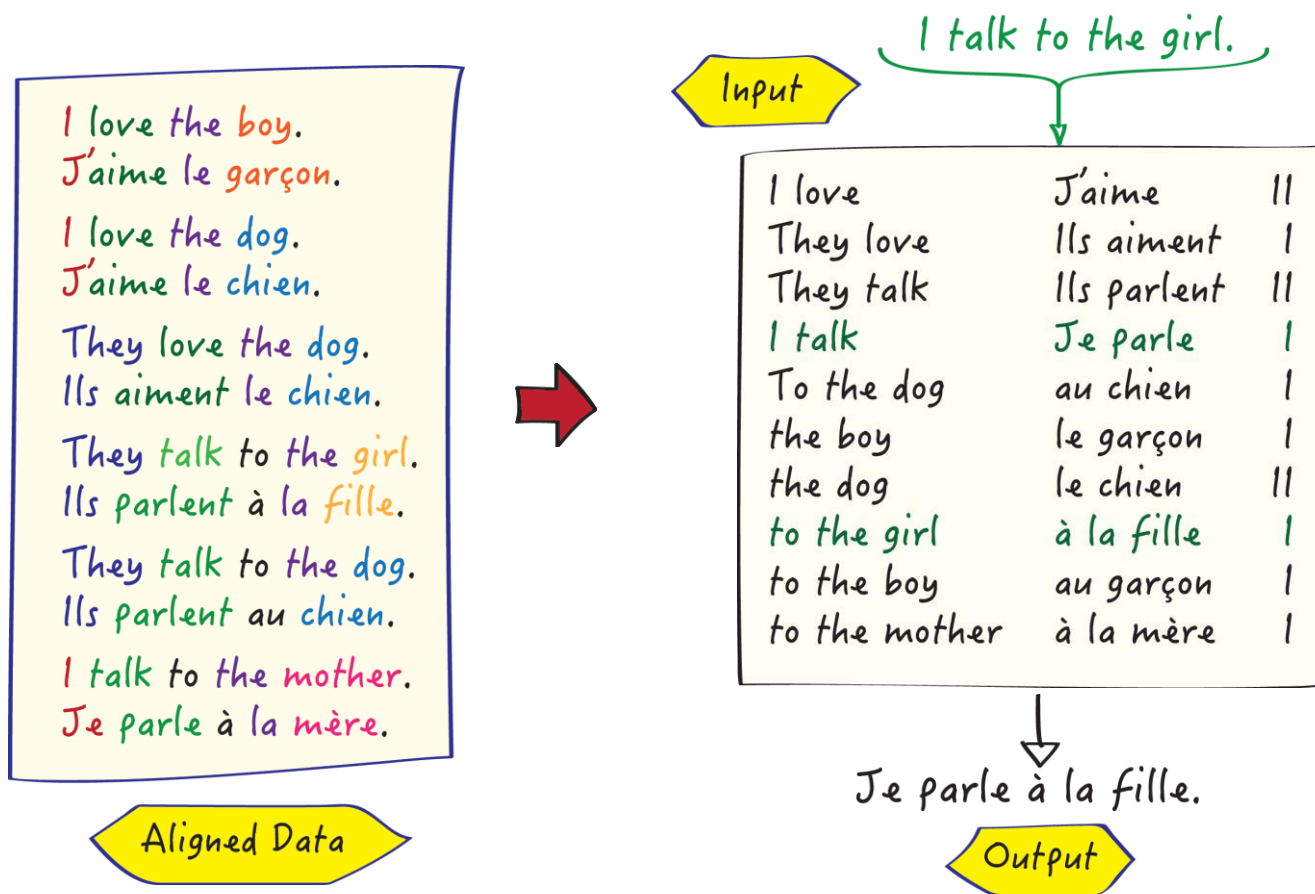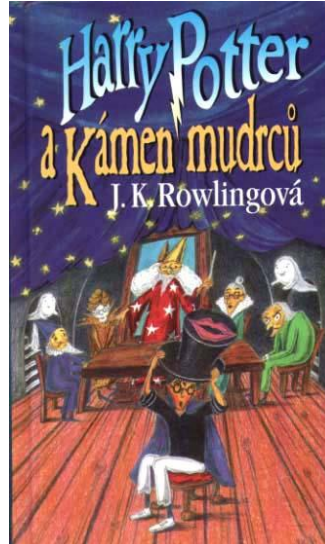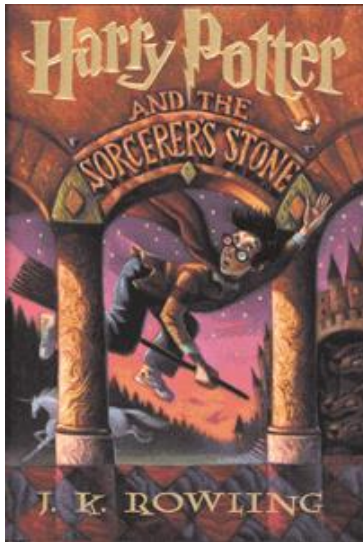I talk to the mother.
Je parle à la mère.

**Aligned Data**

## The Language Model:

- What is good target language?
- Which words can follow which words and which can't … the grammar
- Learnt from the data …

  - *Je parle* is good …
  - *J' parlent* is bad …

  - *la fille* is good …
  - *le fille* is bad …

- *Je parle à la fille >> J' parlent à le fille*

# Statistical Machine Translation

$$e_{best} = \arg\max_e P(e|f)$$
$$= \arg\max_e P(f|e)P(e)$$

Output

Input

Decoding Algorithm

Translation Model

Language Model

- No maths today

- Instead:

- The story of Statistical MT in pictures …

- Its all about **Data** …

# Phrase-Based SMT

- So far: translating single words
- Loses context: such as agreement (*le fille* …) etc.
- To some extent "repaired" by language model

- A better model:
- Not just translations of single words
- But also phrase translations:

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

# Statistical Machine Translation

# Phrase Based - Statistical Machine Translation

# Phrase Based - Statistical Machine Translation

- Much better than word-based SMT!
- Standard technology: Google, Microsoft, Baidu, Global Localisation & Translation Industry

- Moses Open Source PB-SMT
- Most widely used SMT system
- Research funded by EC
- Used by EC DGT's MT@EC

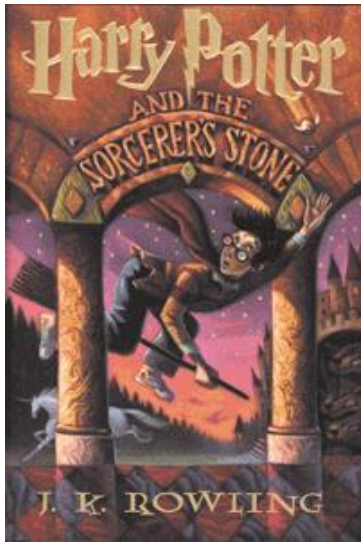# Machine Translation and Data

- Statistical Machine Translation is all about data
- SMT learns how to translate from data
- Data
  - translations (bilingual data)
  - Monolingual data (target language text)
  - Dictionaries, terminology, ontologies, named entities
- Like people SMT is good at what it has learned

# CEF.AT and Data

- CEF.AT needs the right kind of data
- National governments, public administration, public services, NGOs
- CEF provide services for multilingual engagement with national citizens, EU citizens and other customers of public administration

# ELRC

- Help us make CEF.AT a success
  - Services for Europe's citizens
  - Services for you
  - Support multi-linguality

- Help us find the right kind of data

- Supporting our language is supporting Europe and vice versa