



UPPSALA  
UNIVERSITET

# Machine Translation 5LN426 and 5LN711

*Sara Stymne*  
*Uppsala University*

*Slides mainly from Jörg Tiedemann*





# Outline for Today

Motivation

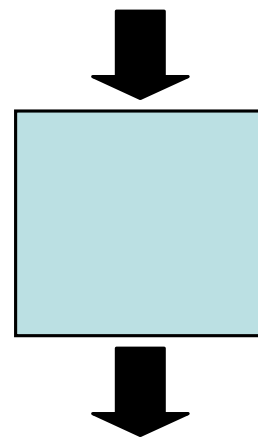
Overview of the course

Classical MT approaches



# Machine Translation

美国关岛国际机场及其办公室均接获一名自称沙地阿拉伯富商拉登等发出的电子邮件，威胁将会向机场等公众地方发动生化袭击後，关岛经保持高度戒备。



The U.S. island of Guam is maintaining a high state of alert after the Guam airport and its offices both received an e-mail from someone calling himself the Saudi Arabian Osama bin Laden and threatening a biological/chemical attack against public places such as the airport.



# Why Machine Translation?

<b>MANDARIN</b>	<b>885,000,000</b>
<b>SPANISH</b>	<b>332,000,000</b>
<b>ENGLISH</b>	<b>322,000,000</b>
<b>BENGALI</b>	<b>189,000,000</b>

<b>TURKISH</b>	<b>59,000,000</b>
<b>URDU</b>	<b>58,000,000</b>
<b>MIN NAN (China)</b>	<b>49,000,000</b>
<b>JINYU (China)</b>	<b>45,000,000</b>

<b>HINDI</b>	<b>182,000,000</b>
<b>PORTUGUESE</b>	<b>170,000,000</b>
<b>RUSSIAN</b>	<b>170,000,000</b>
<b>JAPANESE</b>	<b>125,000,000</b>
<b>GERMAN</b>	<b>98,000,000</b>



<b>GUJARATI</b>	<b>44,000,000</b>
<b>POLISH</b>	<b>44,000,000</b>
<b>ARABIC</b>	<b>42,500,000</b>
<b>UKRAINIAN</b>	<b>41,000,000</b>

<b>WU (China)</b>	<b>77,175,000</b>
<b>JAVANESE</b>	<b>75,500,800</b>
<b>KOREAN</b>	<b>75,000,000</b>
<b>FRENCH</b>	<b>72,000,000</b>
<b>VIETNAMESE</b>	<b>67,662,000</b>

<b>ITALIAN</b>	<b>37,000,000</b>
<b>XIANG (China)</b>	<b>36,015,000</b>
<b>MALAYALAM</b>	<b>34,022,000</b>
<b>HAKKA (China)</b>	<b>34,000,000</b>

<b>TELUGU</b>	<b>66,350,000</b>
<b>YUE (China)</b>	<b>66,000,000</b>
<b>MARATHI</b>	<b>64,783,000</b>
<b>TAMIL</b>	<b>63,075,000</b>

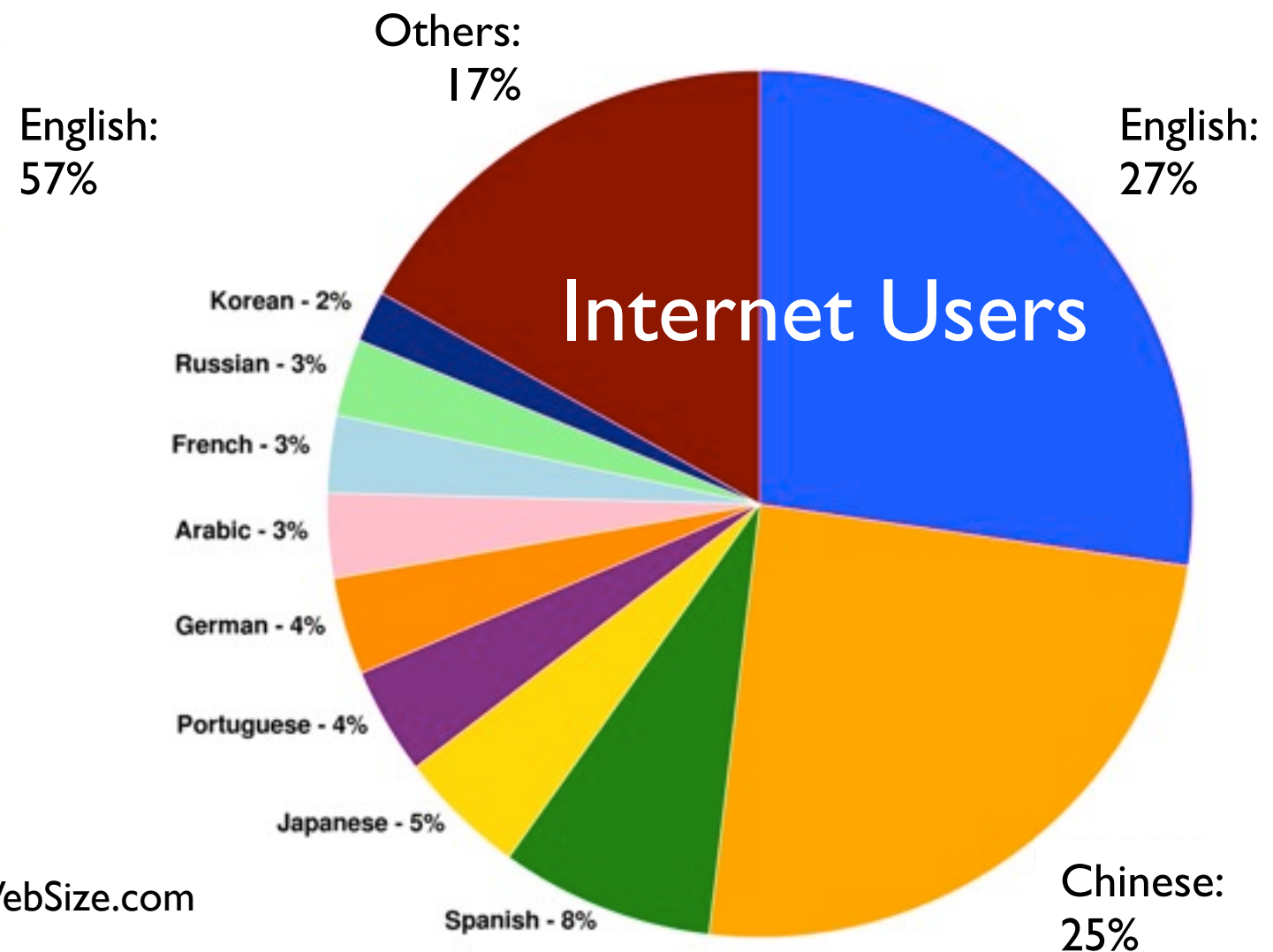
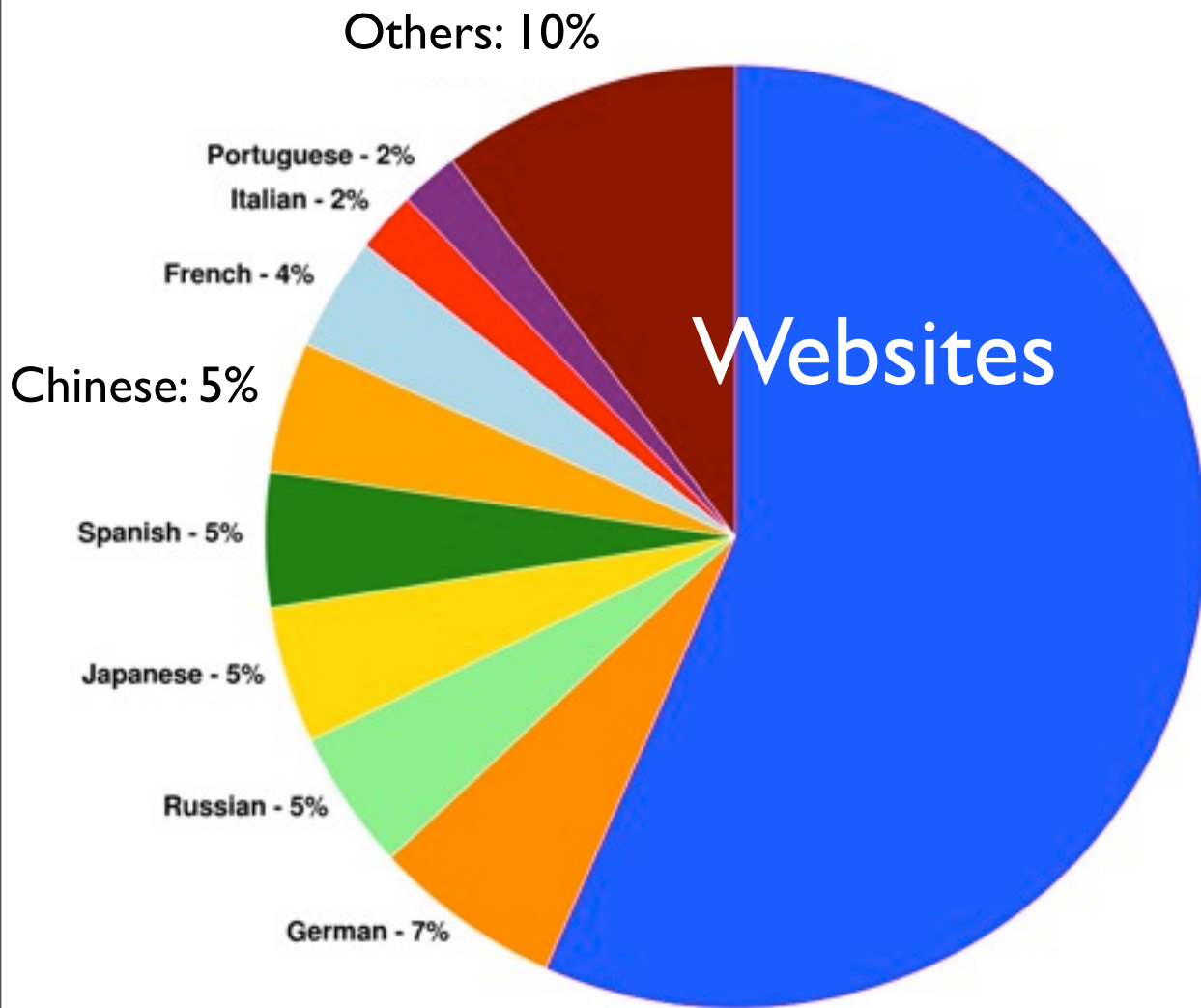
<b>KANNADA</b>	<b>33,663,000</b>
<b>ORIYA</b>	<b>31,000,000</b>
<b>PANJABI</b>	<b>30,000,000</b>
<b>SUNDA</b>	<b>27,000,000</b>

Source: Ethnologue



# Why Machine Translation?

- > 2 billion Internet users
- > 550 million registered domains
- > 12 billion indexed web pages



Sources: W3Techs.com, Internet World Stats, WorldWideWebSize.com



# Why Machine Translation?

- Translation is expensive
- On-line demand for translation (on-the-fly)
- Globalization, growing export
- Lots of language pairs
- Political issues (UN, EU, minority languages, ...)
- Tourism, movies, news
- ...



# MT is a Tough Challenge (and Fun)

Translation errors may be quite severe:

- Doctor's office: *Specialist in women and other diseases*
- Pub: *Ladies are requested not to have children in the bar*
- Hotel: *Please leave your values at the front desk*
- Chinese dining hall: *Translation server error*



MT is not a solved problem ... but constantly improves?

- Input: *Vem vann Allsvenskan i fjol?*
- Google 2010: *Who stole headlines last year?*
- Google 2013: *Who won the Championship last year?*
- Google 2016: *Who won the Olympics last year?*





# MT and Other Language Technology

Translate

From: English



To: Swedish

Translate



Finnish English Swedish Detect language

attribute



language identification



Swedish German English

attribut

speech synthesis



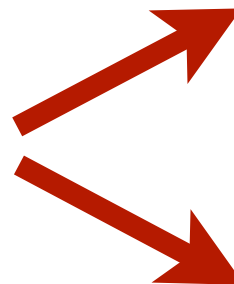
noun

- egenskap property, attribute, capacity, trait, qualification, quality
- attribut attribute
- bestämning determination, qualifier, modifier, adjunct, attribute
- kännetecken feature, characteristic, attribute, badge, token, earmark

verb

- tillskriva ascribe, attribute, accredit, assign to, set down, impute
- attribuera attribute

part of speech



synonyms







# MT is a Cool Research Topic

## How does human language work?

- What are the differences between languages?
- How can we preserve meaning when translating?

## Complex but natural task

- MT is not a solved problem
- MT is a useful end-user application

## Combines various aspects of computational linguistics

- analyze text or speech
- understand/transfer meaning
- generate text or speech



# What is the Problem with MT?

## Unrealistic expectations

- “MT is a waste of time because you will never make a machine that can translate Shakespeare”
- MT is useless because it may translate “*The spirit is willing but the flesh is weak*” into the Russian equivalent of “*The vodka is good, but the steak is lousy*”

## Unexpected (not humanlike) errors

- German Input: *Fussball ist langweilig. **Tore** gibt es selten.*
- Google 2012: *Fotboll är tråkigt. **Gates** är sällsynta.*
- Google 2016: *Fotboll är tråkigt . **Mål** är sällsynta .*



# What are the problems?

- Source language ambiguity
- Cross-lingual divergences
- Target language variation



# Source language ambiguity

## Get (English)

- *I'll **get** a cup of coffee*
- *I didn't **get** the joke*
- *I **get** up at 8am*
- *I **get** nervous*
- *Yeah, I **get** around*

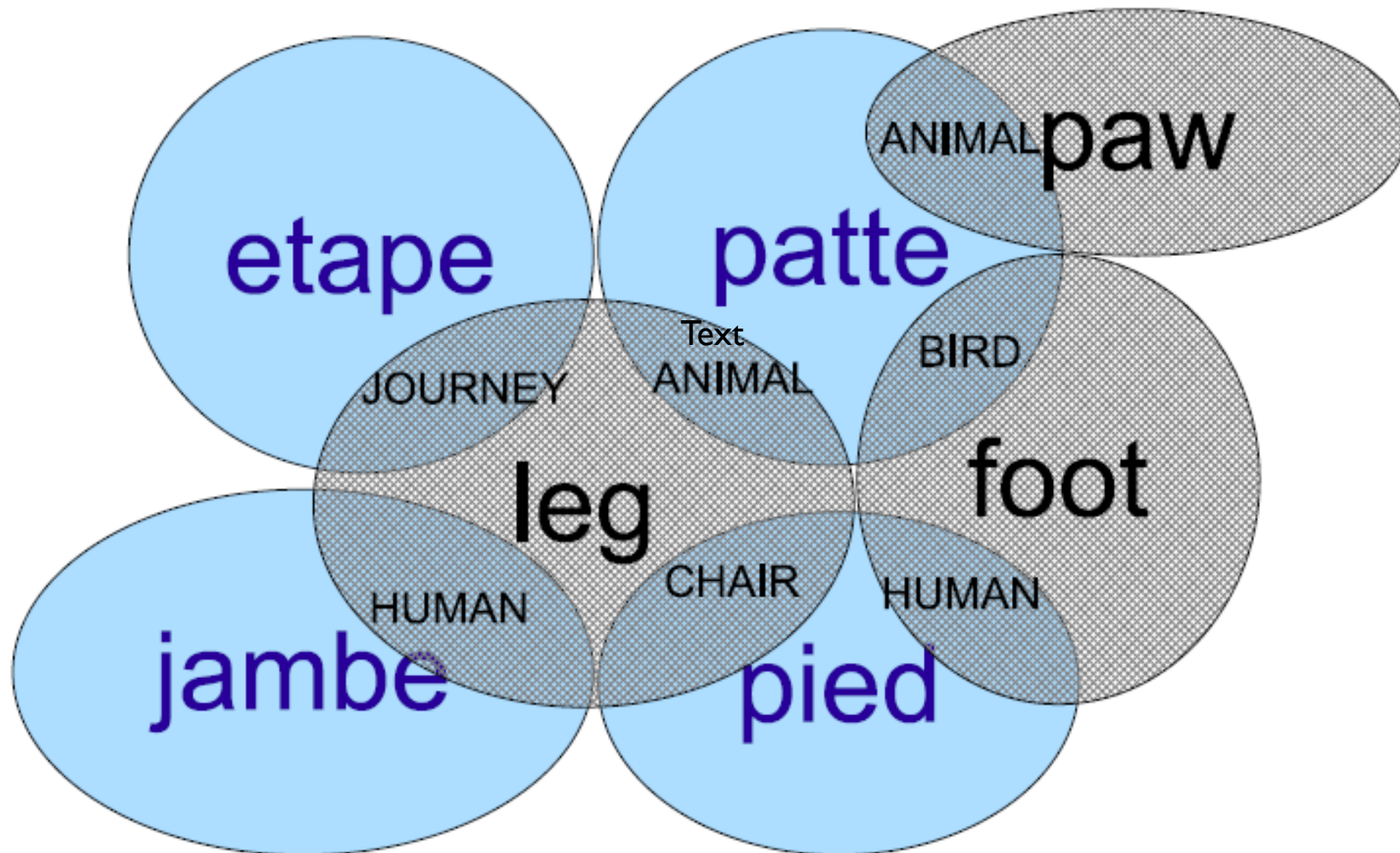
## Var (Swedish)

- was, were (verb)
- each, every (pron)
- where, apiece (adv)
- pus (noun)

> Ambiguity is usually solved in context



# Lexical Ambiguities across Languages



From Jurafsky and Martin



# Language Divergences and Mismatches

## Systematic differences between the 2 languages

- morphology (isolating vs polysynthetic, agglutinative vs fusional)
- syntax (SVO, SOV, VSO, argument structure, pro-drop)

## Idiosyncratic and lexical differences

- differences in lexical ambiguity
- lexical gaps
- differences in tempus, aspect, voice
- different idiomatic/fixed expressions
- ...



# Verb Frame Divergences

## Categorial

- *Kim var förkyld – Kim had a cold*

## Conflation

- *Kim snyter sig – Kim blows her nose*

## Structural

- *Kim sätter sig upp mot Bo – Kim defies Bo*

## Head swapping

- *Kim packar klart – Kim finishes packing*

## Thematic

- *Me gustan uvas – I like grapes*





# Variation in Target Language

## Redundancy of natural languages

- *translate "Vid avslutad kurs ..."*
  - *On completion of the course ...*
  - *After completion of the course ...*
  - *Having completed the course ...*
  - *After finishing the course ...*
  - *Once the course has been completed ...*
  - *...*

Which one is best? How do we decide that?



# In-domain MT with Related Languages

## Example from the book:

### **French input**

Nous savons très bien que les Traités actuels ne suffisent pas et qu'il sera nécessaire à l'avenir de développer une structure plus efficace et différente pour l'Union, une structure plus constitutionnelle qui indique clairement quelles sont les compétences des États membres et quelles sont les compétences de l'Union.

### **Statistical machine translation**

We know very well that the current treaties are not enough and that in the future it will be necessary to develop a different and more effective structure for the union, a constitutional structure which clearly indicates what are the responsibilities of the member states and what are the competences of the union.

### **Human translation**

We know all too well that the present Treaties are inadequate and that the Union will need a better and different structure in future, a more constitutional structure which clearly distinguishes the powers of the Member States and those of the Union.



# MT Between Less Related Languages

Also from the book:

## **Chinese input**

伦敦每日快报指出,两台记载黛安娜王妃一九九七年巴黎死亡车祸调查资料的手提电脑,被从前大都会警察总长的办公室里偷走.

## **Statistical machine translation**

The London Daily Express pointed out that the death of Princess Diana in 1997 Paris car accident investigation information portable computers, the former city police chief in the offices of stolen.

## **Human translation**

London's Daily Express noted that two laptops with inquiry data on the 1997 Paris car accident that caused the death of Princess Diana were stolen from the office of a former metropolitan police commissioner.



# How do Humans do it?

## Human translators need

- to **understand** the source language
- to know how to **speak** the target language (well)
- knowledge about the **topic** of the text to be translated
- knowledge about culture, values, traditions and expectations of speakers in both languages

## Corresponding NLP challenges

- Natural language understanding
- Language generation
- Topic detection and domain adaptation



# Is it Possible at all?

Balance MT quality and input restrictions, depending on task

general purpose browsing quality	post-editing editing quality	sublanguage publishing quality
fully automatic <b>Gisting</b>	computer-aided translation ( <b>CAT/</b> <b>MT</b> )	fully automatic <b>FAHQMT</b>
on-line service	localization, ...	domain-specific tasks



# What exactly is MT?

- MT = automatic translation from one language (source language) to another (target language) using computers
- MT  $\neq$  translation memories and bilingual dictionaries
- MT - usually sentence-by-sentence translation
- MT often refers to translation of written text (cf speech-to-speech translation)
- Semi-automatic: CAT = computer aided translation



# Computer-Assisted Translation Tools

A range of tools to support translators

Translation memories

- A database that stores previously translated sentences/segments
- When translating a new segment, it searches for a matching segment, to display
- Fuzzy matching, it finds similar segments if no full match, and highlights the differences
  - The translator edits this segment, if good enough
  - A score is shown that indicates how similar the matched segment is
- Some TM software has integration with MT





UPPSALA  
UNIVERSITET

# Course overview



# Course Overview (5LN426, 5LN711)

## Lectures

- Introduction of main MT approaches
- MT Evaluation
- Basics of Statistical MT and Word-Based Models
- Phrase-Based SMT
- Tree-based SMT; Document-Level Models
- Seminars: Advanced topics in SMT  
given by master students

## Labs

- Practical sessions and assignments
- 4 written reports, 1 oral
- Performed in pairs, signup by email to Sara and Aaron



# Course Overview (5LN426, 5LN711)

## Bachelor students

- 3 individual assignments
  - 2 assignments on MT/SMT, will be given at least 2 weeks before each deadline
  - 1 literature assignment, choose and summarize a research article

## Master students

- Group project, 3-4 students
- Groups will be created based on wishes for topics to work on
- Give a seminar on your topic
- Perform a practical project and write a report
- Individual reflection report



# Examination

## Bachelor program

- Lab sessions with assignments, no grades
- Three graded individual reports
- VG:VG on at least 2 individual assignments

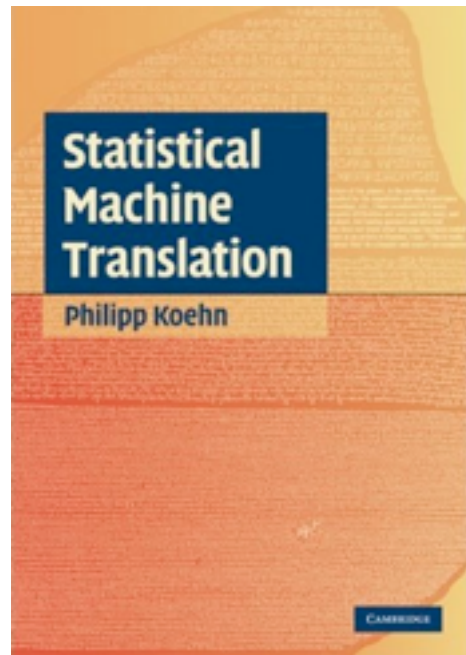
## Master program

- Lab sessions with assignments, no grades
- Group work, graded (in total)
- Individual reflection report, graded
- Final grades based on both



# Course Information

Website: <http://stp.lingfil.uu.se/~sara/kurser/MT16/>



## Literature:

- Philipp Koehn: Statistical Machine Translation
- Daniel Jurafsky and James H. Martin: Speech and language processing
- Other (on-line) material including research articles

## Lab sessions:

- STP Linux account



# Teachers

## Teachers:

- Sara Stymne (Course coordinator, lectures, lab, supervision)
- Aaron Smith (Lecture, labs)
- Christian Hardmeier (Lectures, supervision)
- Fabienne Cap (Lectures, supervision)
- Mats Dahllöf (Examiner)

## Guest lectures:

- Anna Sågvall-Hein, Convertus
  - Machine translation provider
- Nils-Erik Lindström and others, Semantix
  - Translation and interpreting agency



# Road Map

Type	Date	Time	Place	Topic	Reading / Assignments
F	2016-03-30	10-12	6-K1031	Introduction (SS)	Koehn 1; JM 25.1-2; Hutchins; CFMF
F	2016-03-30	14-16	6-K1031	MT evaluation (SS)	Koehn 8; JM 25.9
F	2016-04-04	10-12	2-0076	MT in practice (Convertus) - guest lecture	
L	2016-04-06	10-12	Chomsky	MT in practice (AS)	lab report 1
F	2016-04-11	10-12	2-0076	Introduction to SMT (FC)	Koehn Ch 4, Ch 7, KK97
L	2016-04-13	10-12	Chomsky	Word-based SMT (SS)	lab report 2
L	2016-04-18	10-12	Chomsky	Word-based SMT (SS)	lab report 2
F	2016-04-18	14-16	?	Machine translation at Semantix, a translation provider - guest lecture	
F	2016-04-20	10-12	6-K1031	Parallel Corpora, Alignment (AS)	Koehn 2-4, JT 3-4, KK97, KK99
L	2016-04-25	10-12	Chomsky	Parallel corpora & alignment (AS)	lab report 3
F	2016-04-27	10-12	2-0076	Phrase-based SMT (FC)	Koehn Ch 5
L	2016-05-02	10-12	Chomsky	Phrase-based SMT (AS)	lab report 4
F	2016-05-04	10-12	6-K1031	Decoding (CH)	Koehn Ch 6
L	2016-05-09	10-12	Chomsky	Phrase-based SMT (AS)	lab report 4
F	2016-05-11	10-12	2-0076	Tree-based SMT & MT for morphologically rich languages (SS, FC)	Koehn 10.2, 11
F	2016-05-16	10-12	2-0076	Document-wide decoding & Neural MT (CH)	
L	2016-05-18	10-12	Chomsky	Document-wide decoding lab (AS)	oral lab report 5
S	2016-05-23	10-12	2-0076	Seminar - master student presentations	
S	2016-05-25	10-12	6-K1031	Seminar - master student presentations	





# Deadlines

## Deadlines (All, Master, Bachelor)

- April 13, lab 1 (A)
- April 15, topic selection (M)
- April 26, lab 2 (A)
- May 3, lab 3 (A)
- May 10, ass. 1 (B)
- May 16, ass 2a (B)
- May 17, lab 4 (A)
- May 18/25, lab 5 (A)
- May 23/25, seminar presentation (M)
- June 3, ass. 2b+3 (B)
- June 3, project report (M)
- June 3, reflection report (M)
- June 3, backup lab deadline (A)

Note that there are also many deadlines in the information extraction course! You need to plan your time carefully!



# Changes from last year

## Changes

- New course coordinator and one new teacher: practical reasons
- Additional guest lecture: course improvement
- Labs in pairs: practical reasons, and better interaction
- Bachelors:
  - Assignments instead of exam: based on evaluation and better fit with syllabus
- Masters
  - Group project instead of individual project: practical reasons, better chance for interaction between students, chance to do more advanced projects



UPPSALA  
UNIVERSITET

# Classical Translation Models



# MT as Decoding ("Engineer's MT")

When I look at an article in Russian, I say:  
*This is really written in English,*  
*but it has been coded in some strange symbols.*  
*I will now proceed to decode.*

[Weaver, 1947, 1949]





# Linguistically Motivated MT

Assume that human translators work like this:

- **Understand** the original message
- **Transfer** meaning to the target language context
- **Produce** a grammatical message in the target language

## Transfer-based Machine Translation

- Source language **analysis**
- **Transfer** abstract representation
- **Generation** of target language text

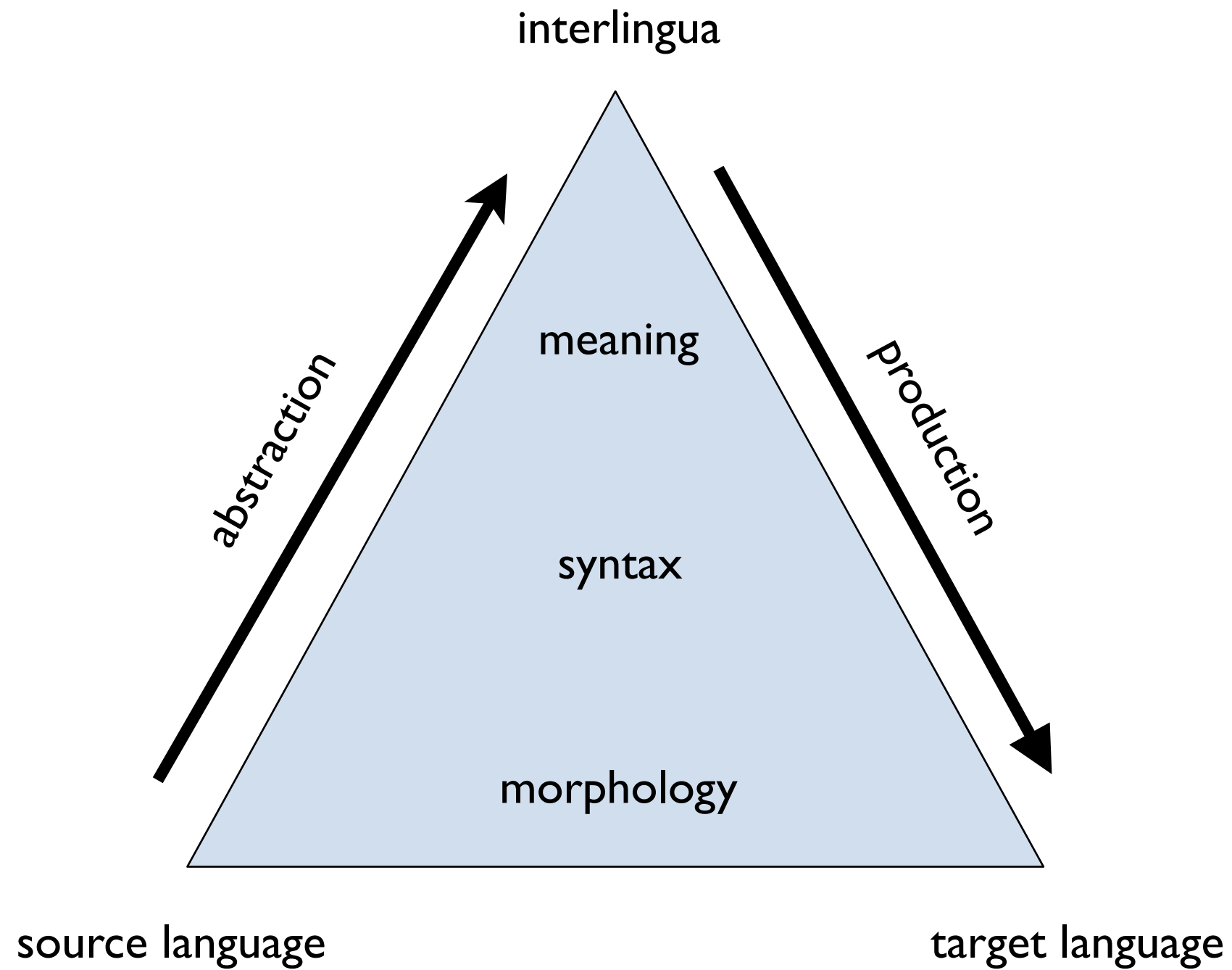


# Early optimism

- Show video!
- Georgetown-IBM demo, 1954
- English-Russian MT
  - 250 words
  - 6 grammar rules
- Much media coverage!



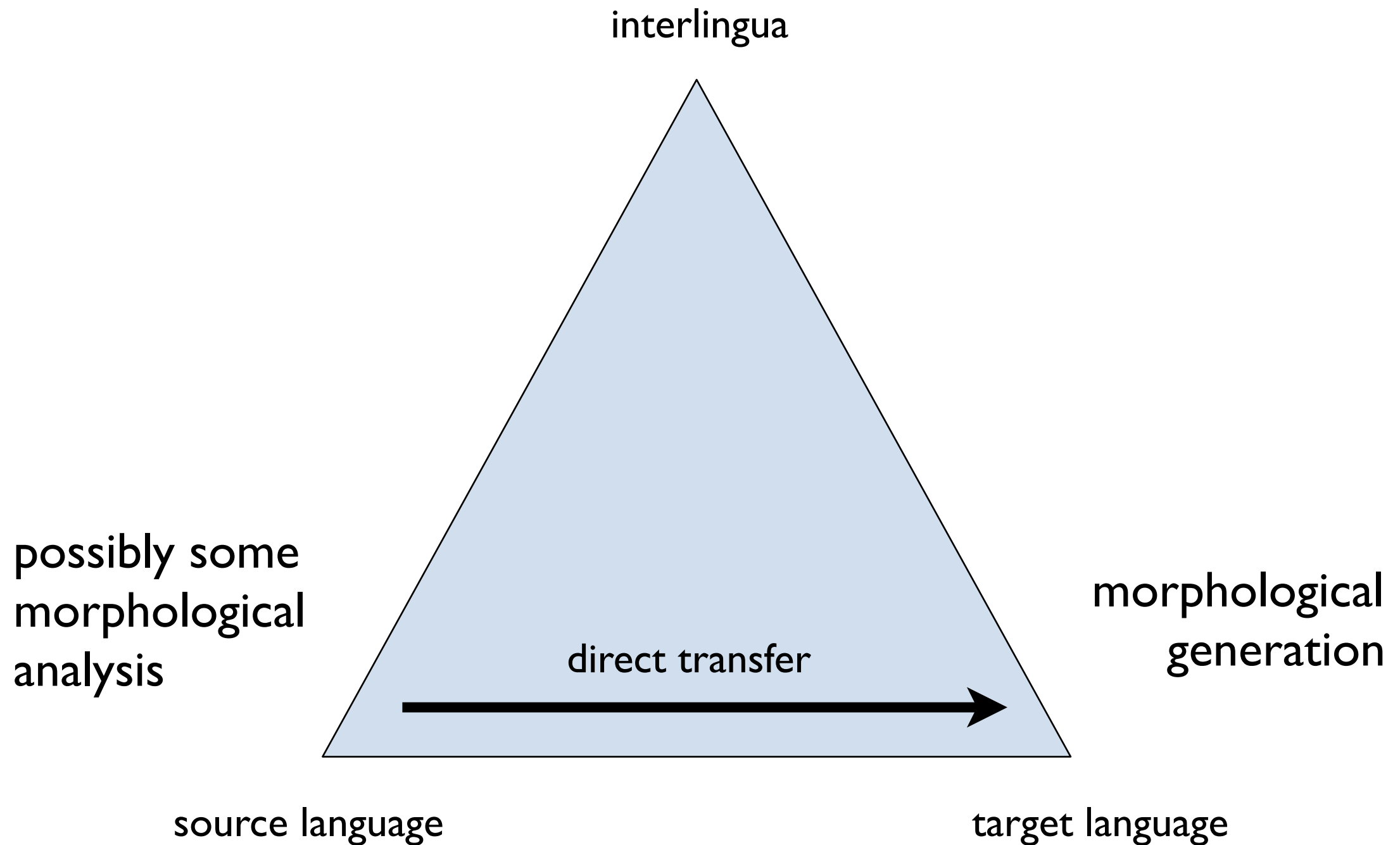
# The Vauquois Triangle





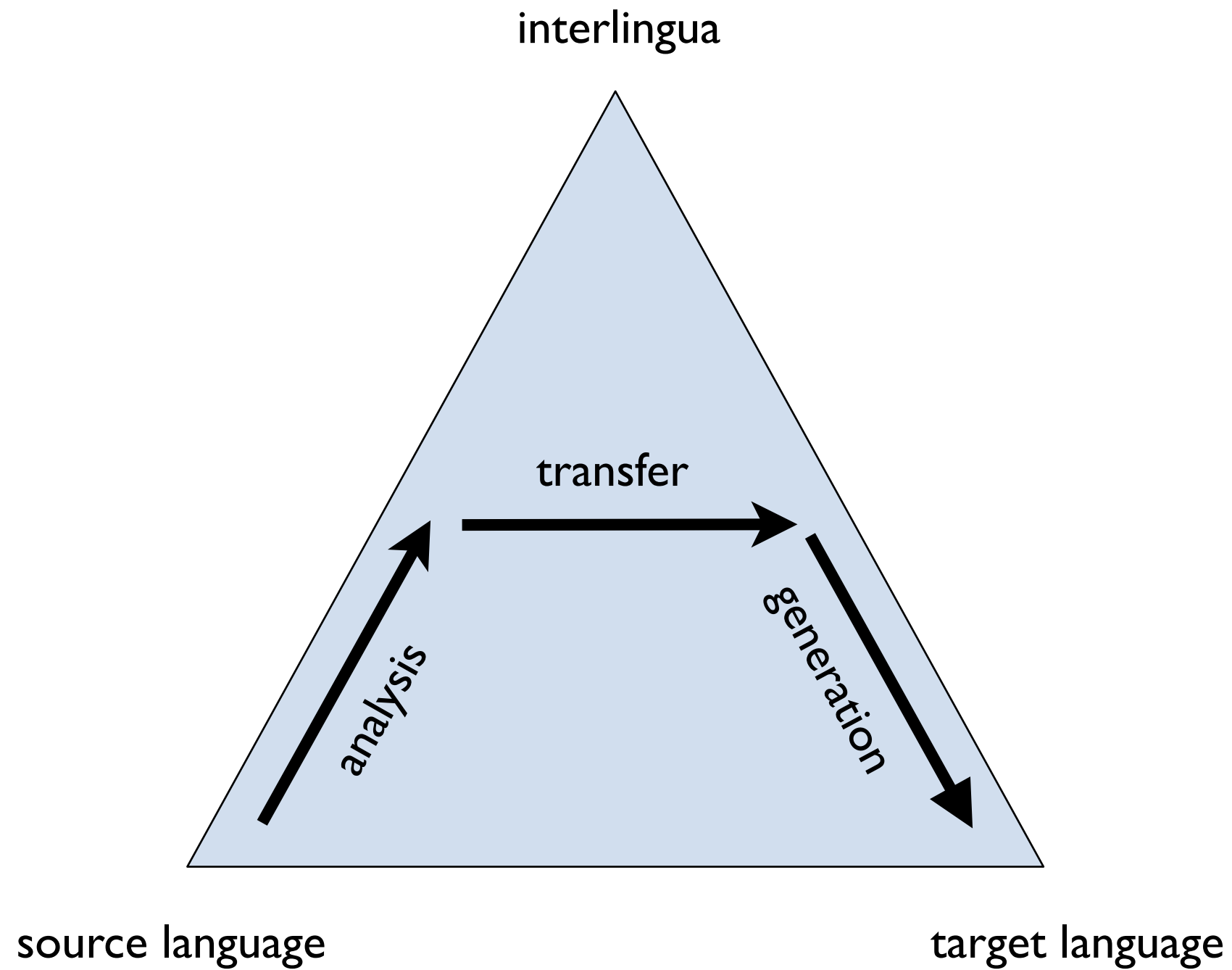


# “Direct Machine Translation”





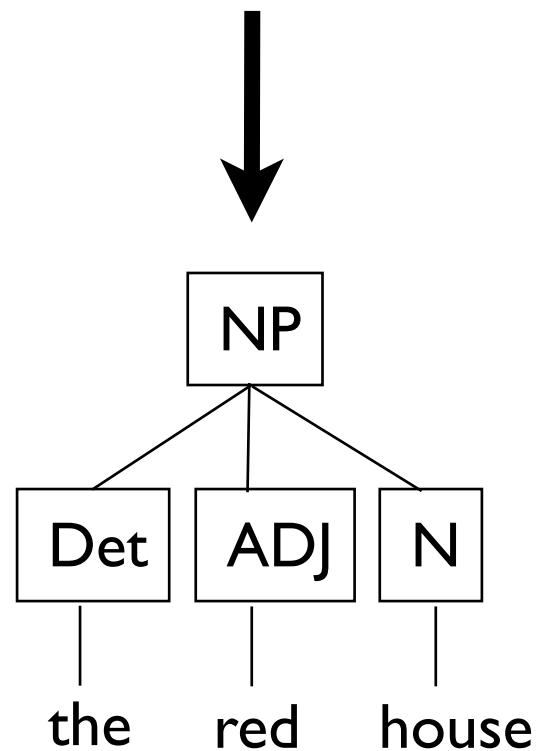
# Transfer-Based Systems: 3 Steps





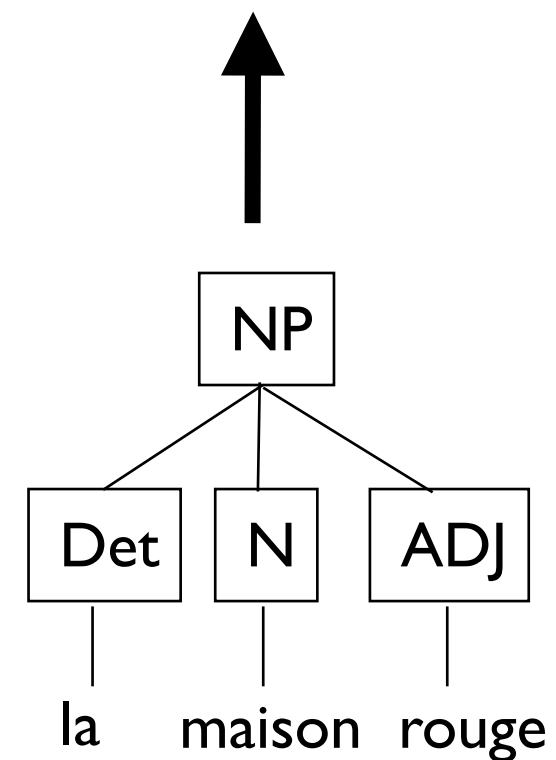
# Transfer-Based Machine Translation

The red house .



N → house  
ADJ → red  
Det → the  
NP → Det ADJ N

La maison rouge.



N → maison  
ADJ → rouge  
Det → la  
NP → Det N ADJ



# Transfer-Based Machine Translation

## What do we need?

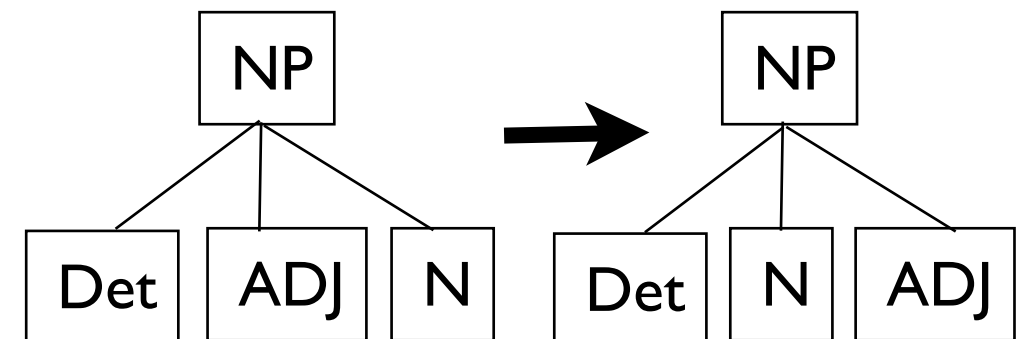
- source language parser
- transfer rules
- target language generator

## Assumption

- Languages are very similar on a high level of abstraction

## Transfer

- large bilingual dictionaries
- structural transfer rules





# Transfer-Based Machine Translation

## Assumptions

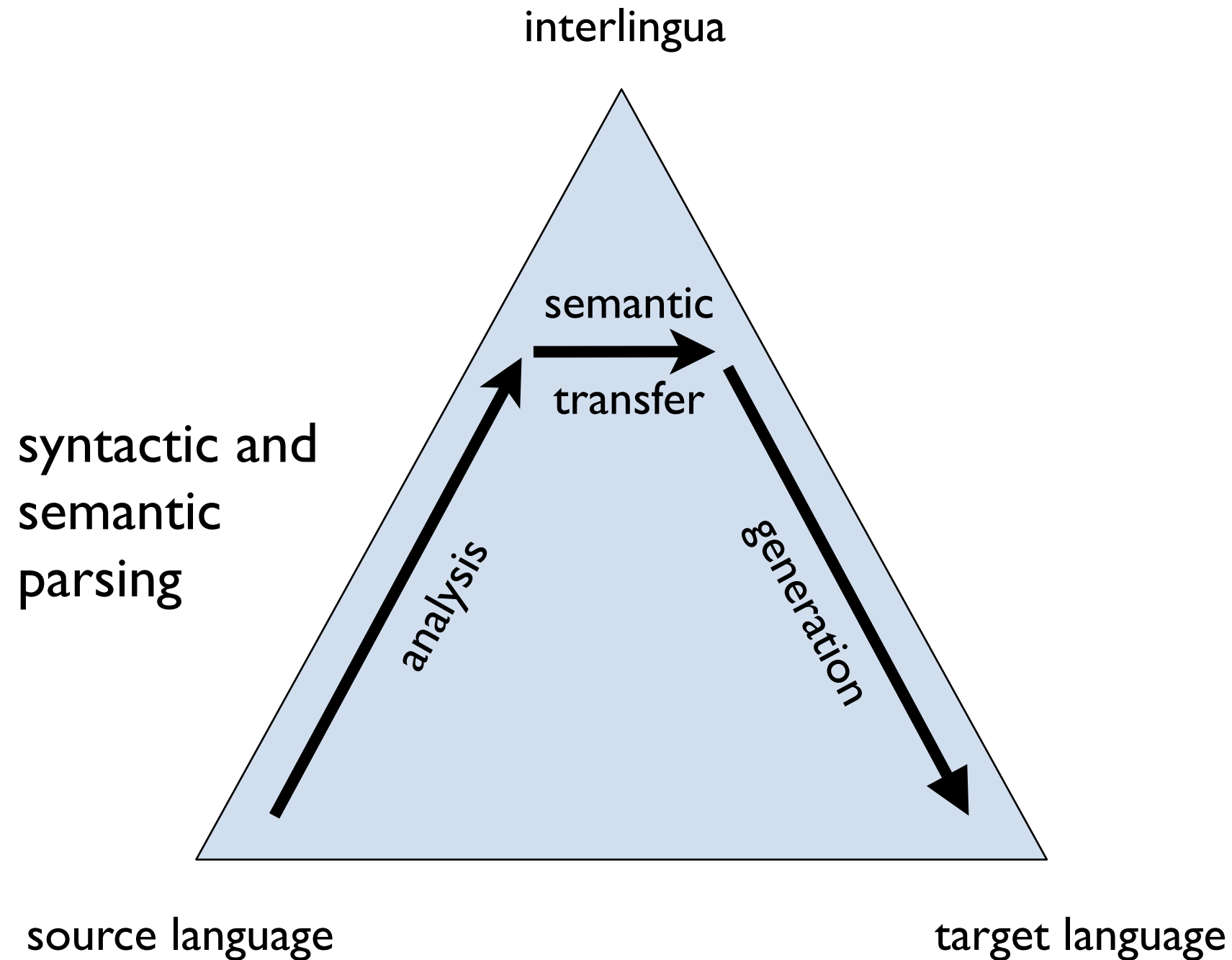
- strong language regularities
- only a few exceptions
- languages = compact grammars + lexicon
- robust and accurate parsing and generation

## Grammar development (3 independent modules)

- rules defined by experts (traditional approach)
- induced from data

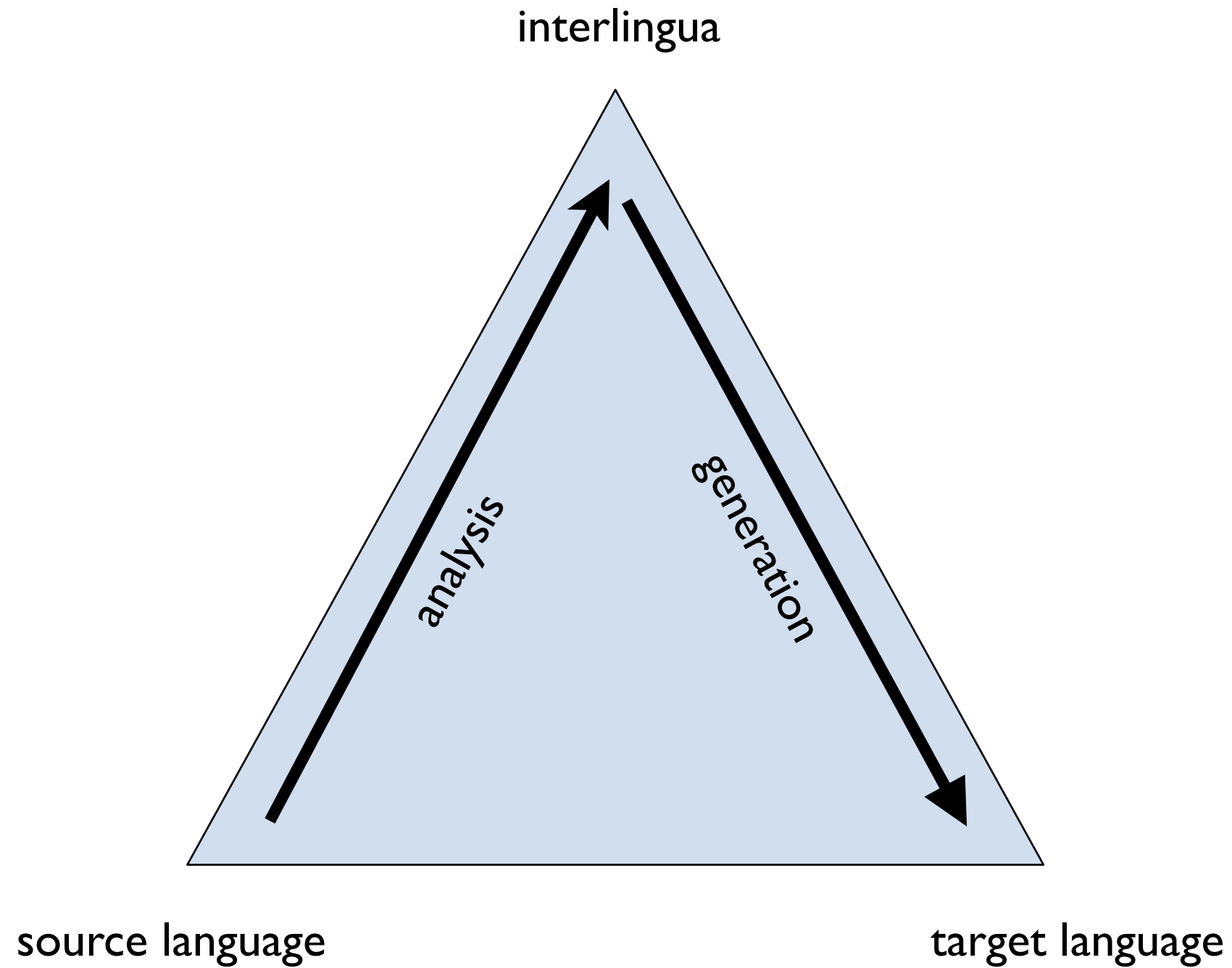


# Transfer-Based Machine Translation





# Interlingua-Based Models





# Interlingua-Based Models

## Assumptions

- All languages can be generated from the same abstract meaning representation
- All aspects of language can be captured by an interlingua

## Advantage:

- no transfer
- new language = new analysis and generation modules, no transfer modules

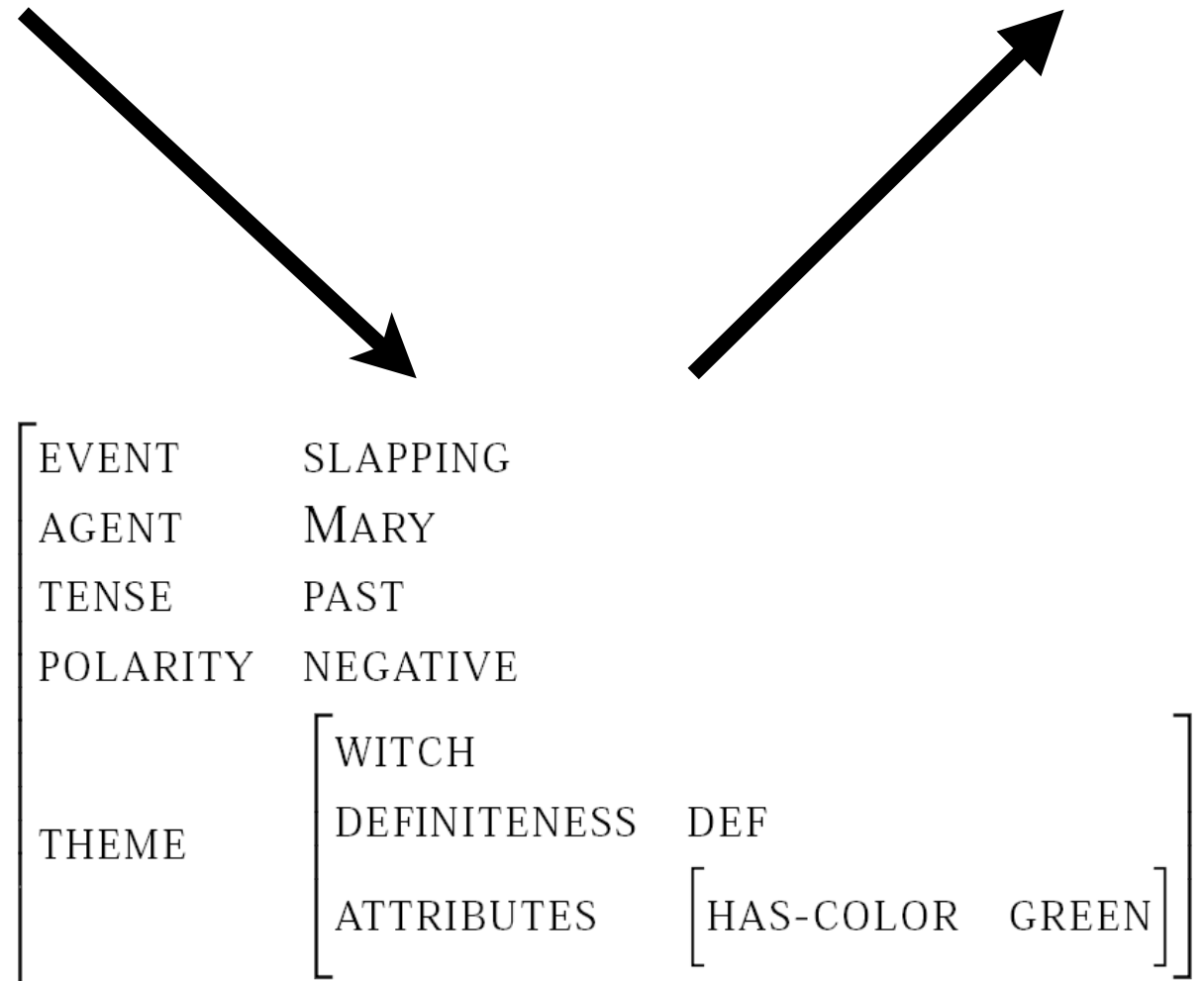




# Interlingua-Based Models

Mary did not slap the  
green witch.

Maria no dió una bofetada  
a la bruja verde





# Expert-Driven Systems

Linguistic grammar formalisms

Handcrafted rules

- $VP \rightarrow PP[+Goal]V \Rightarrow VP \rightarrow V PP[+Goal]$
- ...

Preference mechanisms

- more specific rules first (covering exceptions over more general abstract rules)
- prefer phrases over single word entries in dictionaries
- prefer domain-specific terminology over general vocab.



# Expert-Driven Systems

## ALPAC report (1966)

- MT quality is too low
- No advantage of MT over human translation
- Almost all funding was stopped in the U.S.

## Problems:

- Robustness and translation speed
- Expensive development
- Not very flexible (new domains, languages ...)
- Static, categorical models, but languages are dynamic and ambiguous
- Slow development



# Learning to Translate

## Induce translation knowledge from data

- existing translations
- existing language resources

## Learning Transfer Systems

- data-driven parsing
- induced transfer rules
- language modeling for generation

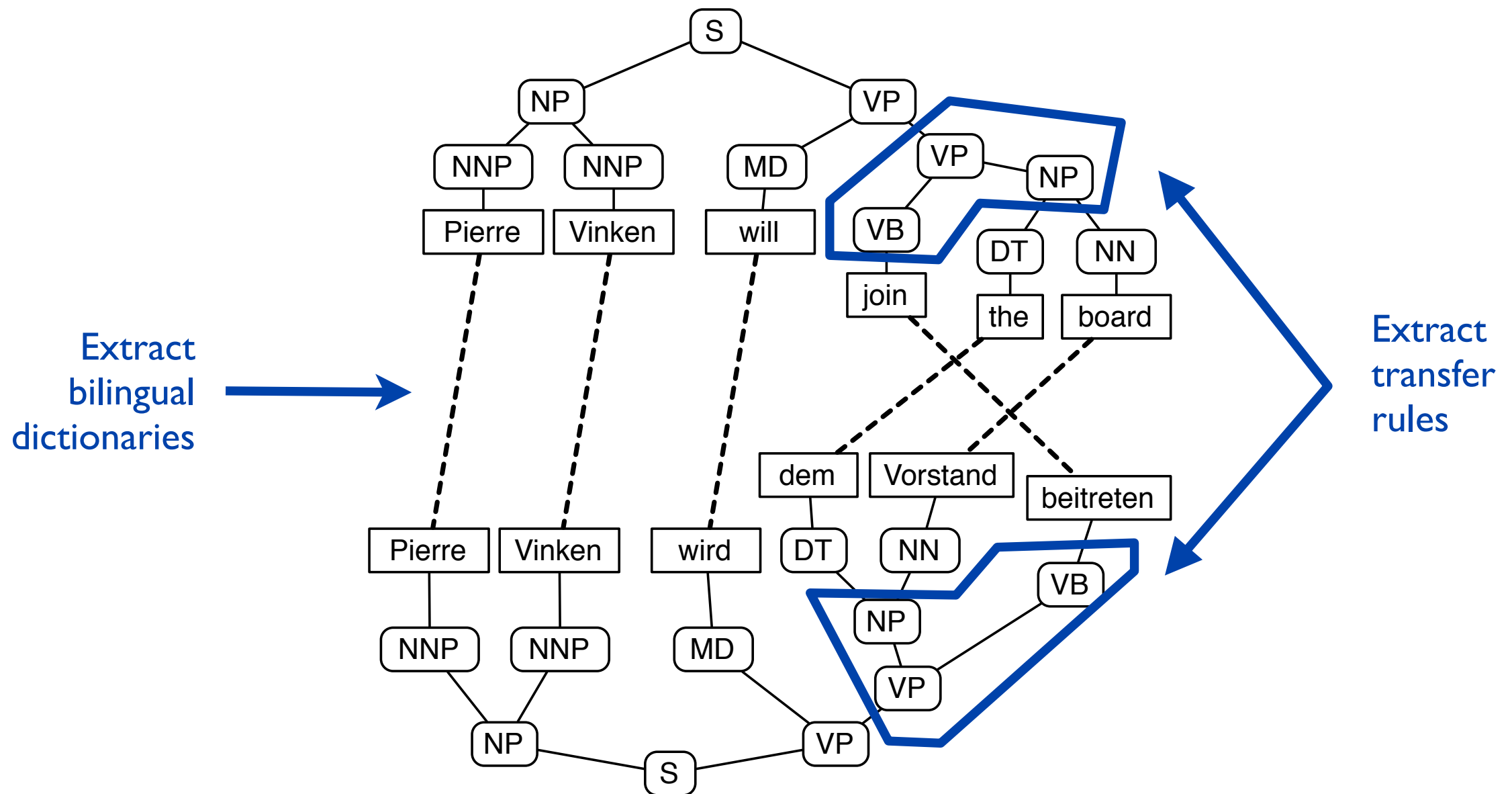
## Global Data-Driven MT

- surface-to-surface MT
- abstraction within a global translation model



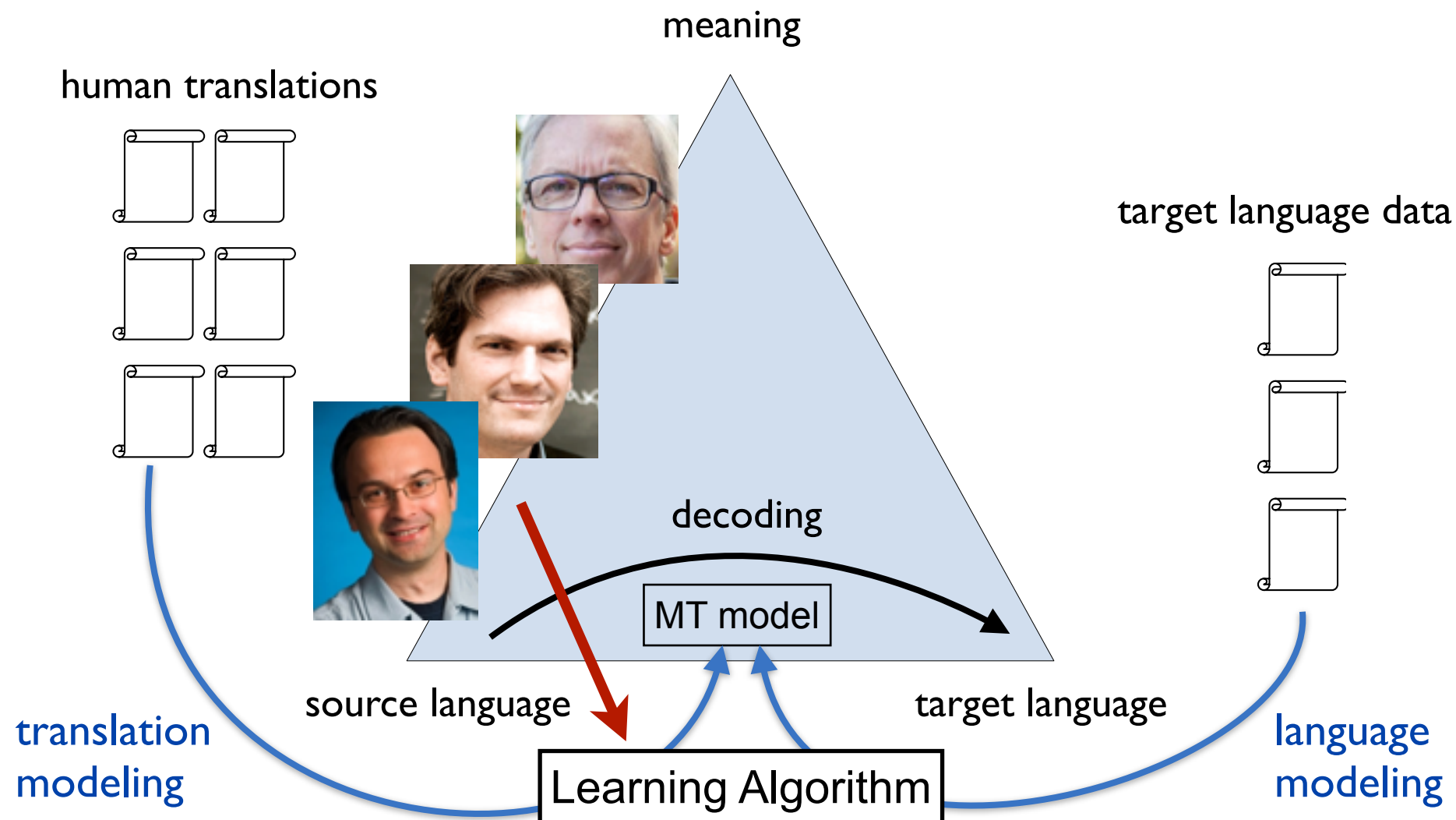
# Learning Transfer-Based MT

Possibility: Induce transfer rules from aligned data





# Data-Driven Machine Translation





# Statistical Machine Translation

1947: MT as decoding (Warren Weaver)

1988: Word-based models

1999: Public implementation of alignment models (GIZA)

2003: Phrase-based SMT

2004: Public phrase-based decoder (Pharaoh)

2005: Hierarchical models

2007: Moses (end-to-end SMT toolbox)

2014: Neural machine translation

along with many tools, much more data and better **computers**



# Advantages of Data-Driven MT

## Human Translations Naturally Appear

- no need for artificial annotation
- can be provided by non-experts

## Implicit Linguistics

- translation knowledge is in the data
- distributional relations within and across languages

## Constant Learning is Possible

- feed with new data as they appear
- quickly adapt to new domains and language pairs





# Take-Home Messages

Machine Translation is important

Machine Translation is difficult

Main Rule-Based Approaches

- Transfer-Based MT
- Interlingua-Based MT
- Direct Translation Systems

Expert-Driven Systems with Hand-Crafted Rules

Data-Driven Systems based on Example Data



# Outlook

This afternoon: MT evaluation

Next week: MT in practice

- Guest lecture, Convertus (Commercial MT solutions in Uppsala)
- Lab 1: evaluation

Then:

- Introduction to SMT
- Lab 2: word-based models
- Guest lecture, Semantix