

# *Machine Translation for Related Languages*

Anoop Kunchukuttan

<https://www.cse.iitb.ac.in/~anoopk>

*Microsoft AI and Research, Hyderabad*



*BITS Pilani, Hyderabad Campus, February 2019*

# Outline

- Introduction to Statistical Machine Translation
- Introduction to Neural Machine Translation
- Machine Translation for Related Languages
- Multilingual Learning

# *Automatic conversion of text/speech from one natural language to another*

*Be the change you want to see in the world*

*वह परिवर्तन बनो जो संसार में देखना चाहते हो*



**Government:** administrative requirements, education, security.

**Enterprise:** product manuals, customer support

**Social:** travel (signboards, food), entertainment (books, movies, videos)

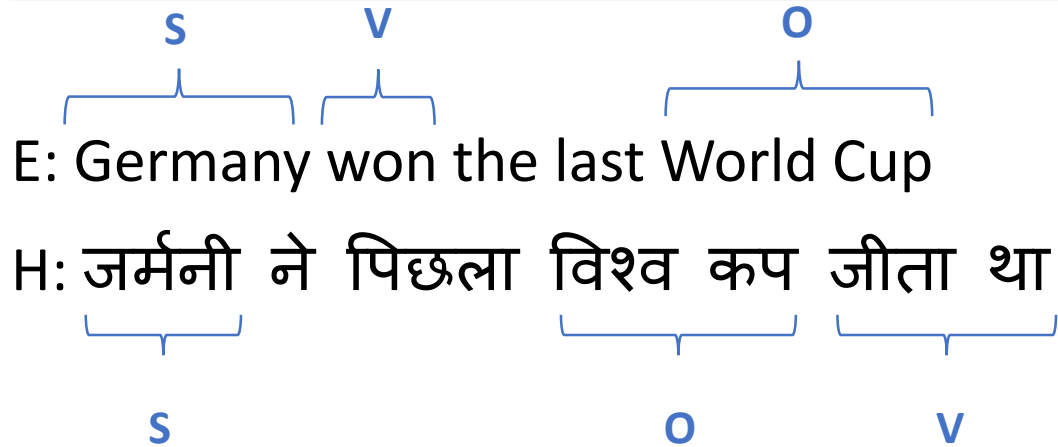
## **Translation under the hood**

- Cross-lingual Search
- Cross-lingual Summarization
- Building multilingual dictionaries

*Any multilingual NLP system will involve some kind of machine translation at some level*

# What is Machine Translation?

## Word order: SOV (Hindi), SVO (English)

  
E: Germany won the last World Cup  
H: जर्मनी ने पिछला विश्व कप जीता था

## Free (Hindi) vs rigid (English) word order

पिछला विश्व कप जर्मनी ने जीता था (correct)

The last World Cup Germany won (grammatically incorrect)

The last World Cup won Germany (meaning changes)

*Language Divergence → the great diversity among languages of the world*

*The central problem of MT is to bridge this language divergence*

# *Why is Machine Translation difficult?*

- **Ambiguity**

- Same word, multiple meanings: मंत्री (minister or chess piece)
- Same meaning, multiple words: जल, पानी, नीर (water)

- **Word Order**

- Underlying deeper syntactic structure
- Phrase structure grammar?
- Computationally intensive

- **Morphological Richness**

- Identifying basic units of words

# *Why should you study Machine Translation?*

- One of the most challenging problems in Natural Language Processing
- Pushes the boundaries of NLP
- Involves analysis as well as synthesis
- Involves all layers of NLP: morphology, syntax, semantics, pragmatics, discourse
- *Theory and techniques in MT are applicable to a wide range of other problems like transliteration, speech recognition and synthesis, and other NLP problems.*

# *Approaches to build MT systems*

Knowledge based, Rule-based MT

*Transfer-based*

*Interlingua based*

Data-driven, Machine Learning based MT

*Example-based*

*Statistical*

*Neural*

# Statistical Machine Translation



## Let's formalize the translation process

We will model translation using a **probabilistic model**. Why?

- We would like to have a measure of confidence for the translations we learn
- We would like to model uncertainty in translation

$E$ : target language

$F$ : source language

$e$ : source language sentence

$f$ : target language sentence

Best  
translation

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

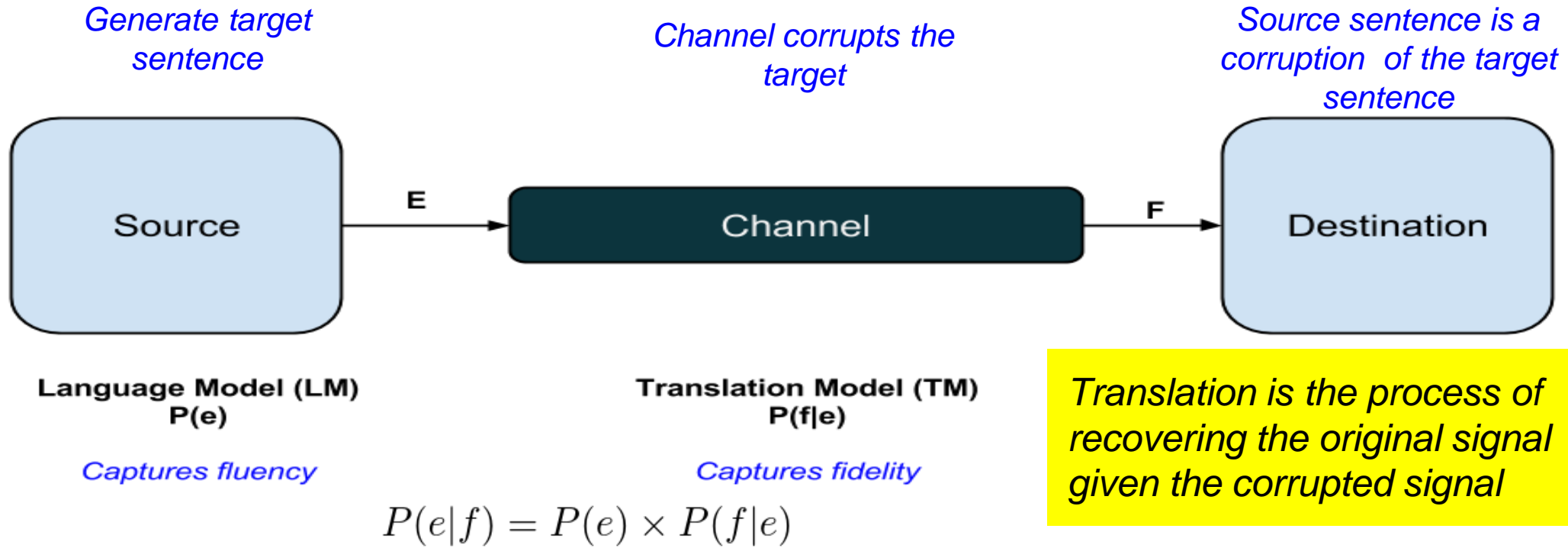
How do we  
**model** this  
quantity?

**Model:** a simplified and idealized understanding of a physical process

We must first explain the process of translation

A very general framework  
for many NLP problems

We explain translation using the **Noisy Channel Model**



Why use this counter-intuitive way of explaining translation?

- Makes it easier to mathematically represent translation and learn probabilities
- **Fidelity** and **Fluency** can be modelled separately

*Let's assume we know how to learn n-gram language models*

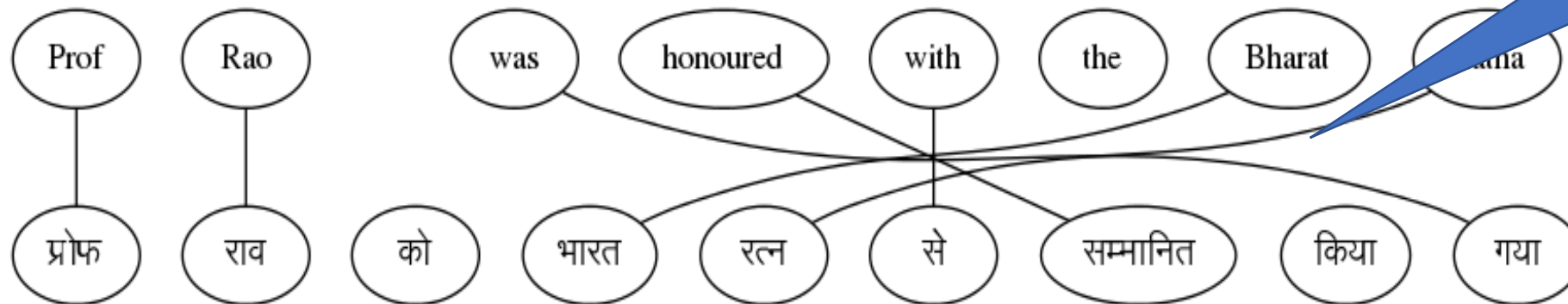
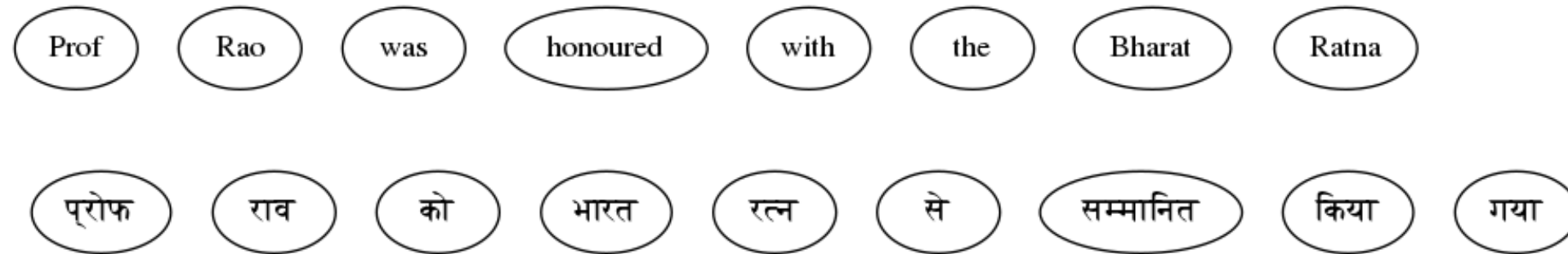
*Let's see how to learn the translation model  $\rightarrow P(f|e)$*

***To learn sentence translation probabilities,***

***$\rightarrow$  we first need to learn word-level translation probabilities***

*That is the task of word alignment*

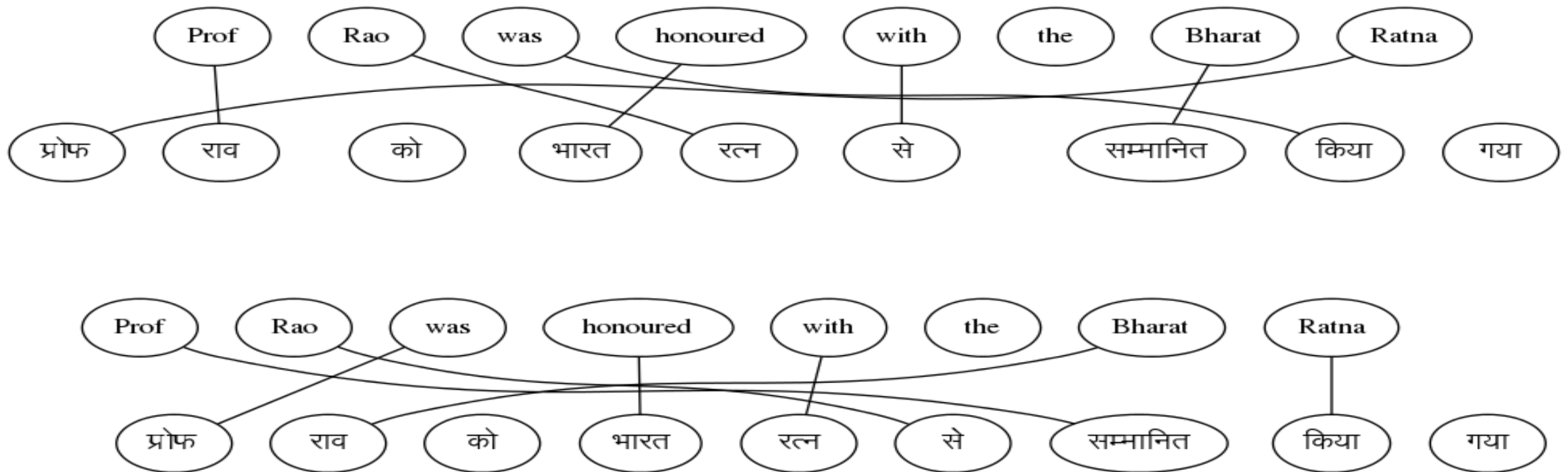
*Given a parallel sentence pair, find word level correspondences*



This set of links for a sentence pair is called an 'ALIGNMENT'

*But there are multiple possible alignments*

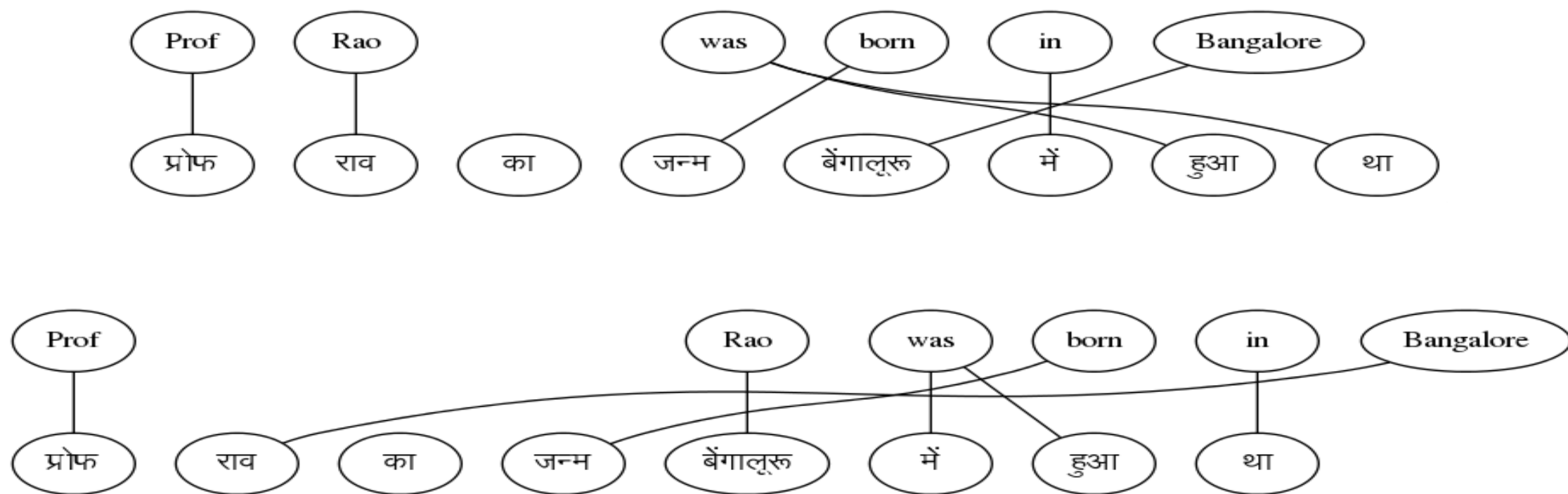
### Sentence 1



*With one sentence pair, we cannot find the correct alignment*

*Can we find alignments if we have multiple sentence pairs?*

## Sentence 2



*Yes, let's see how to do that ...*

## Parallel Corpus

A boy is sitting in the kitchen	एक लडका रसोई में बैठा है
A boy is playing tennis	एक लडका टेनिस खेल रहा है
A boy is sitting on a round table	एक लडका एक गोल मेज पर बैठा है
Some men are watching tennis	कुछ आदमी टेनिस देख रहे हैं
A girl is holding a black book	एक लडकी ने एक काली किताब पकड़ी है
Two men are watching a movie	दो आदमी चलचित्र देख रहे हैं
A woman is reading a book	एक औरत एक किताब पढ़ रही है
A woman is sitting in a red car	एक औरत एक काले कार में बैठी है

## Parallel Corpus

A boy is <b>sitting</b> in the kitchen	एक लडका रसोई में <b>बैठा</b> है
A boy is playing <b>tennis</b>	एक लडका <b>टेनिस</b> खेल रहा है
A boy is <b>sitting</b> on a round table	एक लडका एक गोल मेज पर <b>बैठा</b> है
Some men <b>are watching tennis</b>	कुछ आदमी <b>टेनिस देख रहे हैं</b>
A girl is holding a black book	एक लडकी ने एक काली किताब पकड़ी है
Two men <b>are watching</b> a movie	दो आदमी चलचित्र <b>देख रहे हैं</b>
A woman is reading a book	एक औरत एक किताब पढ़ रही है
A woman is <b>sitting</b> in a red car	एक औरत एक काले कार में <b>बैठा</b> है

### Key Idea

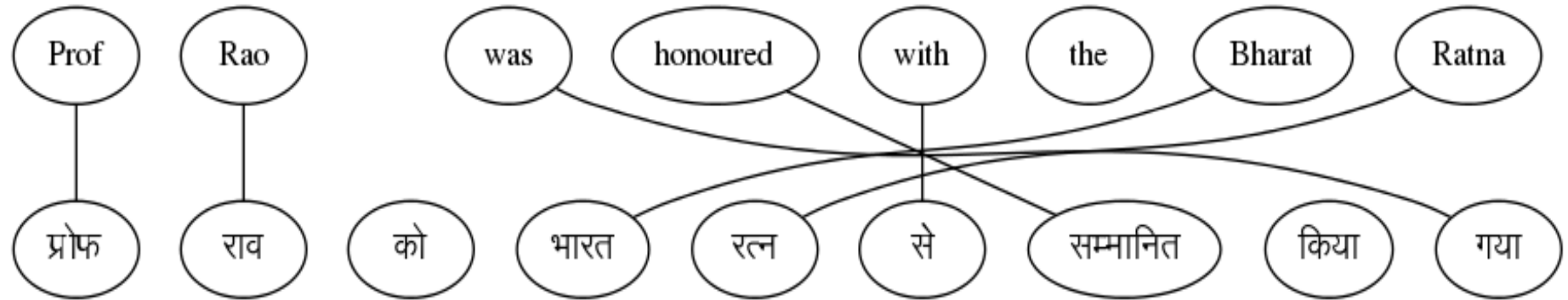
*Co-occurrence of translated words*

*Words which occur together in the parallel sentence are likely to be translations (higher  $P(f|e)$ )*

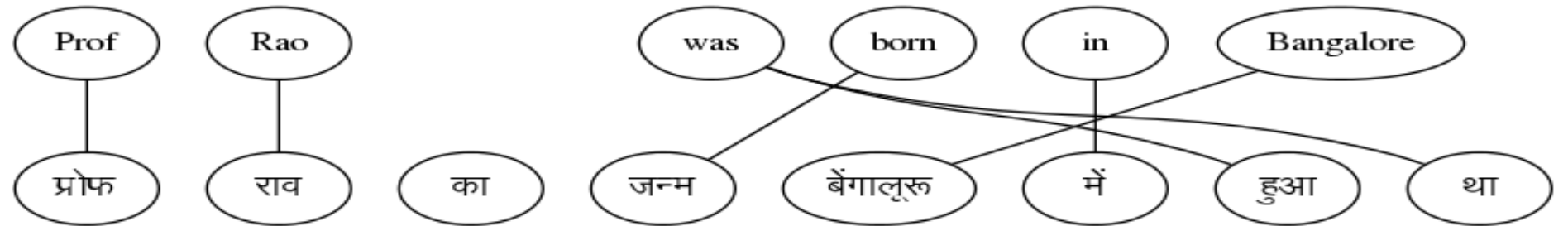


*If we knew the alignments, we could compute  $P(f|e)$*

Sentence 1



Sentence 2



$$P(f|e) = \frac{\#(f, e)}{\#(*, e)}$$

$$P(\text{Prof} | \text{प्रोफ}) = \frac{2}{2}$$

$\#(a, b)$ : number of times word  $a$  is aligned to word  $b$

*But, we can find the best alignment only if we know the word translation probabilities*

*The best alignment is the one that maximizes the sentence translation probability*

$$P(\mathbf{f}, \mathbf{a} | \mathbf{e}) = P(a) \prod_{i=1}^{i=m} P(f_i | e_{a_i}) \quad \longrightarrow \quad \mathbf{a}^* = \operatorname{argmax}_{\mathbf{a}} \prod_{i=1}^{i=m} P(f_i | e_{a_i})$$

*This is a chicken and egg problem! How do we solve this?*

# *We can solve this problem using a two-step, iterative process*

*Start with random values for word translation probabilities*

*Step 1: Estimate alignment probabilities using word translation probabilities*

*Step 2: Re-estimate word translation probabilities*

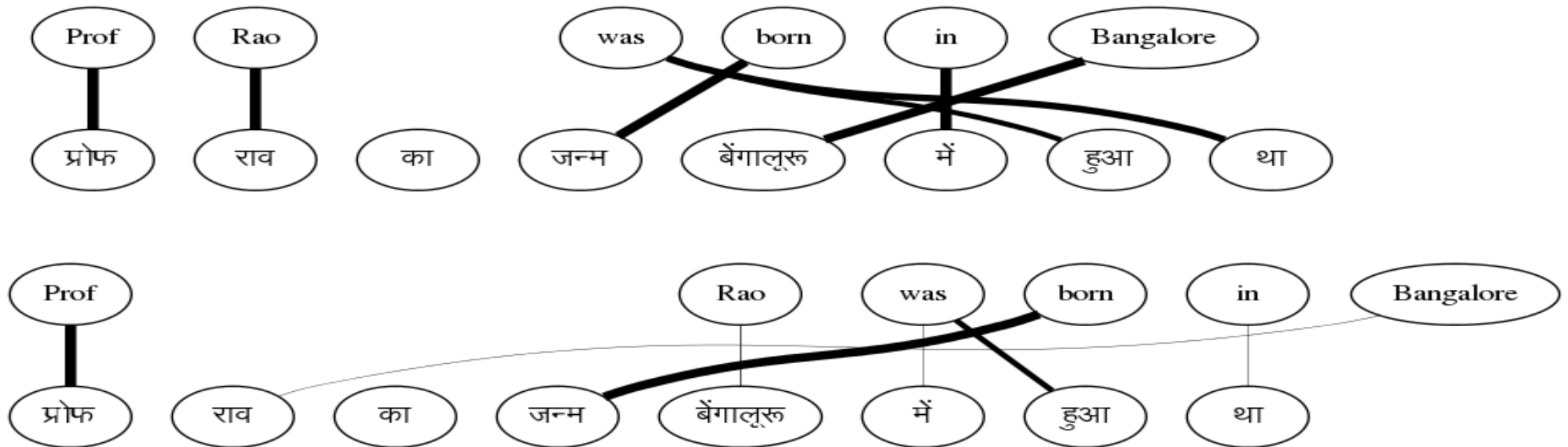
- We don't know the best alignment*
- So, we consider all alignments while estimating word translation probabilities*
- Instead of taking only the best alignment, we consider all alignments and weigh the word alignments with the alignment probabilities*

$$P(f|e) = \frac{\text{expected } \#(f, e)}{\text{expected } \#(*, e)}$$

*Repeat Steps (1) and (2) till the parameters converge*

*At the end of the process ...*

## Sentence 2



**Expectation-Maximization Algorithm:** guaranteed to converge, maybe to local minima  
Hence we need to good initialization and training regimens.

# IBM Models

- IBM came up with a series of increasingly complex models
- Called Models 1 to 5
- Differed in assumptions about alignment probability distributions
- Simpler models are used to initialize the more complex models
- This pipelined training helped ensure better solutions

# Phrase Based SMT

Why stop at learning word correspondences?

KEY IDEA → Use “Phrase” (Sequence of Words) as the basic translation unit

*Note: the term ‘phrase’ is not used in a linguistic sense*

The Prime Minister of India	भारत के प्रधान मंत्री bhArata ke pradhAna maMtrl India of Prime Minister
is running fast	तेज भाग रहा है teja bhAg rahA hai fast run -continuous is
honoured with	से सम्मानित किया se sammanita kiyA with honoured did
Rahul lost the match	राहुल मुकाबला हार गया rAhula mukAbala hAra gayA Rahul match lost

# Benefits of PB-SMT

Local Reordering → Intra-phrase re-ordering can be memorized

The Prime Minister of India	भारत के प्रधान मंत्री bhaarat ke pradhaan maMtrl India of Prime Minister
-----------------------------	--

Sense disambiguation based on local context → Neighbouring words help make the choice

heads towards Pune	पुणे की ओर जा रहे हैं pune ki or jaa rahe hai Pune towards go –continuous is
heads the committee	समिति की अध्यक्षता करते हैं Samiti kii adhyakshata karte hai committee of leading - verbalizer is

# Benefits of PB-SMT (2)

## Handling institutionalized expressions

- Institutionalized expressions, idioms can be learnt as a single unit

hung assembly	त्रिशंकु विधानसभा trishanku vidhaansabha
Home Minister	गृह मंत्री gruh mantrii
Exit poll	चुनाव बाद सर्वेक्षण chunav baad sarvekshana

- Improved Fluency

- The phrases can be arbitrarily long (even entire sentences)



# Mathematical Model

Let's revisit the decision rule for SMT model

$$\begin{aligned} \mathbf{e}_{\text{best}} &= \operatorname{argmax}_{\mathbf{e}} p(\mathbf{e}|\mathbf{f}) \\ &= \operatorname{argmax}_{\mathbf{e}} p(\mathbf{f}|\mathbf{e}) p_{\text{LM}}(\mathbf{e}) \end{aligned}$$

Let's revisit the translation model  $p(\mathbf{f}|\mathbf{e})$

- Source sentence can be segmented in  $\mathbf{I}$  phrases
- Then,  $p(\mathbf{f}|\mathbf{e})$  can be decomposed as:

$$p(\bar{\mathbf{f}}_1^I | \bar{\mathbf{e}}_1^I) = \prod_{i=1}^I \phi(\bar{f}_i | \bar{e}_i) d(\text{start}_i - \text{end}_{i-1} - 1)$$

Distortion  
probability

Phrase Translation  
Probability

$\text{start}_i$  :start position in  $\mathbf{f}$  of  $i^{\text{th}}$  phrase of  $\mathbf{e}$   
 $\text{end}_i$  :end position in  $\mathbf{f}$  of  $i^{\text{th}}$  phrase of  $\mathbf{e}$

# Learning The Phrase Translation Model

Involves Structure + Parameter Learning:

- Learn the **Phrase Table**: the central data structure in PB-SMT

The Prime Minister of India	भारत के प्रधान मंत्री
is running fast	तेज भाग रहा है
the boy with the telescope	दूरबीन से लड़के को
Rahul lost the match	राहुल मुकाबला हार गया

- Learn the **Phrase Translation Probabilities**

Prime Minister of India	भारत के प्रधान मंत्री India of Prime Minister	0.75
Prime Minister of India	भारत के भूतपूर्व प्रधान मंत्री India of former Prime Minister	0.02
Prime Minister of India	प्रधान मंत्री Prime Minister	0.23

# Learning Phrase Tables from Word Alignments

- Start with word alignments
- Word Alignment : reliable input for phrase table learning
  - high accuracy reported for many language pairs
- Central Idea: A consecutive sequence of aligned words constitutes a “phrase pair”

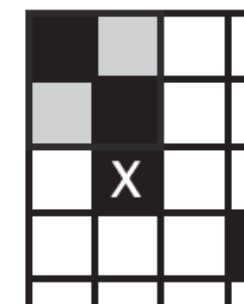
	Prof	C.N.R.	Rao	was	honoured	with	the	Bharat	Ratna
प्रोफेसर	■								
सी.एन.आर		■	■						
राव			■						
को									
भारतरत्न								■	■
से							■		
सम्मानित					■	■			
किया									
गया									

Which phrase pairs to include in the phrase table?

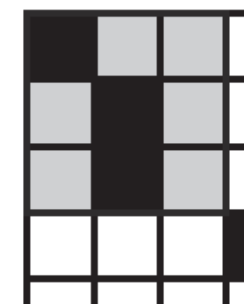
	Prof	C.N.R.	Rao	was	honoured	with	the	Bharat	Ratna
प्रोफेसर									
सी.एन.आर									
राव									
को									
भारतरत्न									
से									
सम्मानित									
किया									
गया									



consistent



inconsistent



consistent



Source: SMT, Phillip Koehn

Professor CNR	प्रोफेसर सी.एन.आर
Professor CNR Rao	प्रोफेसर सी.एन.आर राव
Professor CNR Rao was	प्रोफेसर सी.एन.आर राव
Professor CNR Rao was	प्रोफेसर सी.एन.आर राव को
honoured with the Bharat Ratna	भारतरत्न से सम्मानित
honoured with the Bharat Ratna	भारतरत्न से सम्मानित किया
honoured with the Bharat Ratna	भारतरत्न से सम्मानित किया गया
honoured with the Bharat Ratna	को भारतरत्न से सम्मानित किया गया

# Discriminative Training of PB-SMT

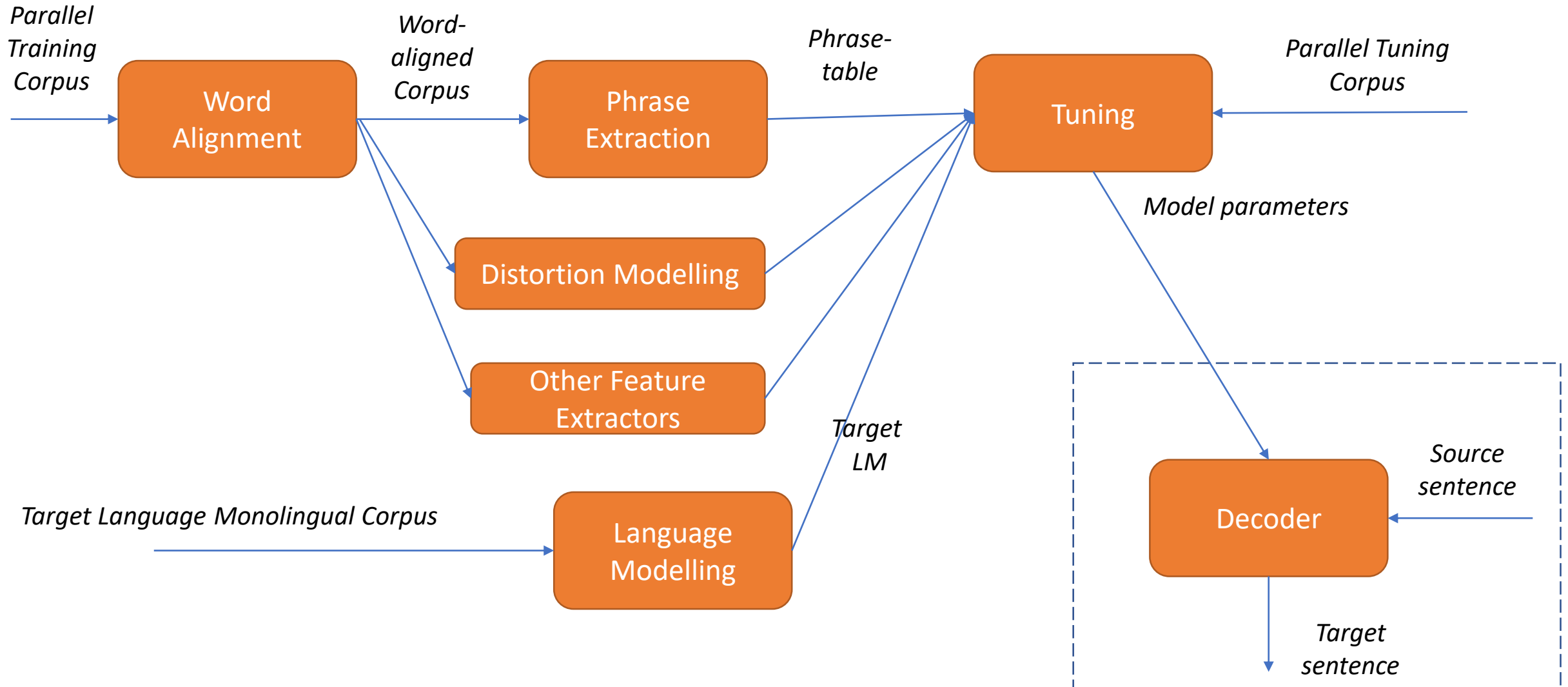
- Directly model the posterior probability  $p(\mathbf{e} | \mathbf{f})$
- Use the Maximum Entropy framework

$$P(\mathbf{e} | \mathbf{f}) = \exp \left( \sum_i \lambda_i h_i(f_1^I, e_1^J) \right)$$

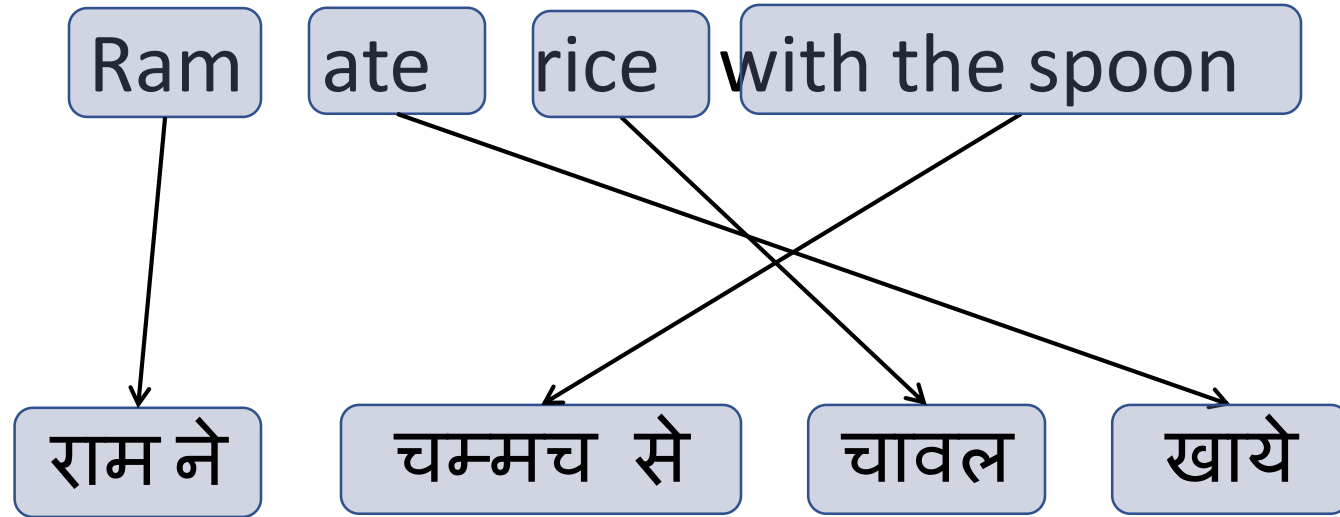
$$e^* = \arg \max_{e_i} \sum_i \lambda_i h_i(f_1^I, e_1^J)$$

- $h_i(\mathbf{f}, \mathbf{e})$  are feature functions ,  $\lambda_i$ 's are feature weights
- Benefits:
  - *Can add arbitrary features to score the translations*
  - Can assign different weight for each features
  - Assumptions of generative model may be incorrect
  - *Feature weights  $\lambda_i$  are learnt during tuning*

# Typical SMT Pipeline



# Decoding

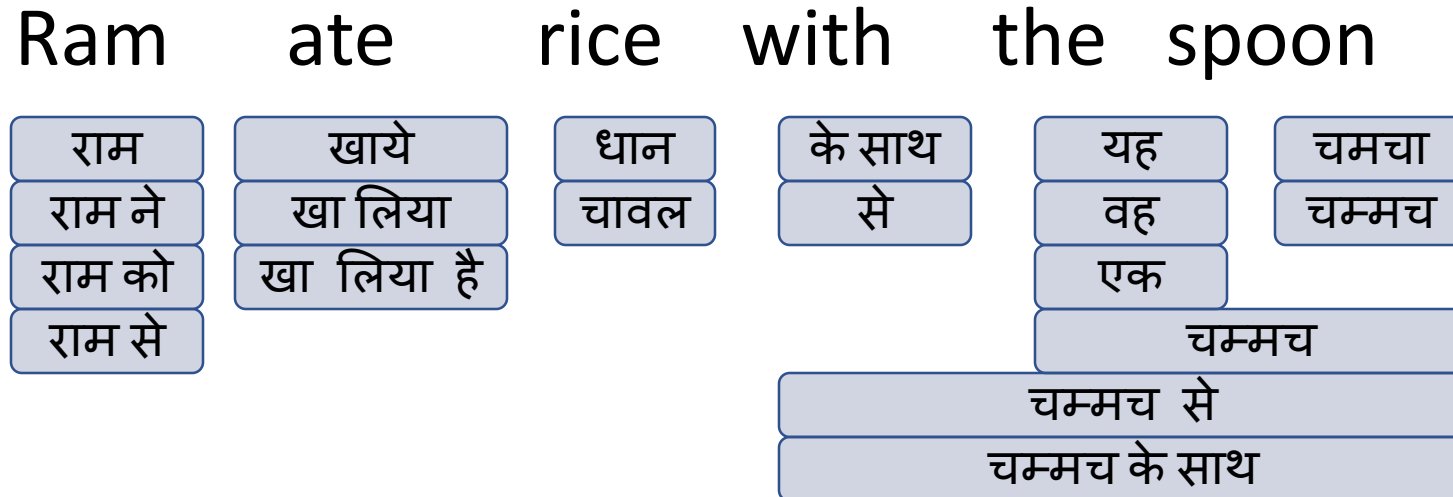


Searching for the best translations in the space of all translations

$$e^* = \arg \max_{e_i} \sum_i \lambda_i h_i(f_1^I, e_1^J)$$

# Decoding is challenging

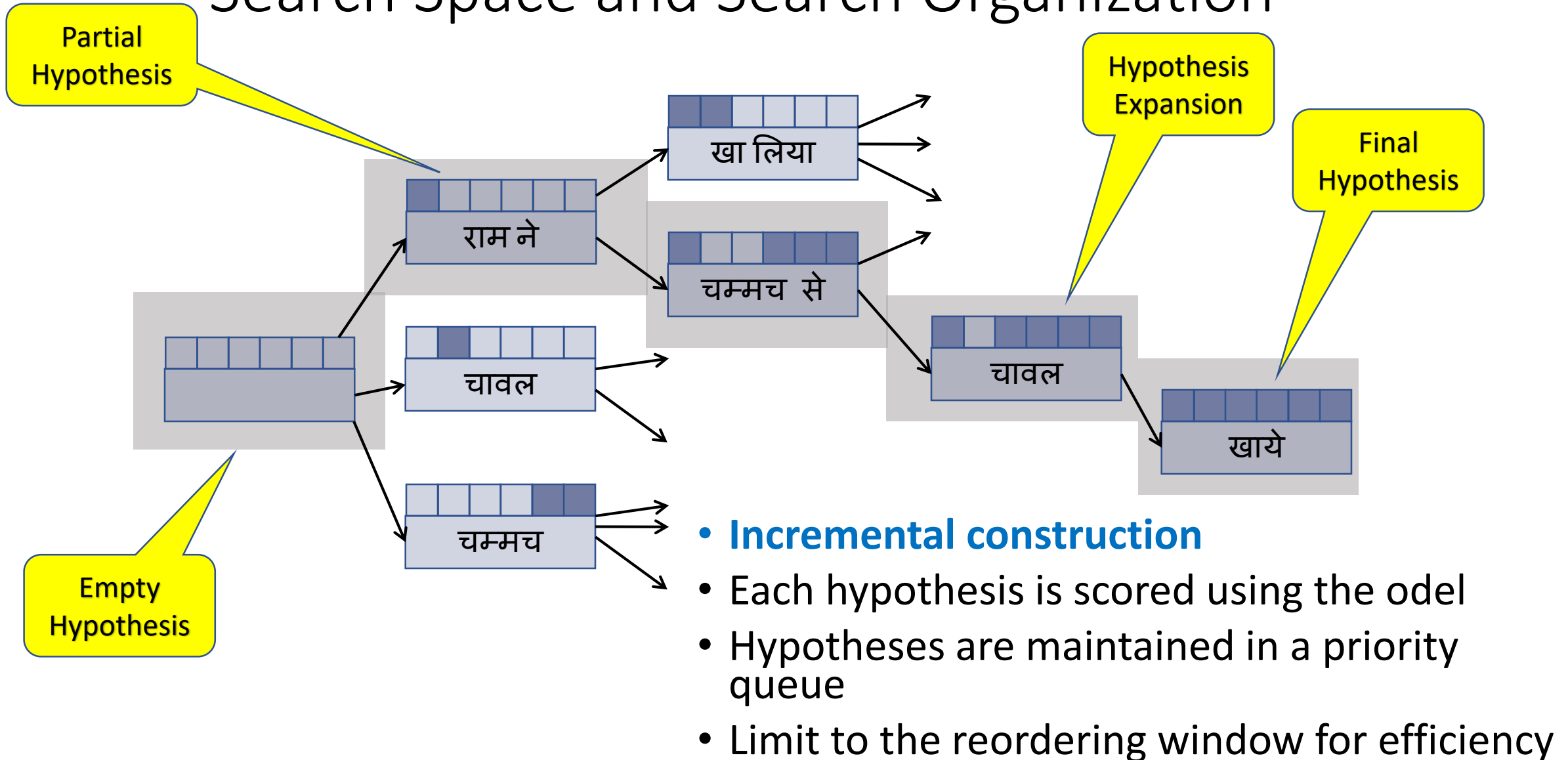
- We picked the phrase translation that made sense to us
- The computer has less intuition
- Phrase table may give many options to translate the input sentence
- Multiple possible word orders



An NP complete search problem → Needs a heuristic search method



# Search Space and Search Organization



*We have looked at a basic phrase-based SMT system*

*This system can learn word and phrase translations from parallel corpora*

But many important linguistic phenomena need to be handled

- **Divergent Word Order**
- Rich morphology
- Named Entities and Out-of-Vocabulary words

# Getting word order right

*Phrase based MT is not good at learning word ordering*

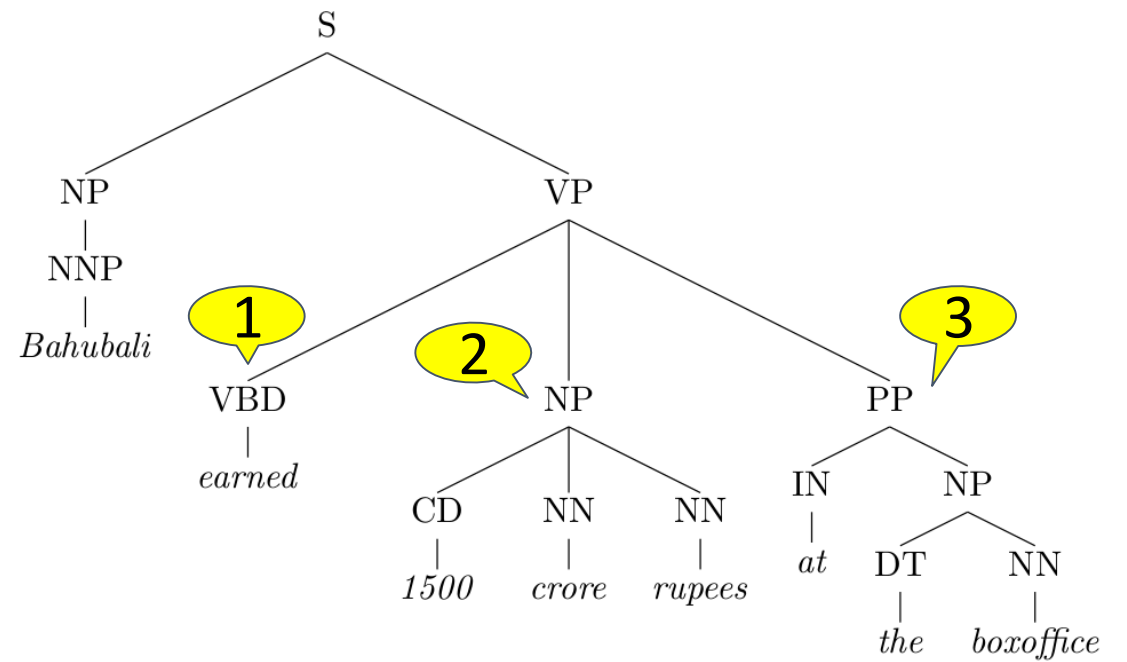
*Solution: Let's help PB-SMT with some preprocessing of the input*

*Change order of words in input sentence to match order of the words in the target language*

Let's take an example

*Bahubali earned more than 1500 crore rupee sat the boxoffice*

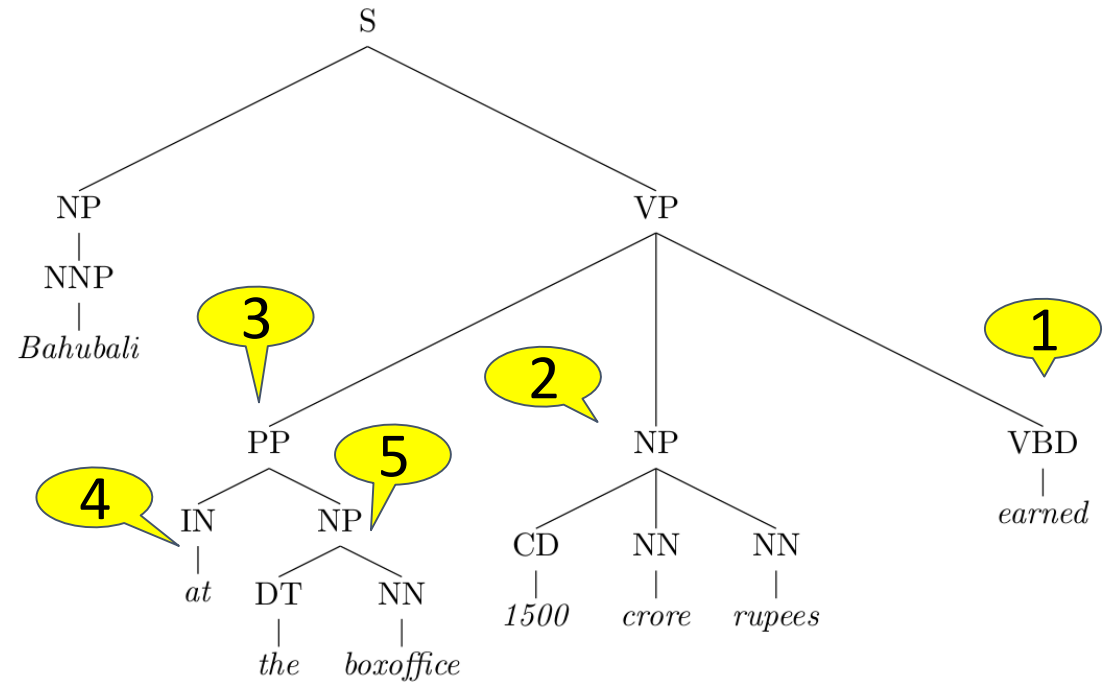
*Parse the sentence to understand its syntactic structure*



*Apply rules to transform the tree*

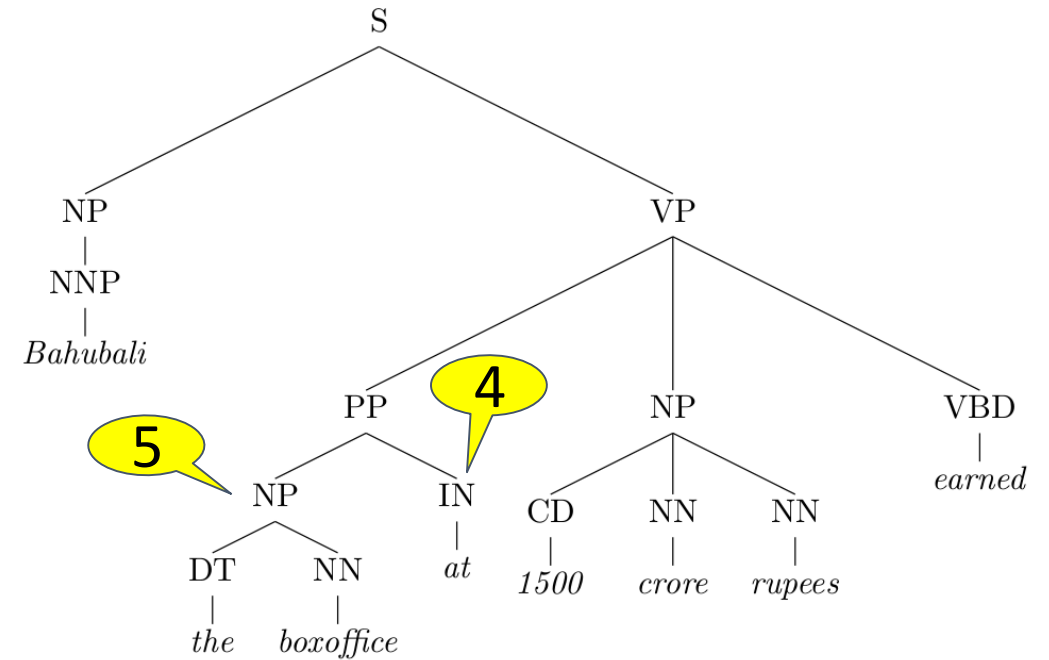
**VP → VBD NP PP ⇒ VP → PP NP VBD**

This rule captures  
Subject-Verb-Object to Subject-  
Object-Verb divergence



*Prepositions in English become postpositions in Hindi*

PP → IN NP ⇒ PP → NP IN



*The new input to the machine translation system is*  
*Bahubali the boxoffice at 1500 crore rupees earned*

*Now we can translate with little reordering*

*बाहुबली ने बॉक्सओफिस पर 1500 करोड रुपए कमाए*

*These rules can be  
written manually or  
learnt from parse trees*

# Addressing Rich Morphology

Inflectional forms of the Marathi word घर

घर	house
घरात	in the house
घरावरती	on the house
घराखाली	below the house
घरामध्ये	in the house
घरामागे	behind the house
घराचा	of the house
घरामागचा	that which is behind the house
घरासमोर	in front of the house
घरासमोरचा	that which is in front of the house
घरांसमोर	in front of the houses

Hindi words with the suffix वाद

साम्यवाद	communism
समाजवाद	socialism
पूंजीवाद	capitalism
जातीवाद	casteism
साम्राज्यवाद	imperialism

*The corpus should contains all variants to learn translations*

*This is infeasible!*

***Language is very productive, you can combine words to generate new words***

# Addressing Rich Morphology

Inflectional forms of the Marathi word घर

घर	house
घर ा त	in the house
घर ा वरती	on the house
घर ा खाली	below the house
घर ा मध्ये	in the house
घर ा मागे	behind the house
घर ा चा	of the house
घर ा माग चा	that which is behind the house
घर ा समोर	in front of the house
घर ा समोर चा	that which is in front of the house
घर ा ं समोर	in front of the houses

Hindi words with the suffix वाद

साम्य वाद	communism
समाज वाद	socialism
पूंजी वाद	capitalism
जाती वाद	casteism
साम्राज्य वाद	imperialism

- *Break the words into its component morphemes*
- *Learn translations for the morphemes*
- *Far more likely to find morphemes in the corpus*

# Handling Names and OOVs

Some words not seen during train will be seen at test time  
These are *out-of-vocabulary (OOV)* words

**Names** are one of the most important category of OOVs  
⇒ There will always be names not seen during training

How do we translate names like *Sachin Tendulkar* to Hindi?

What we want to do is map the Roman characters to Devanagari to they sound the same when read → सचिन तेंदुलकर

→ We call this process '**transliteration**'

Can be seen as a simple translation problem at character level with no re-ordering

*sachin* → सच िन



# Evaluation of MT output

- How do we judge a good translation?
- Can a machine do this?
- Why should a machine do this?
  - Because human evaluation is time-consuming and expensive!
  - Not suitable for rapid iteration of feature improvements

# What is a good translation?

Evaluate the quality with respect to:

- **Adequacy:** How good the output is in terms of preserving content of the source text
- **Fluency:** How good the output is as a well-formed target language entity

**For example,** I am attending a lecture

मैं एक व्याख्यान बैठा हूँ  
*Main ek vyaakhyan baitha hoon*  
*I a lecture sit (Present-first person)*  
*I sit a lecture : Adequate but not fluent*

मैं व्याख्यान हूँ  
*Main vyakhyan hoon*  
*I lecture am*  
*I am lecture: Fluent but not adequate.*

# Human Evaluation

## Direct Assessment

How do you rate your Olympic experience?

— Reference

How do you value the Olympic experience?

— Candidate translation

### Adequacy:

Is the meaning translated correctly?

5 = All

4 = Most

3 = Much

2 = Little

1 = None

### Fluency:

Is the sentence grammatically valid?

5 = Flawless

4 = Good

3 = Non-native

2 = Disfluent

1 = Incomprehensible

## Ranking Translations

[Appraise](#) [Overview](#) [Status](#) cfedermann ▾

Până la mijlocul lui iulie, procentul a urcat la 40%. La începutul lui august, era 52%.

— Source

By mid-July, it was 40 percent. In early August, it was 52 percent.

— Reference

Best ← Rank 1 ● Rank 2 ● Rank 3 ● Rank 4 ● Rank 5 ● → Worst

Until the middle of July, the percentage rose to 40%.

Best ← Rank 1 ● Rank 2 ● Rank 3 ● Rank 4 ● Rank 5 ● → Worst

Until mid-July, the percentage rose to 40%.

Best ← Rank 1 ● Rank 2 ● Rank 3 ● Rank 4 ● Rank 5 ● → Worst

By mid-July, the percentage climbed to 40 per cent.

Best ← Rank 1 ● Rank 2 ● Rank 3 ● Rank 4 ● Rank 5 ● → Worst

Until mid-July, the percentage climbed to 40%.

Best ← Rank 1 ● Rank 2 ● Rank 3 ● Rank 4 ● Rank 5 ● → Worst

Until the middle of July, the figure climbed to 40%.

$$\text{score}(S_i) = \frac{1}{|\{S\}|} \sum_{S_j \neq S_i} \frac{\text{wins}(S_i, S_j)}{\text{wins}(S_i, S_j) + \text{wins}(S_j, S_i)}$$

# Automatic Evaluation

*Human evaluation is not feasible in the development cycle*

*Key idea of Automatic evaluation:*

*The closer a machine translation is to a professional human translation, the better it is.*

- Given: A corpus of good quality human reference translations
- Output: A numerical “translation closeness” metric
- Given (ref,sys) pair, score =  $f(\text{ref}, \text{sys}) \rightarrow \mathbb{R}$

where,

sys (candidate Translation): Translation returned by an MT system

ref (reference Translation): ‘Perfect’ translation by humans

Multiple references are better

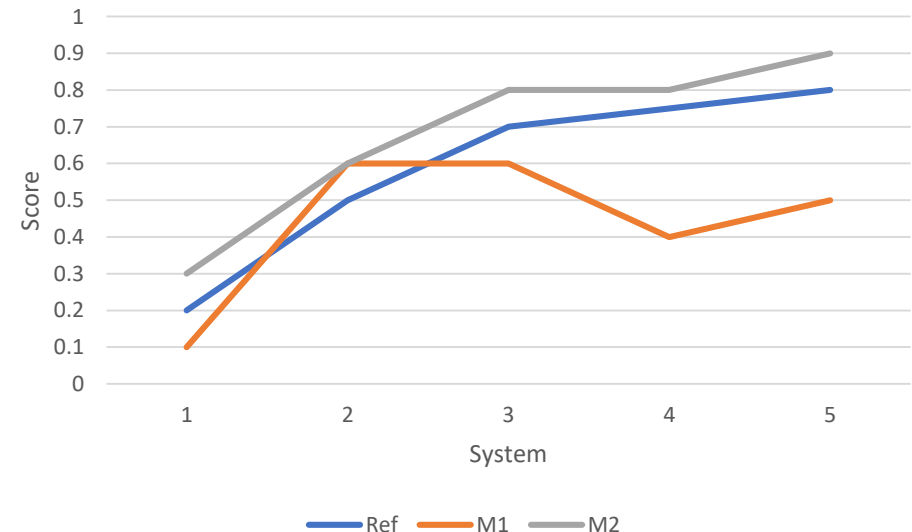
# Some popular automatic evaluation metrics

- BLEU (Bilingual Evaluation Understudy)
- TER (Translation Edit Rate)
- METEOR (Metric for Evaluation of Translation with Explicit Ordering)

How good is an automatic metric?



How well does it correlate with human judgment?



# Neural Machine Translation

**SMT, Rule-based MT and Example based MT** manipulate **symbolic representations** of knowledge

Every word has an atomic representation,  
which can't be further analyzed

**No notion of similarity or relationship between words**

- Even if we know the translation of `home`, we can't translate `house` if it is an OOV

home	0
water	1
house	2
tap	3

1	0	0	0
0	1	0	0
0	0	1	0
0	0	0	1

**Difficult to represent new concepts**

- We cannot say anything about 'mansion' if it comes up at test time
- Creates problems for language model as well  $\Rightarrow$  whole area of smoothing exists to overcome this problem

Symbolic representations are **discrete representations**

- **Generally computationally expensive** to work with discrete representations
- e.g. Reordering requires evaluation of an exponential number of candidates

## Neural Network techniques work with **distributed representations**

Every word is represented by a vector of numbers

- No element of the vector represents a particular word
- The word can be understood with all vector elements
- Hence distributed representation
- But less interpretable

**Can define similarity between words**

- Vector similarity measures like cosine similarity
- Since representations of `home` *and* `house`, we may be able to translate `house`

home
Water
house
tap

0.5	0.6	0.7
0.2	0.9	0.3
0.55	0.58	0.77
0.24	0.6	0.4

Word vectors or embeddings

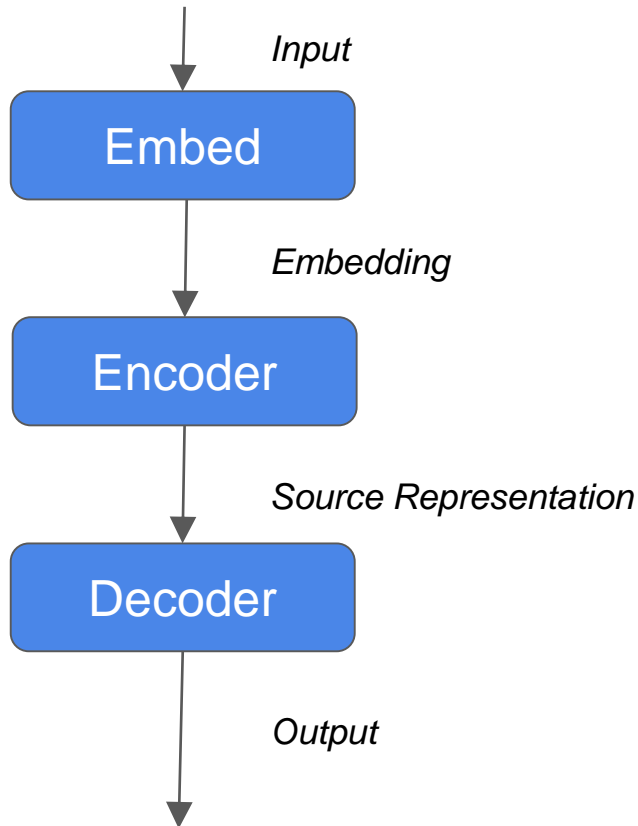
**New concepts can be represented using a vector with different values**

Symbolic representations are **continuous representations**

- **Generally computationally more efficient** to work with continuous values
- Especially optimization problems



# Encode - Decode Paradigm



*Entire input sequence is processed before generation starts  
⇒ In PBSMT, generation was piecewise*

***The input is a sequence of words, processed one at a time***

- *While processing a word, the network needs to know what it has seen so far in the sequence*
- *Meaning, know the history of the sequence processing*
- *Needs a special kind of neural: **Recurrent neural network unit** which can keep state information*

# Encode - Decode Paradigm Explained

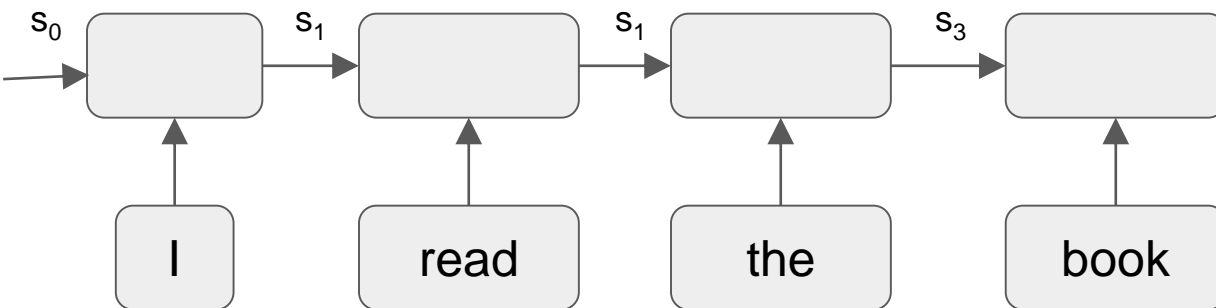
Use two RNN networks: the encoder and the decoder

(1) Encoder processes one sequence at a time

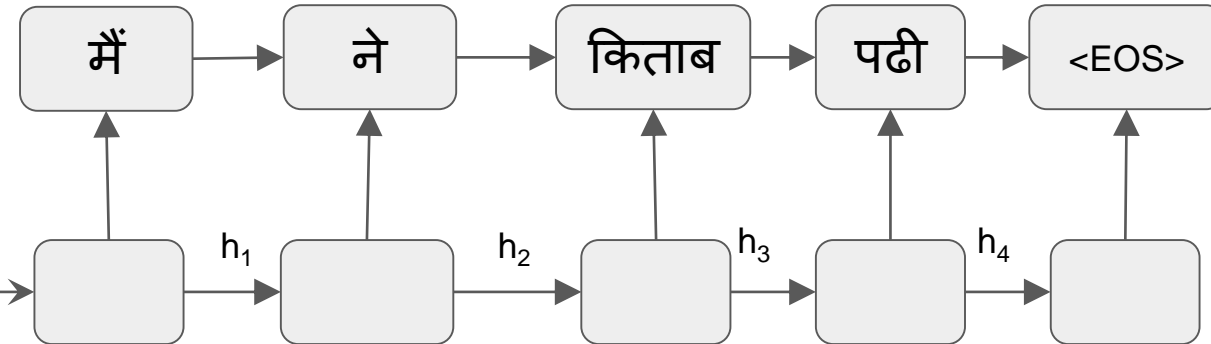
(3) This is used to initialize the decoder state

(4) Decoder generates one element at a time

(5)... continue till end of sequence tag is generated



$h_0$



(2) A representation of the sentence is generated

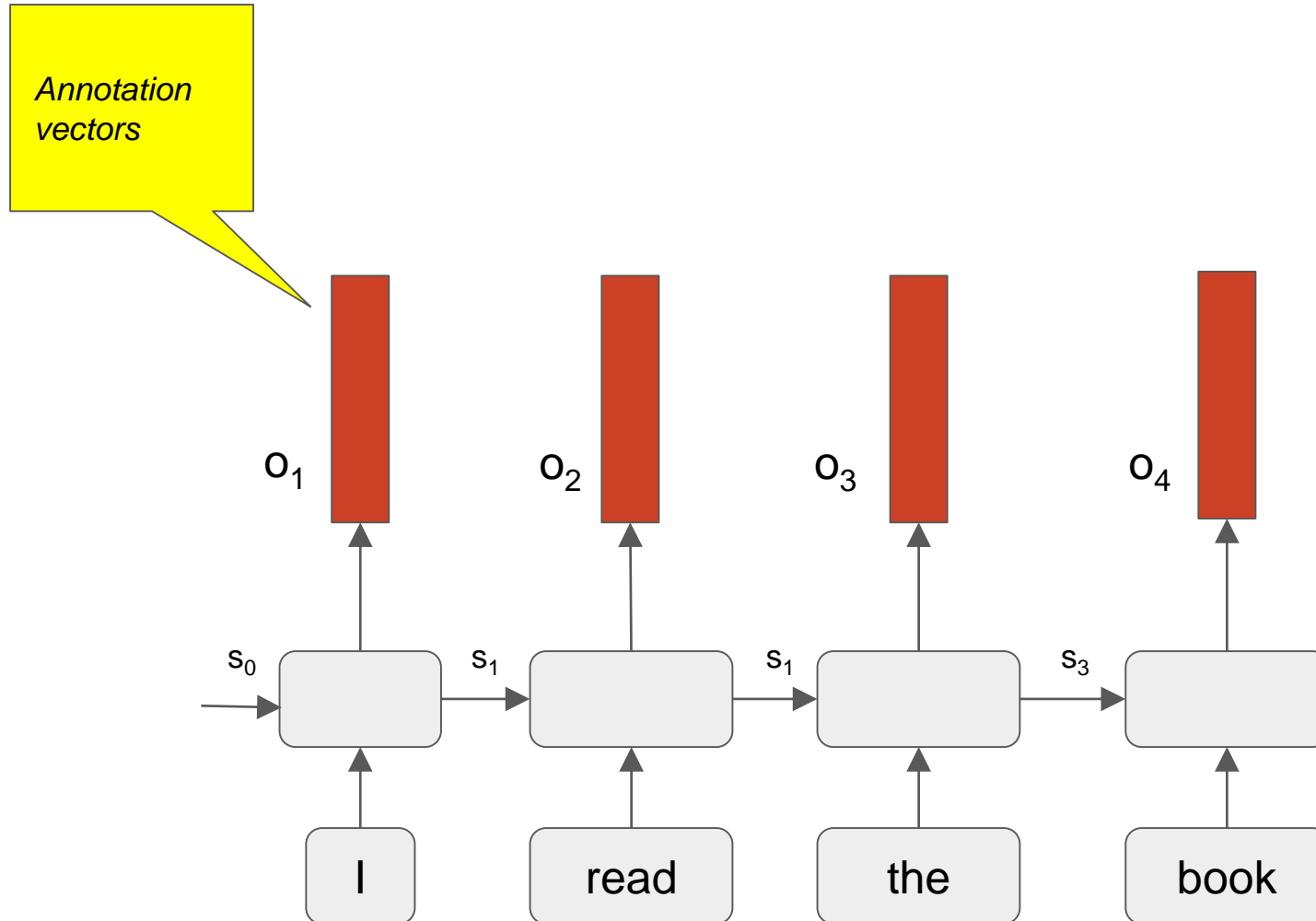
*This approach reduces the entire sentence representation to a single vector*

Two problems with this design choice:

- A single vector is not sufficient to represent to capture all the syntactic and semantic complexities of a sentence
  - *Solution: Use a richer representation for the sentences*
- Problem of capturing long term dependencies: The decoder RNN will not be able to make use of source sentence representation after a few time steps
  - *Solution: Make source sentence information when making the next prediction*
  - *Even better, make **RELEVANT** source sentence information available*

*These solutions motivate the next paradigm*

# Encode - Attend - Decode Paradigm



Represent the source sentence by the **set of output vectors** from the encoder

Each output vector at time  $t$  is a contextual representation of the input at time  $t$

*Note: in the encoder-decode paradigm, we ignore the encoder outputs*

Let's call these encoder output vectors **annotation vectors**

*How should the decoder use the set of annotation vectors while predicting the next character?*

**Key Insight:**

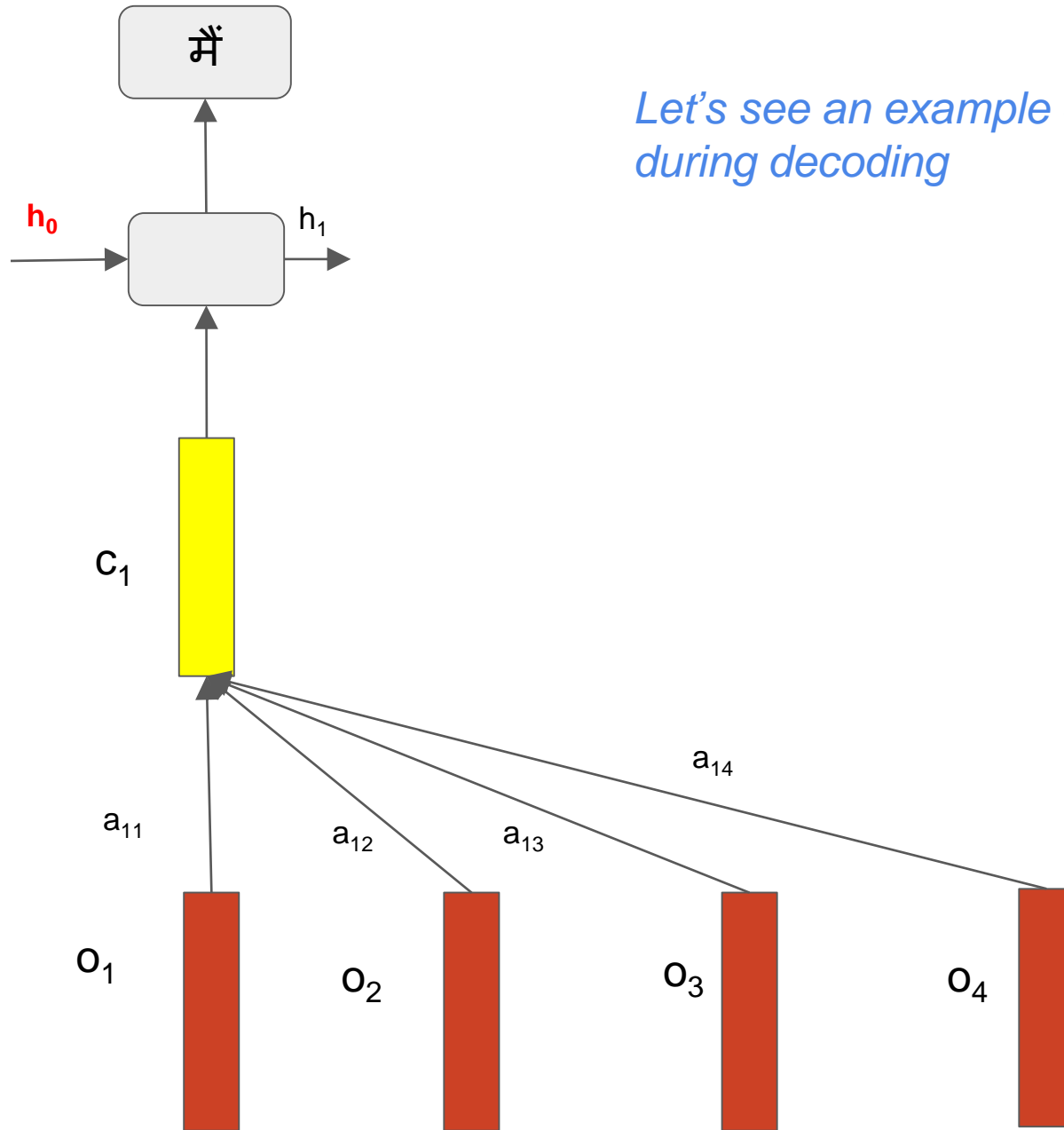
- (1) Not all annotation vectors are equally important for prediction of the next element
- (2) The annotation vector to use next depends on what has been generated so far by the decoder

eg. To generate the 3<sup>rd</sup> target word, the 3<sup>rd</sup> annotation vector (hence 3<sup>rd</sup> source word) is most important

One way to achieve this:

Take a **weighted average of the annotation vectors**, with more weight to annotation vectors which need more **focus or attention**

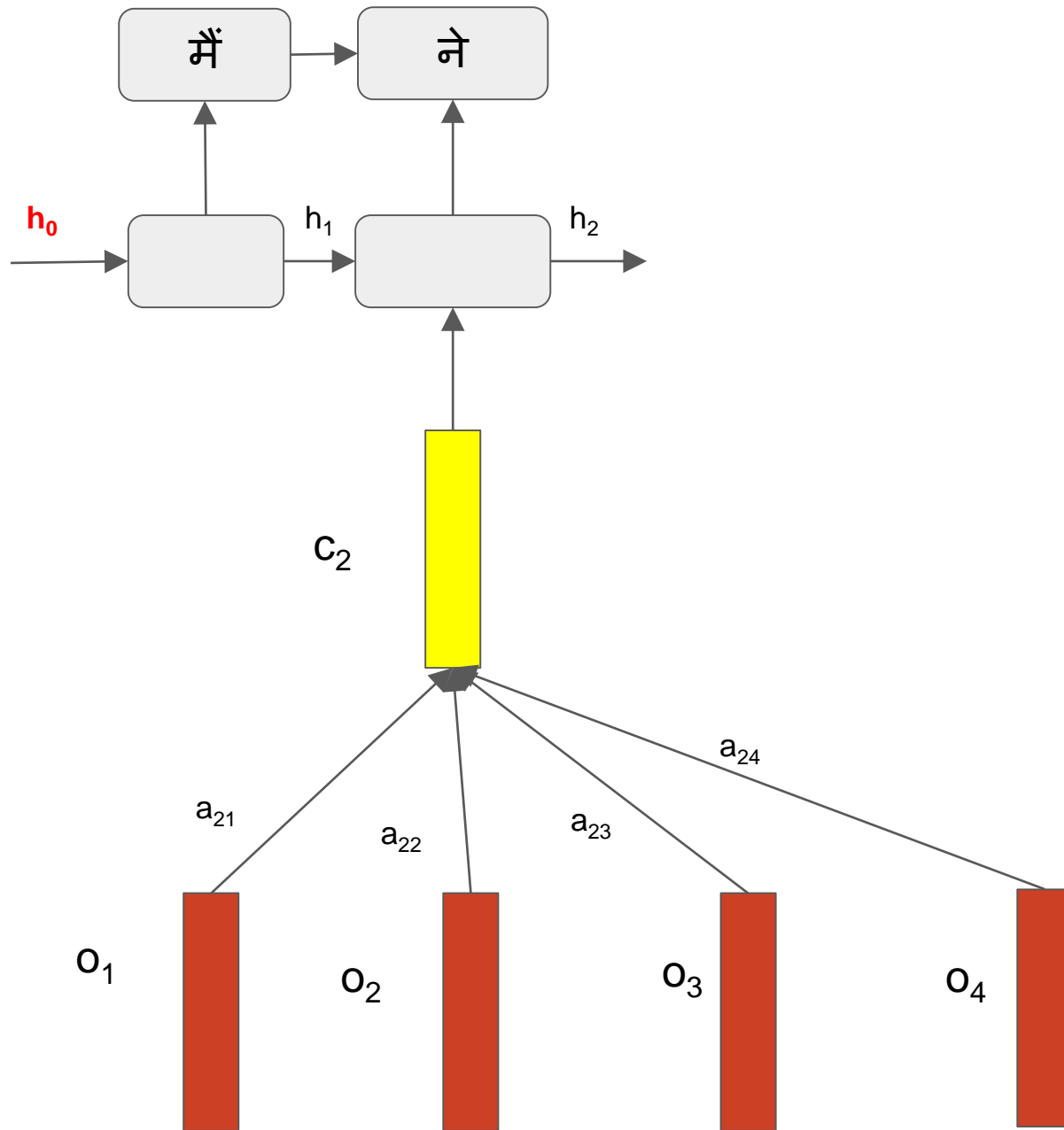
This averaged **context vector** is an input to the decoder

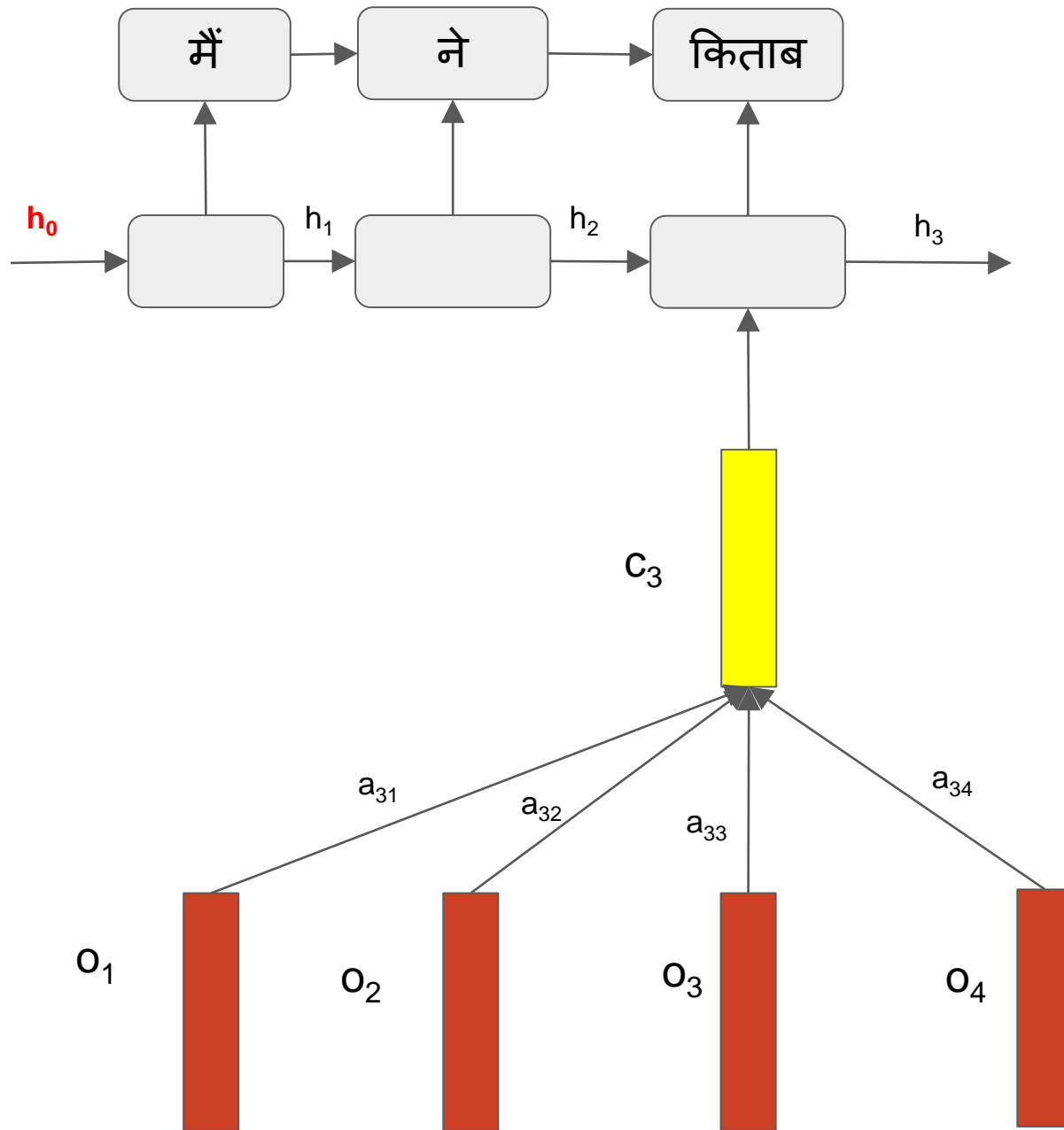


Let's see an example of how the **attention mechanism** works during decoding

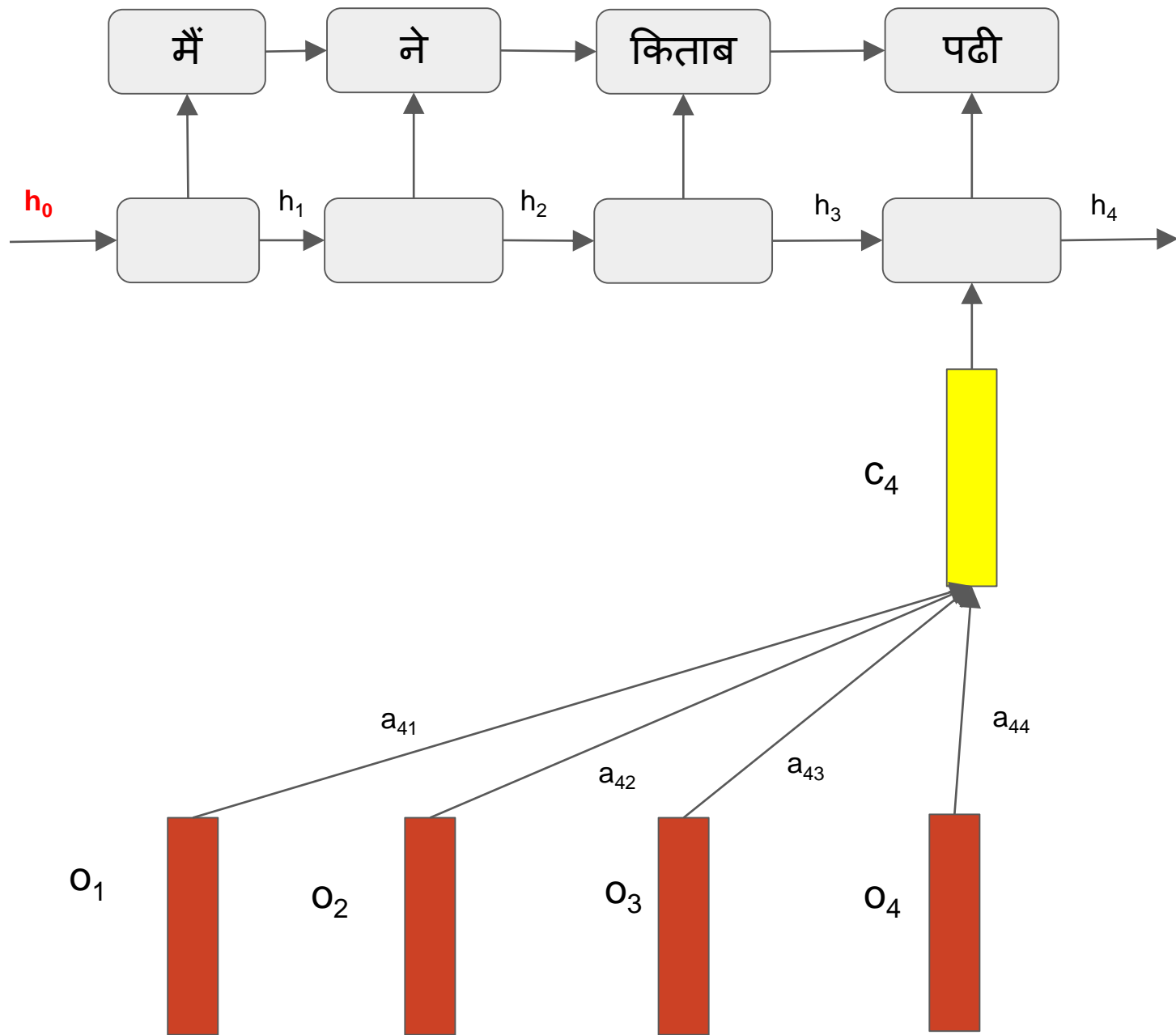
$$c_i = \sum_{j=1}^n a_{ij} o_j$$

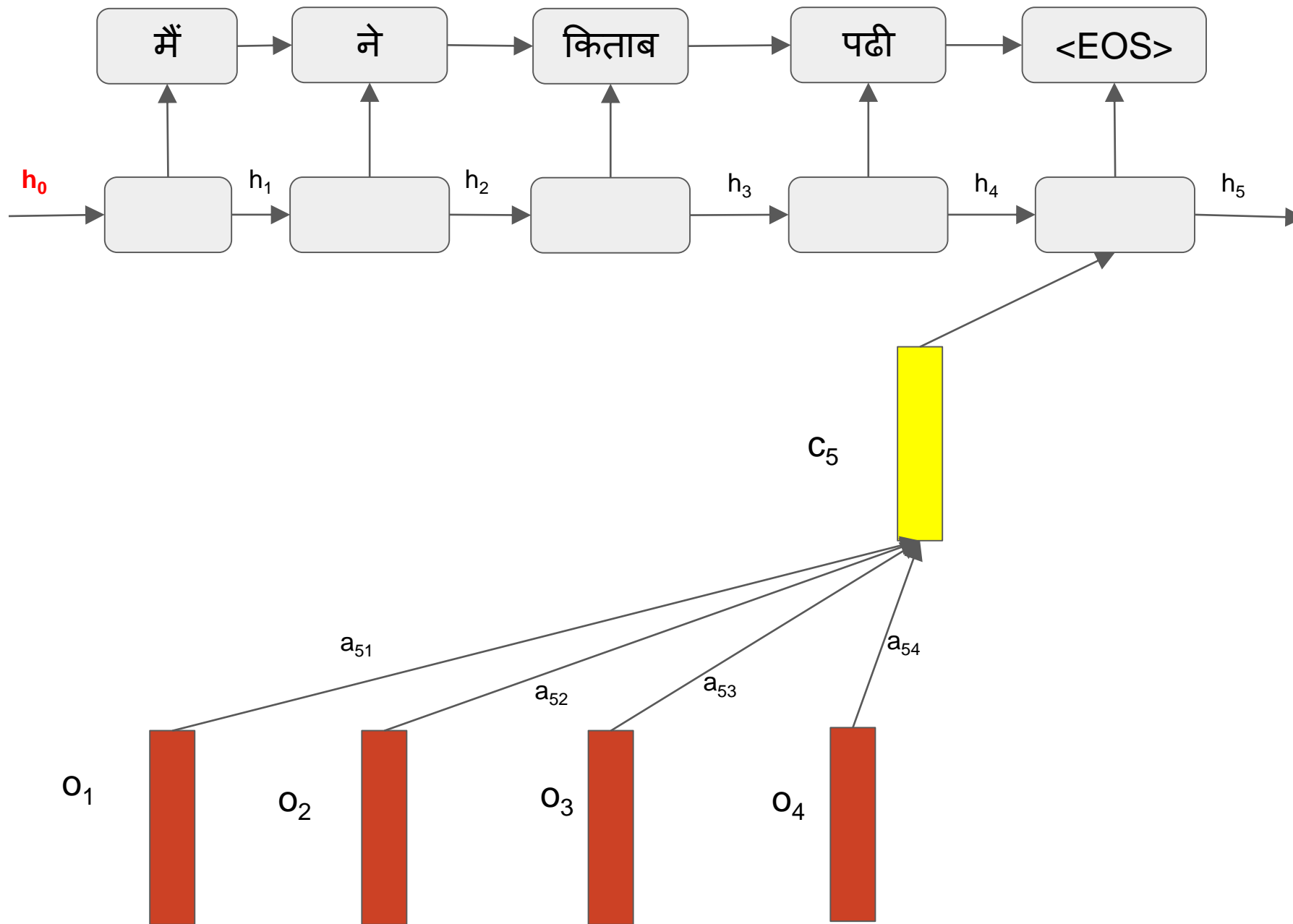
For generation of  $i^{\text{th}}$  output character:  
 $c_i$  : context vector  
 $a_{ij}$  : attention weight for the  $j^{\text{th}}$  annotation vector  
 $o_j$  :  $j^{\text{th}}$  annotation vector











*But we do not know the attention weights?  
How do we find them?*

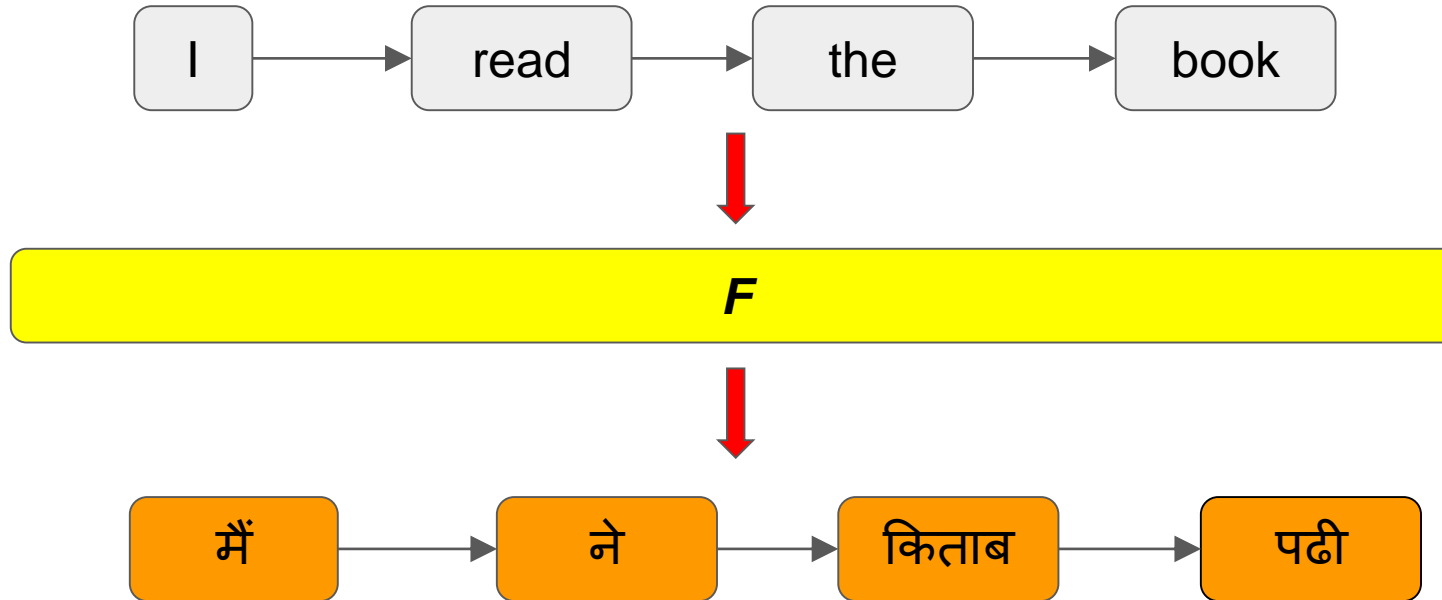
*Let the training data help you decide!!*

**Idea:** Pick the attention weights that maximize the translation accuracy  
(more precisely, decrease training data loss)

- *Note  $\Rightarrow$  no separate language model*
- *Neural MT generates fluent sentences*
- *Quality of word order is better*
- *No combinatorial search required for evaluating different word orders:*
  - *Decoding is very efficient compared to PBSMT*
- *End-to-end training*

We can look at translation as a *sequence to sequence transformation* problem

*Read the entire sequence and predict the output sequence (using function  $F$ )*



- Length of output sequence need not be the same as input sequence
- Prediction at any time step  $t$  has access to the entire input
- A very general framework

## *Sequence to Sequence transformation is a very general framework*

Many other problems can be expressed as sequence to sequence transformation

- *Summarization: Article  $\Rightarrow$  Summary*
- *Question answering: Question  $\Rightarrow$  Answer*
- *Image labelling: Image  $\Rightarrow$  Label*
- *Transliteration: character sequence  $\Rightarrow$  character sequence*

# Machine Translation for Related Languages

# *Related Languages*

*Related by Genealogy*



## *Language Families*

Dravidian, Indo-European, Turkic

*(Jones, Rasmus, Verner, 18<sup>th</sup> & 19<sup>th</sup> centuries, Raymond ed. (2005))*

*Related by Contact*



## *Linguistic Areas*

Indian Subcontinent,  
Standard Average  
European

*(Trubetzkoy, 1923)*

*Related languages may not belong to the same language family!*

# Key Similarities between related languages

भारताच्या स्वातंत्र्यदिनानिमित्त अमेरिकेतील लॉस एन्जल्स शहरात कार्यक्रम आयोजित करण्यात आला

*bhAratAcyA svAta.ntryadinAnimitta ameriketIla lOsA enjalsa shaharAta kAryakrama Ayojita karaNyAta AlA*

Marathi

भारता च्या स्वातंत्र्य दिना निमित्त अमेरिके तील लॉस एन्जल्स शहरा त कार्यक्रम आयोजित करण्यात आला

*bhAratA cyA svAta.ntrya dinA nimitta amerike tIla lOsA enjalsa shaharA ta kAryakrama Ayojita karaNyAta AlA*

Marathi  
segmented

भारत के स्वतंत्रता दिवस के अवसर पर अमरीका के लॉस एन्जल्स शहर में कार्यक्रम आयोजित किया गया

*bhArata ke svata.ntratA divasa ke avasara para amarIkA ke losa enjalsa shahara me.n kAryakrama Ayojita kiyA gayA*

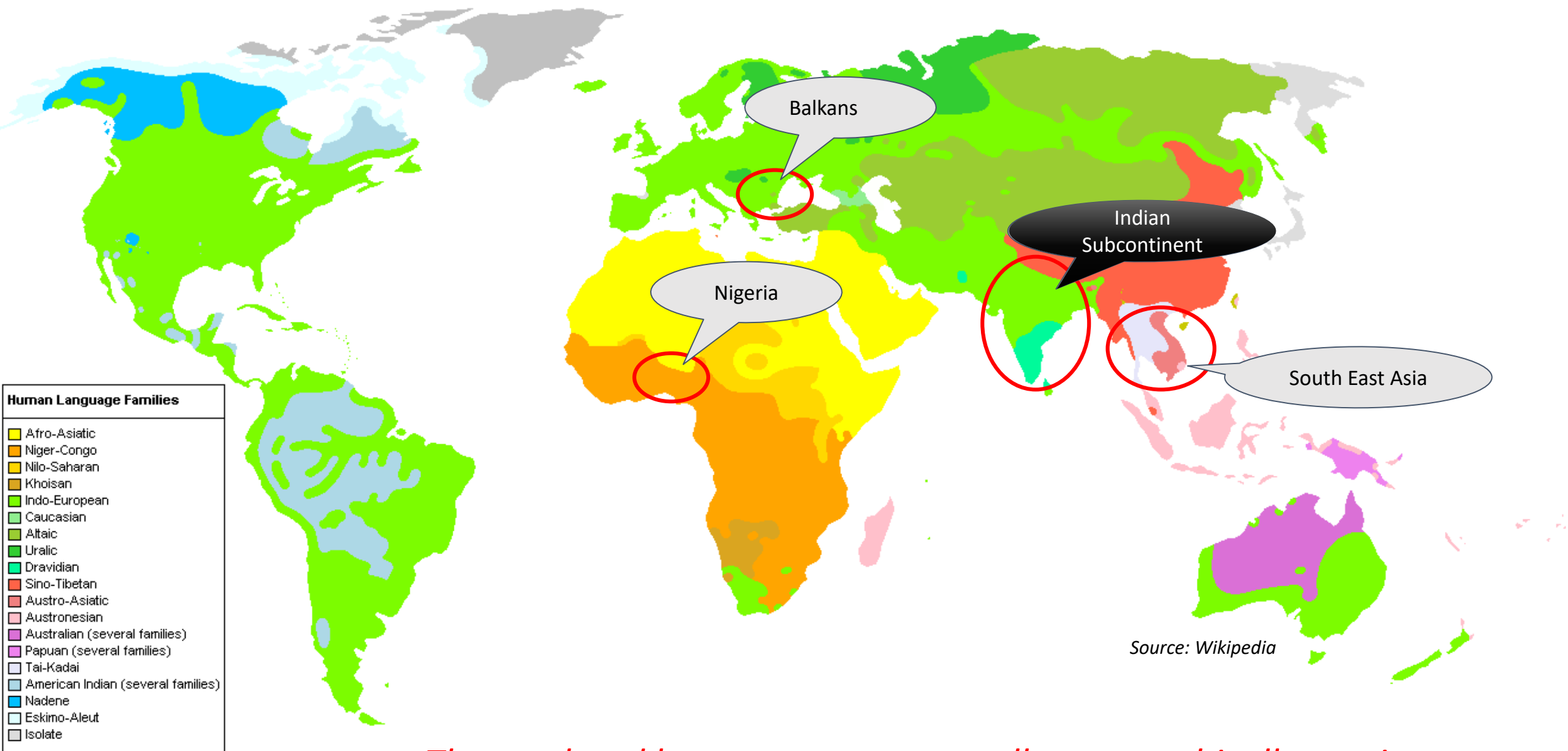
Hindi

**Lexical:** share significant vocabulary (cognates & loanwords)

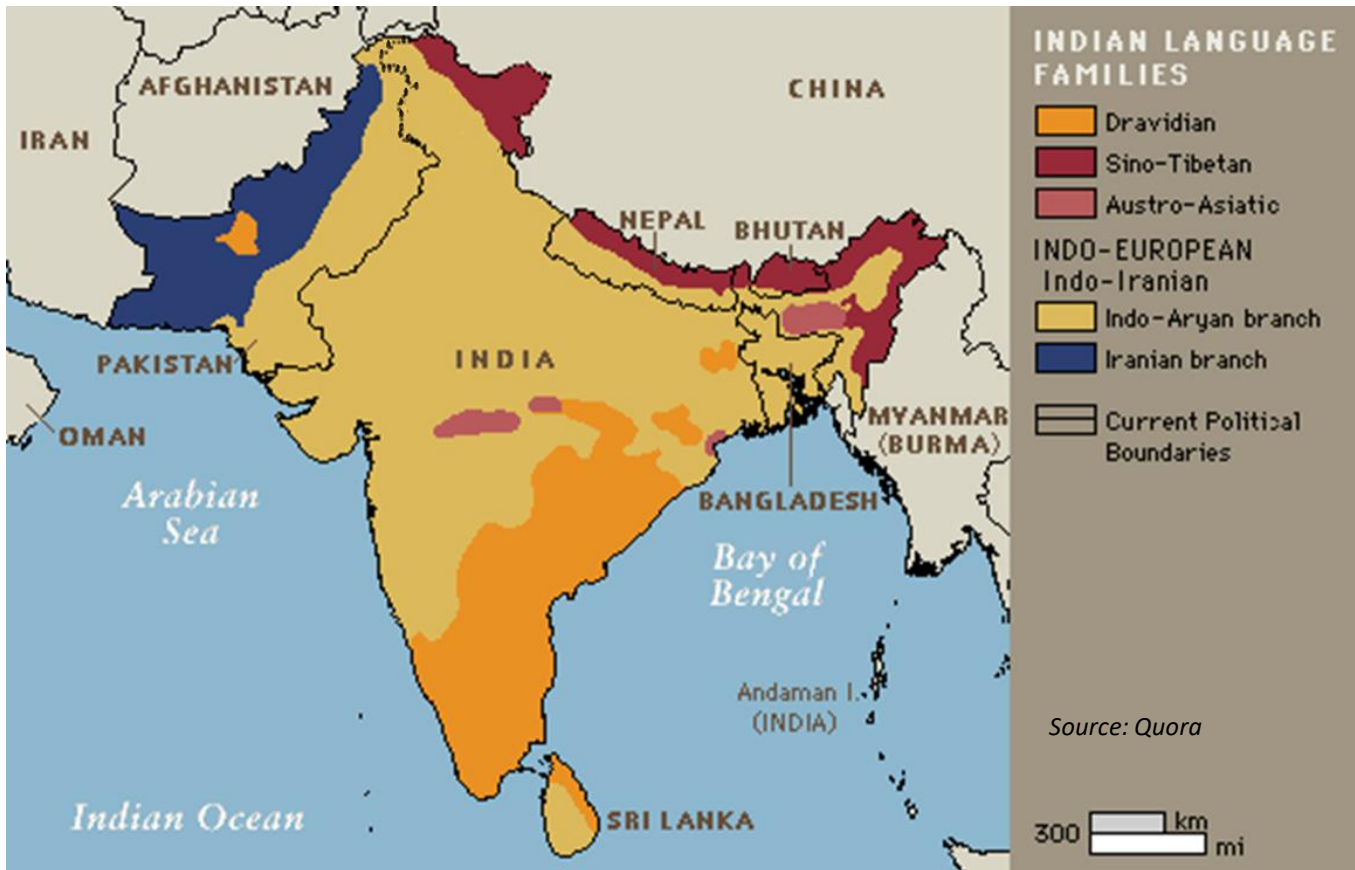
**Morphological:** correspondence between suffixes/post-positions

**Syntactic:** share the same basic word order





*These related languages are generally geographically contiguous*

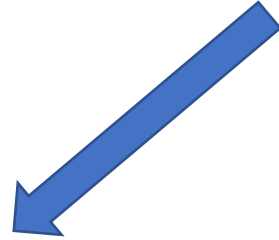


- 5 language families (+ 2 to 3 on the Andaman & Nicobar Islands)
- 22 scheduled languages
- 11 languages with more than 25 million speakers
- Highly multilingual country

*Naturally, lot of communication between such languages  
(government, social, business needs)*

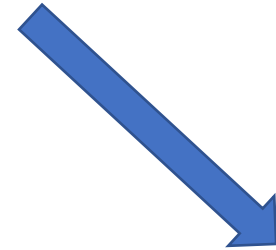


*Most translation requirements also involves related languages*



*Between related languages*

*Hindi-Malayalam  
Marathi-Bengali  
Czech-Slovak*



*Related languages  $\Leftrightarrow$  Link languages*

*Kannada,Gujarati  $\Rightarrow$  English  
English  $\Rightarrow$  Tamil,Telugu*

*We want to be able to handle a large number of such languages*

*e.g. 30+ languages with a speaker population of 1 million + in the Indian subcontinent*

## Aren't “language independent” Statistical/Neural Machine Translation methods sufficient?

- **Implicit assumptions** increase need for:  
(1) Parallel Corpora (2) Linguistic Resources (3) Language specific processing
- ‘*Limited language independence*’ can be achieved between some languages if we can make assumptions that hold across all these languages
- Related languages can serve as a **good level of abstraction** to utilize linguistic regularities:
  - Reduce parallel corpora
  - Reduce linguistic resource requirements
  - Better Generalization

# Utilizing Lexical Similarity

## *Lexically Similar Languages*

(Many words having similar **form** and **meaning**)

- *Cognates*

***a common etymological origin***

<i>roTI (hi)</i>	<i>roTIA (pa)</i>	<i>bread</i>
<i>bhai (hi)</i>	<i>bhAU (mr)</i>	<i>brother</i>

- *Loan Words*

***borrowed without translation***

<i>matsya (sa)</i>	<i>matsyalu (te)</i>	<i>fish</i>
<i>pazha.m (ta)</i>	<i>phala (hi)</i>	<i>fruit</i>

- *Named Entities*

***do not change across languages***

<i>mu.mbal (hi)</i>	<i>mu.mbal (pa)</i>	<i>mu.mbal (pa)</i>
<i>keral (hi)</i>	<i>k.eraLA (ml)</i>	<i>keraL (mr)</i>

- *Fixed Expressions/Idioms*

***MWE with non-compositional semantics***

<i>dAla me.n kuCha kAlA honA</i>	<i>(hi)</i>	<i>Something fishy</i>
<i>dALa mA kAlka kALu hovu</i>	<i>(gu)</i>	

*We will find more matches at the sub-word level*

*Can we use subwords as translation units?*

*Which subword should we use?*

W: राजू, घराबाहेर जाऊ नको .

O: रा जू \_ , \_ घ रा बा हे र \_ जा ऊ \_ न को \_ .

# Simple Units of Text Representation

Transliterate unknown words [Durrani, etal. (2010), Nakov & Tiedemann (2012)]

(a) Primarily used to handle proper nouns

(b) Limited use of lexical similarity

स्वातंत्र्य →  
स्वतंत्रता



*Translation of shared lexically similar words can be seen as kind of transliteration*

Character

[Vilar, etal. (2007), Tiedemann (2009)]

Limited context of character level representation

Limited benefit ....

... just for closely related languages



Character n-gram  $\Rightarrow$  increase in data sparsity

Macedonian - Bulgarian, Hindi-Punjabi, etc.

# Orthographic Syllable *(Kunchukuttan & Bhattacharyya, 2016a)*

(CONSONANT) + VOWEL

Examples: ca, cae, coo, cra, की (kl), प्रे (pre)  
अभिमान → अ भि मा न

## Pseudo-Syllable

True Syllable ⇒ Onset, Nucleus and Coda

Orthographic Syllable ⇒ Onset, Nucleus

- Generalization of **akshara**, the fundamental organizing principle of Indian scripts
- Linguistically motivated, **variable length unit**
- *Number of syllables in a language is finite*



# Byte Pair Encoded (BPE) Unit

*(Kunchukuttan & Bhattacharyya, 2017a; Nguyen and Chang, 2017)*

- *There may be frequent subsequences in text other than syllables*
- *These subsequences may **not be valid linguistic units***
- *But they represent **statistically important patterns** in text*

***How do we identify such frequent patterns?***

Byte Pair Encoding (Sennrich et al, 2016), wordpieces (Wu et al, 2016),  
Huffman encoding based units (Chitnis & DeNero, 2015)

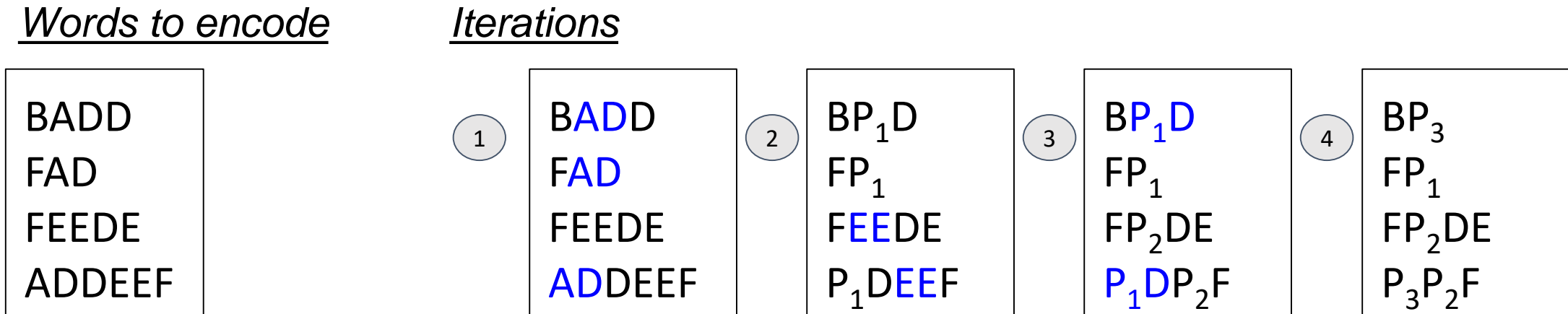
# Byte Pair Encoded (BPE) Unit

*Byte Pair Encoding is a compression technique (Gage, 1994)*

Number of BPE merge operations=3

Vocab: A B C D E F

$P_1=AD$   $P_2=EE$   $P_3=P_1D$



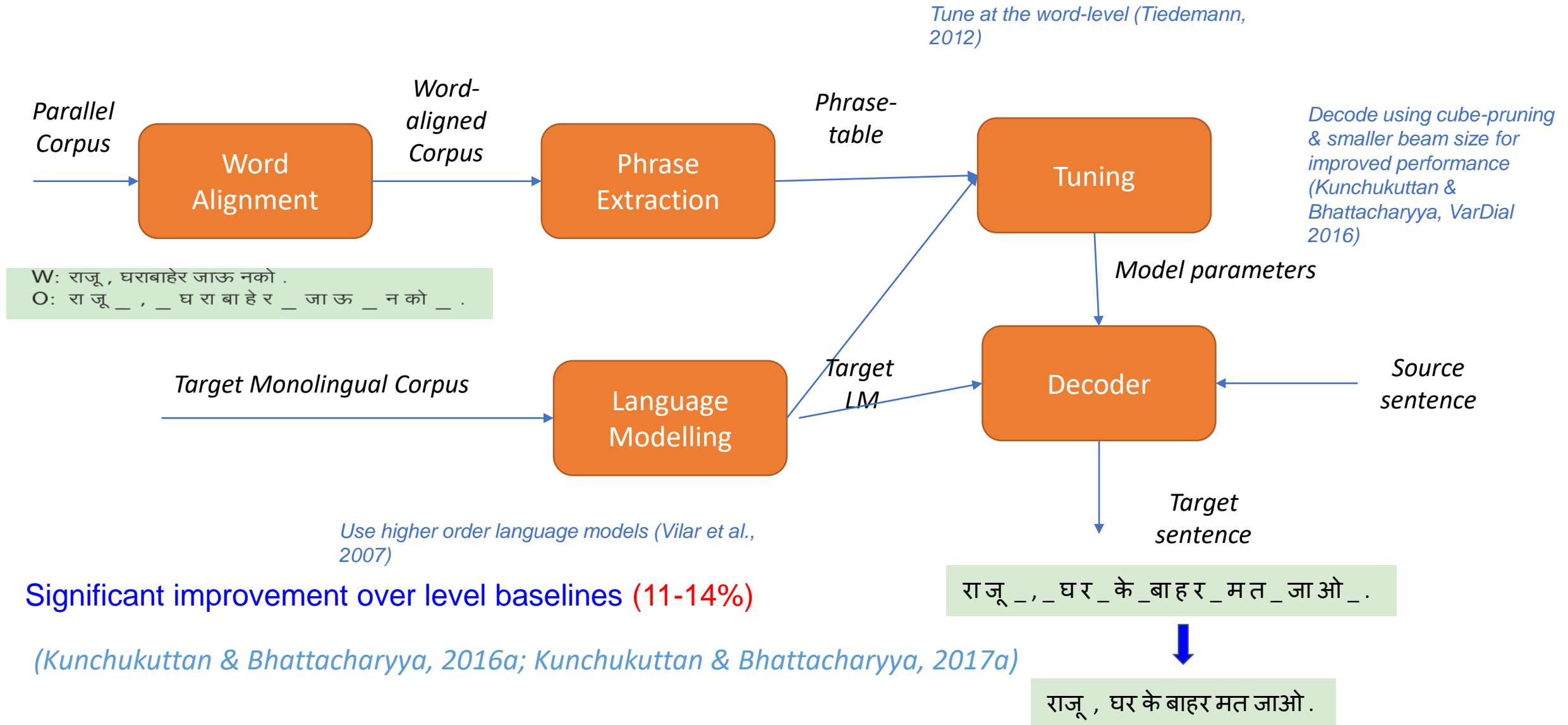
## Data-dependent segmentation

- Inspired from compression theory
- MDL Principle (*Rissanen, 1978*)  $\Rightarrow$  Select segmentation which maximizes data likelihood

# Example of various translation units

Basic Unit	Symbol	Example	Transliteration
Word	W	घरासमोरचा	gharAsamoracA
Morph Segment	M	घरा समोर चा	gharA samora cA
Orthographic Syllable	O	घ रा स मो र चा	gha rA sa mo racA
Character unigram	C	घ र ा स म ो र च ा	gha r A sa m o ra c A
<i>something that is in front of home: ghara=home, samora=front, cA=of</i>			
Various translation units for a Marathi word			

# Adapting SMT for subword-level translation



# NMT with related languages on source side

(Nguyen and Chang, 2017)

We want **Marathi** → **English** translation → but little parallel corpus is available  
We have lot of **Hindi** → **English** parallel corpus

It is cold in Pune	पुण्यात थंड आहे
My home is near the market	माझा घर बाजाराजवळ आहे

I am going home	मैं घर जा रहा हूँ
It rained last week	पिछले हफ्ते बारिश हुई

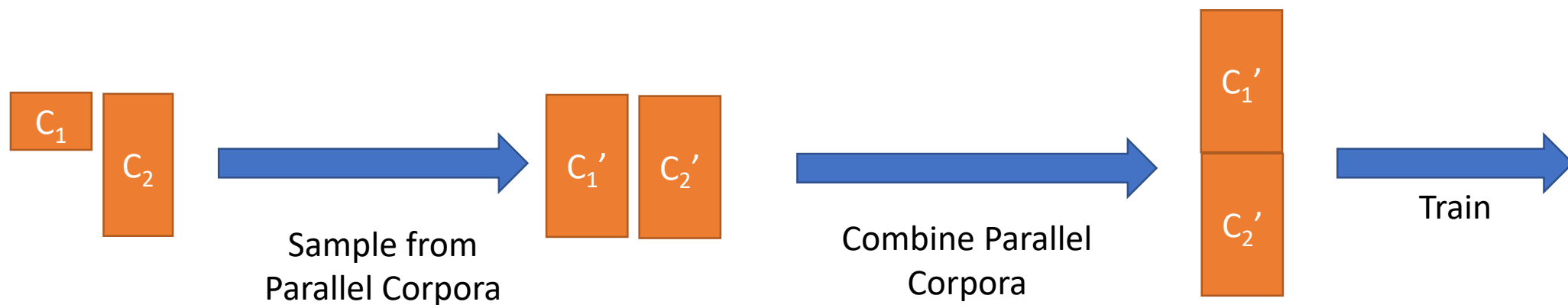


Concat Corpora

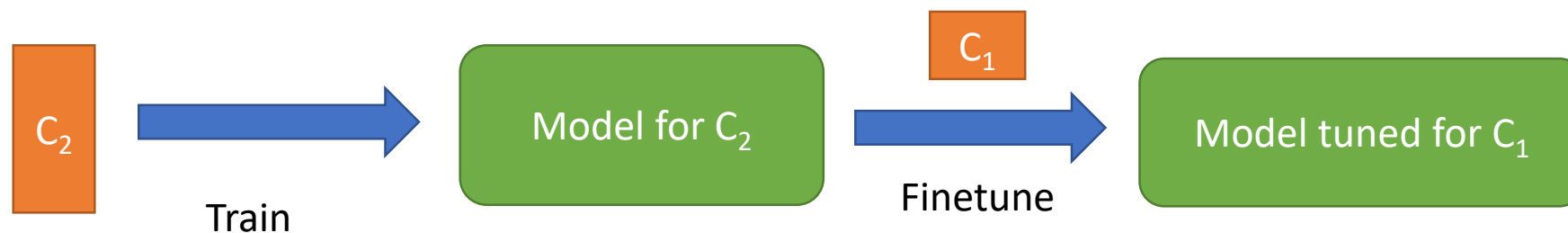
It is cold in Pune	पुण्यात थंड आहे
My home is near the market	माझा घर बाजाराजवळ आहे
I am going home	मैं घर जा रहा हूँ
It rained last week	पिछले हफ्ते बारिश हुई

# Training Multilingual NMT systems

## Method 1



## Method 2



# What if the related languages use different scripts?

(Nguyen and Chang, 2017)

We want **Gujarati** → **English** translation → but little parallel corpus is available  
We have lot of **Marathi** → **English** parallel corpus

I am going home	હુ ઘરે જવ છું
It rained last week	છેલ્લા આઠવડિયા મા વર્સાદ પાડ્યો

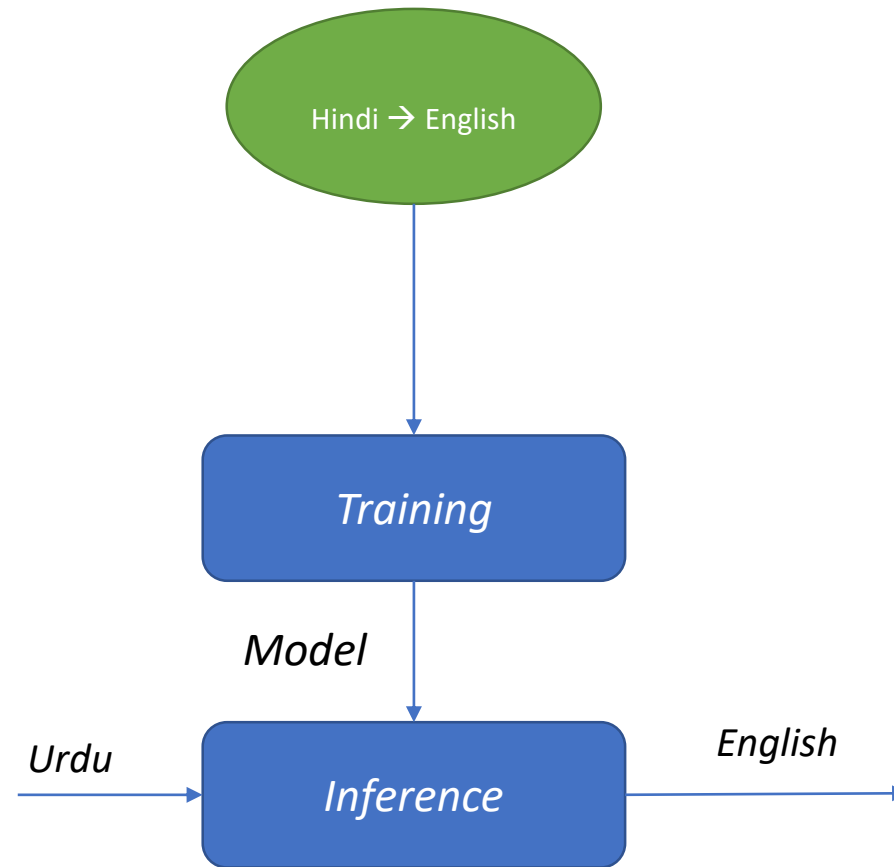
It is cold in Pune	પુણ્યાત થંડ આહે
My home is near the market	માઝા ઘર બાજારાજવલ આહે

Convert Script

Concat Corpora

I am going home	હુ ઘરે જવ છું
It rained last week	છેલ્લા આઠવડિયા મા વર્સાદ પાડ્યો
It is cold in Pune	પુણ્યાત થંડ આહે
My home is near the market	માઝા ઘર બાજારાજવલ આહે

# Zeroshot Translation





*How do we support multiple target languages with a single decoder?*

*A simple trick!*

*Append input with special token indicating the target language*

For English-Hindi Translation

Original Input: *France and Croatia will play the final on Sunday*

Modified Input: *France and Croatia will play the final on Sunday* **<hin>**

*Still an open problem*

# Utilizing Syntactic Similarity

*(Kunchukuttan et al., 2014)*

*Phrase based MT is not good at learning word ordering*

*Solution: Let's help PB-SMT with some preprocessing of the input*

*Change order of words in input sentence to match order of the words in the target language*

Let's take an example

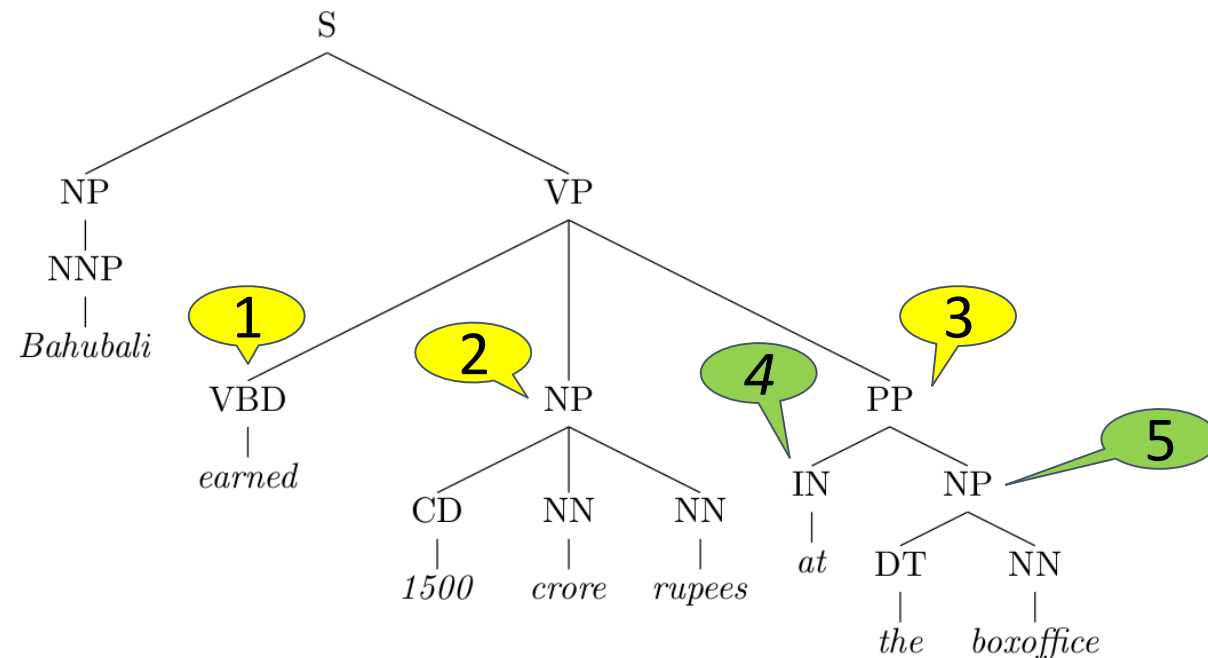
*Bahubali earned more than 1500 crore rupee sat the boxoffice*

Parse the sentence to understand its syntactic structure

Apply rules to transform the tree

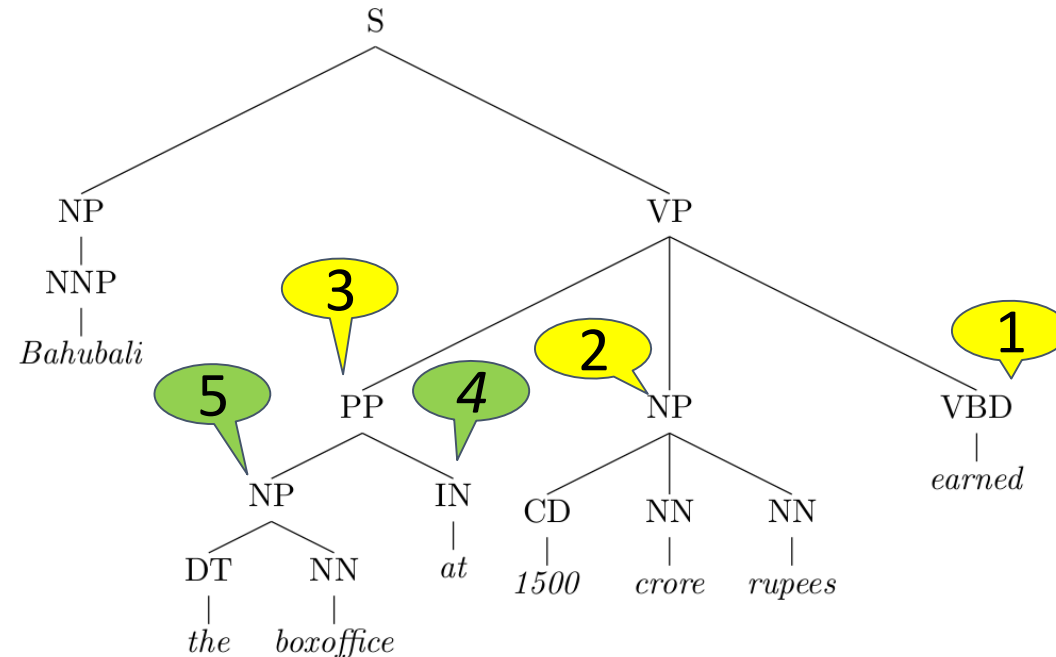
VP → VBD NP PP ⇒ VP → PP NP VBD

PP → IN NP ⇒ PP → NP IN



The new input to the machine translation system is:  
*Bahubali the boxoffice at 1500 crore rupees earned*

Now we can translate with little reordering:  
*बाहुबली ने बॉक्सओफिस पर 1500 करोड रुपए कमाए*



## Can we reuse English-Hindi rules for English-Indian languages?

*All Indian languages have the same basic word order*

	Indo-Aryan						Dravidian		
	pan	hin	guj	ben	mar	kok	tel	tam	mal
Baseline	15.83	21.98	15.80	12.95	10.59	11.07	7.70	6.53	3.91
Generic	17.06	23.70	16.49	13.61	11.05	11.76	7.84	6.82	<b>4.05</b>
Hindi-tuned	<b>17.96</b>	<b>24.45</b>	<b>17.38</b>	<b>13.99</b>	<b>11.77</b>	<b>12.37</b>	<b>8.16</b>	<b>7.08</b>	4.02

*(Kunchukuttan et al., 2014)*

### Generic reordering (Ramanathan et al 2008)

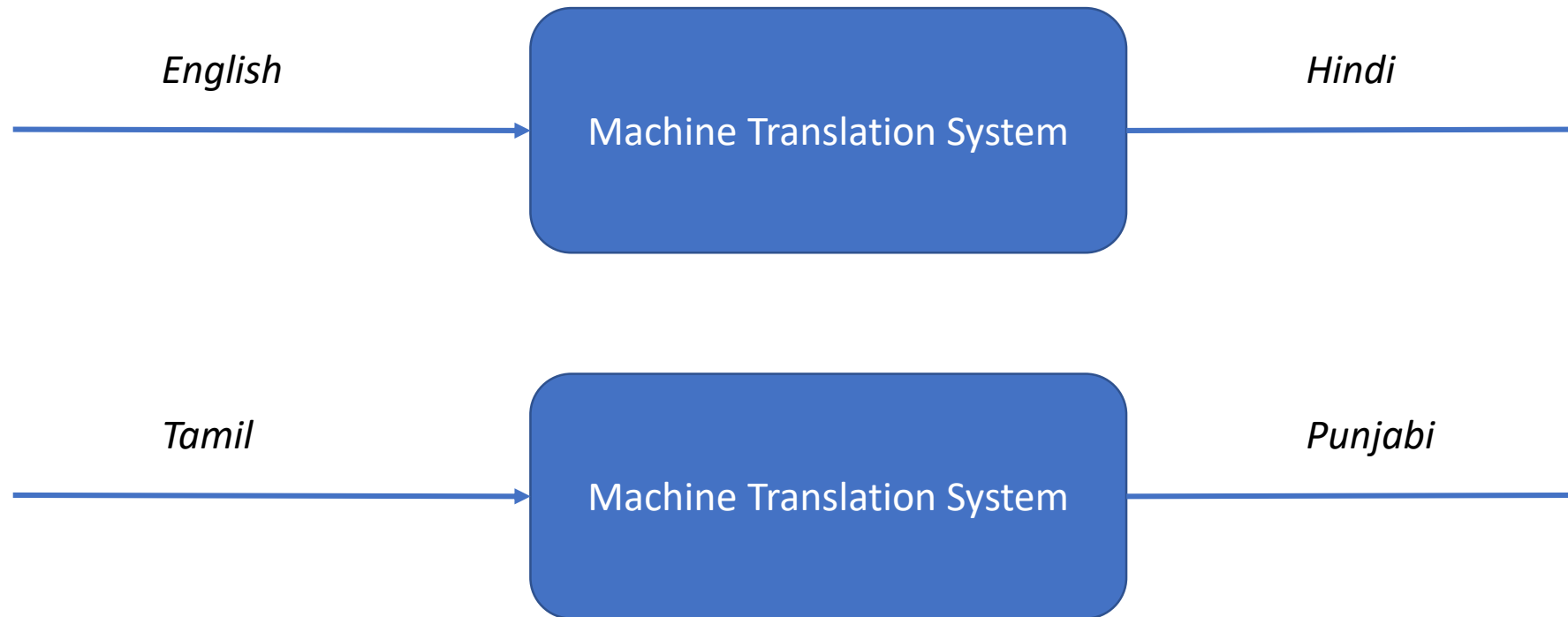
Basic reordering transformation for English → Indian language translation

### Hindi-tuned reordering (Patel et al 2013)

Improvement over the basic rules by analyzing English → Hindi translation output

# Multilingual Learning

*Broad Goal: Build NLP Applications that can work on different languages*



## *Monolingual Applications*

Document Classification  
Sentiment Analysis  
Entity Extraction  
Relation Extraction  
Information Retrieval  
Question Answering  
Conversational Systems

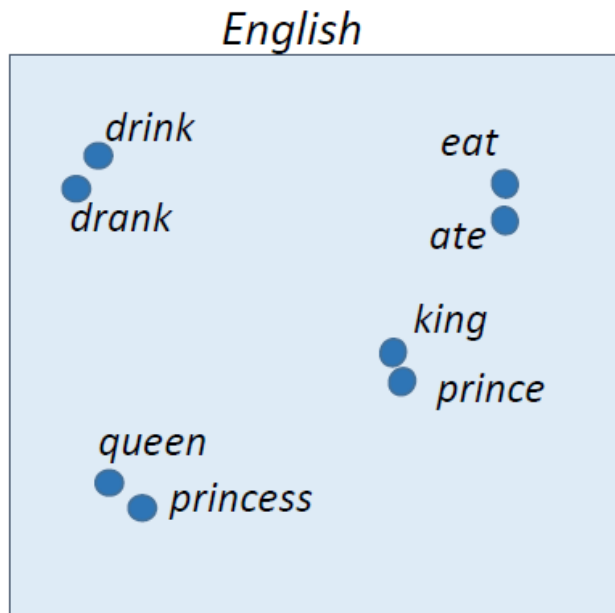
## *Cross-lingual Applications*

Translation  
Transliteration  
**Cross-lingual Applications**  
Information Retrieval  
Question Answering  
Conversation Systems

Code-Mixing  
Creole/Pidgin languages  
Language Evolution  
Comparative Linguistics

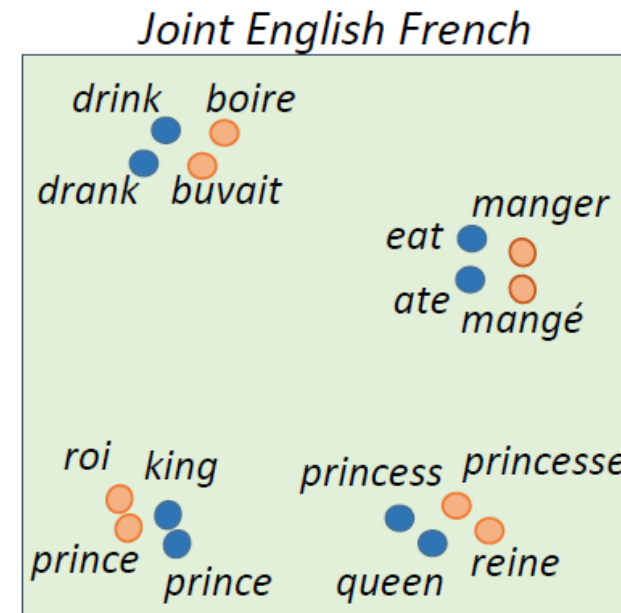
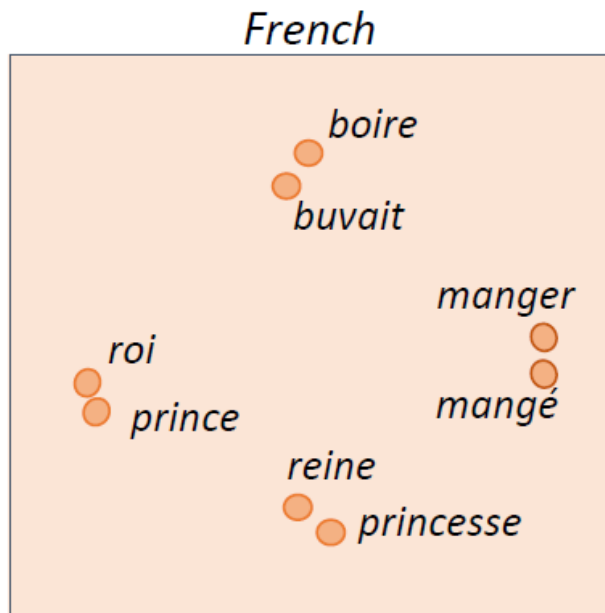
## *Mixed Language Applications*

# Cross Lingual Embeddings



## Monolingual Word Representations

(capture syntactic and semantic similarities between words)



## Multilingual Word Representations

(capture syntactic and semantic similarities between words both within and across languages)

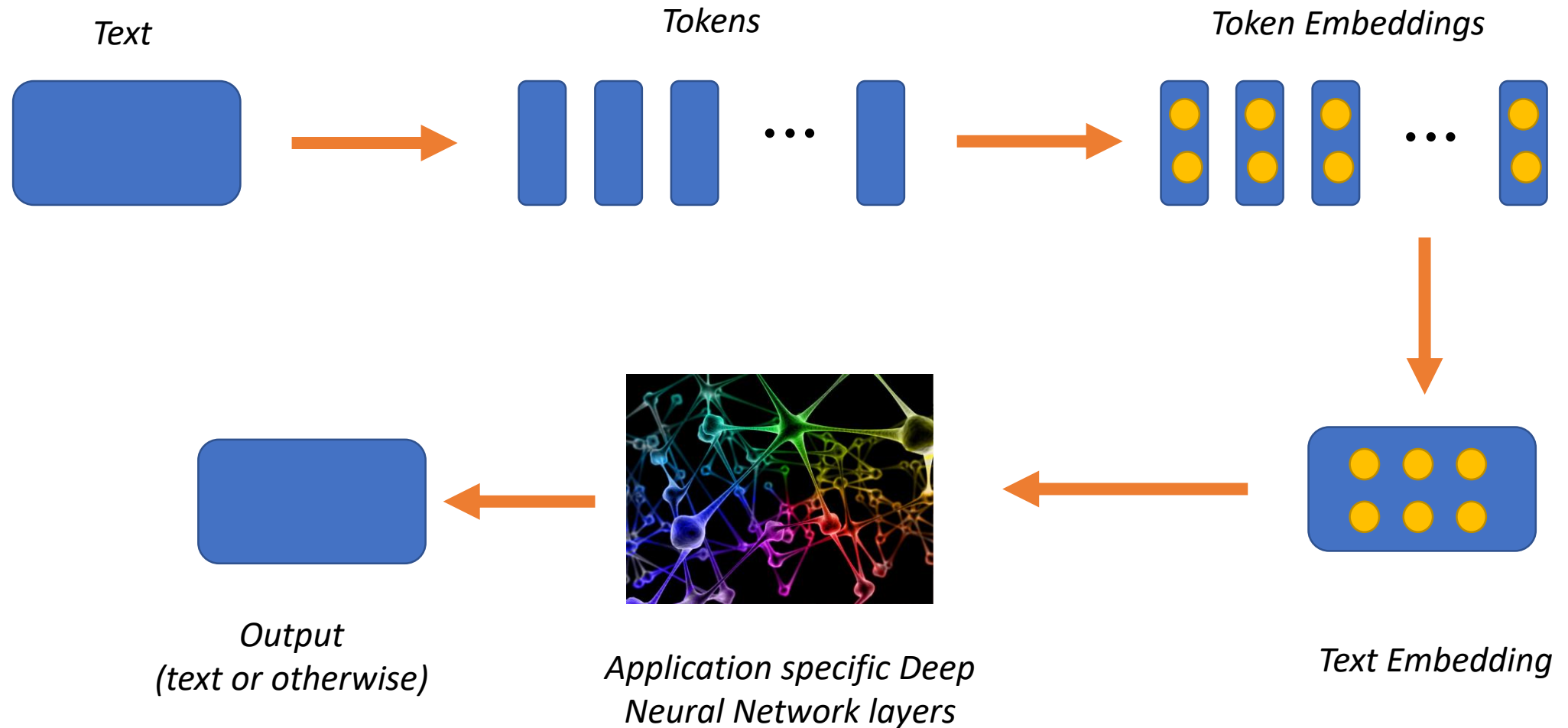
(Source: Khapra and Chandar, 2016)

$$\text{embed}(y) = f(\text{embed}(x))$$

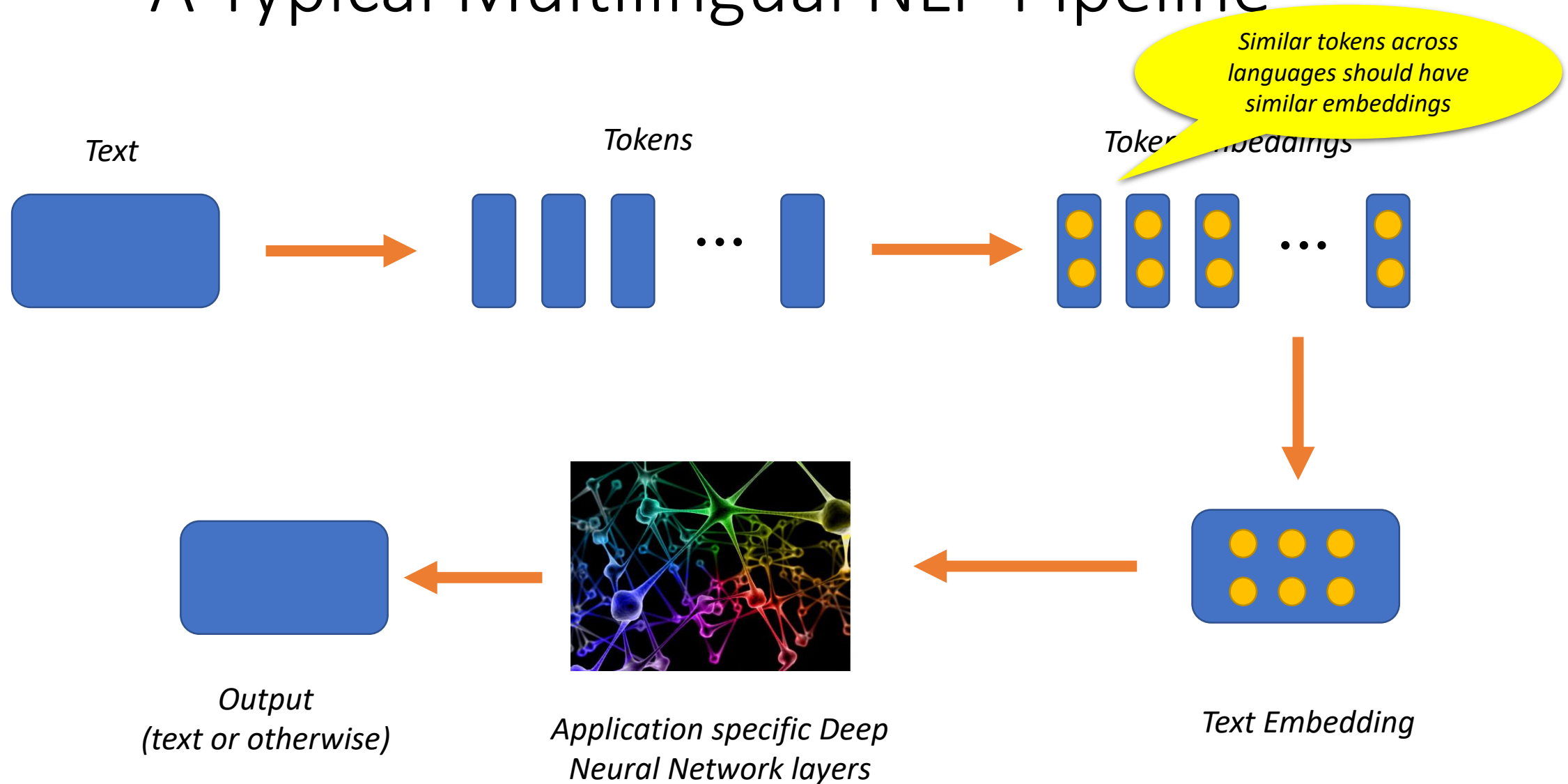
$x, y$  are source and target words  
 $\text{embed}(w)$ : embedding for word  $w$



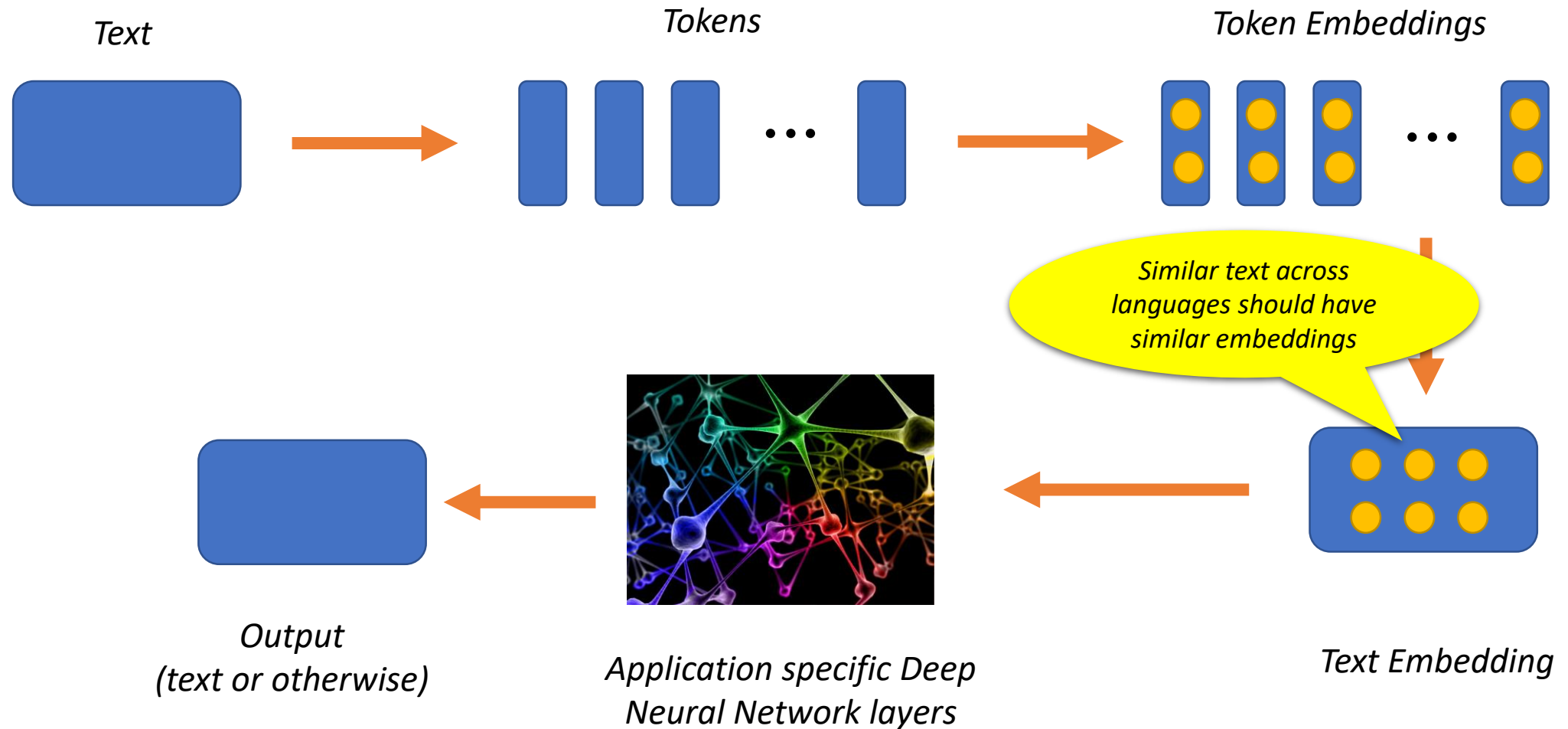
# A Typical Multilingual NLP Pipeline



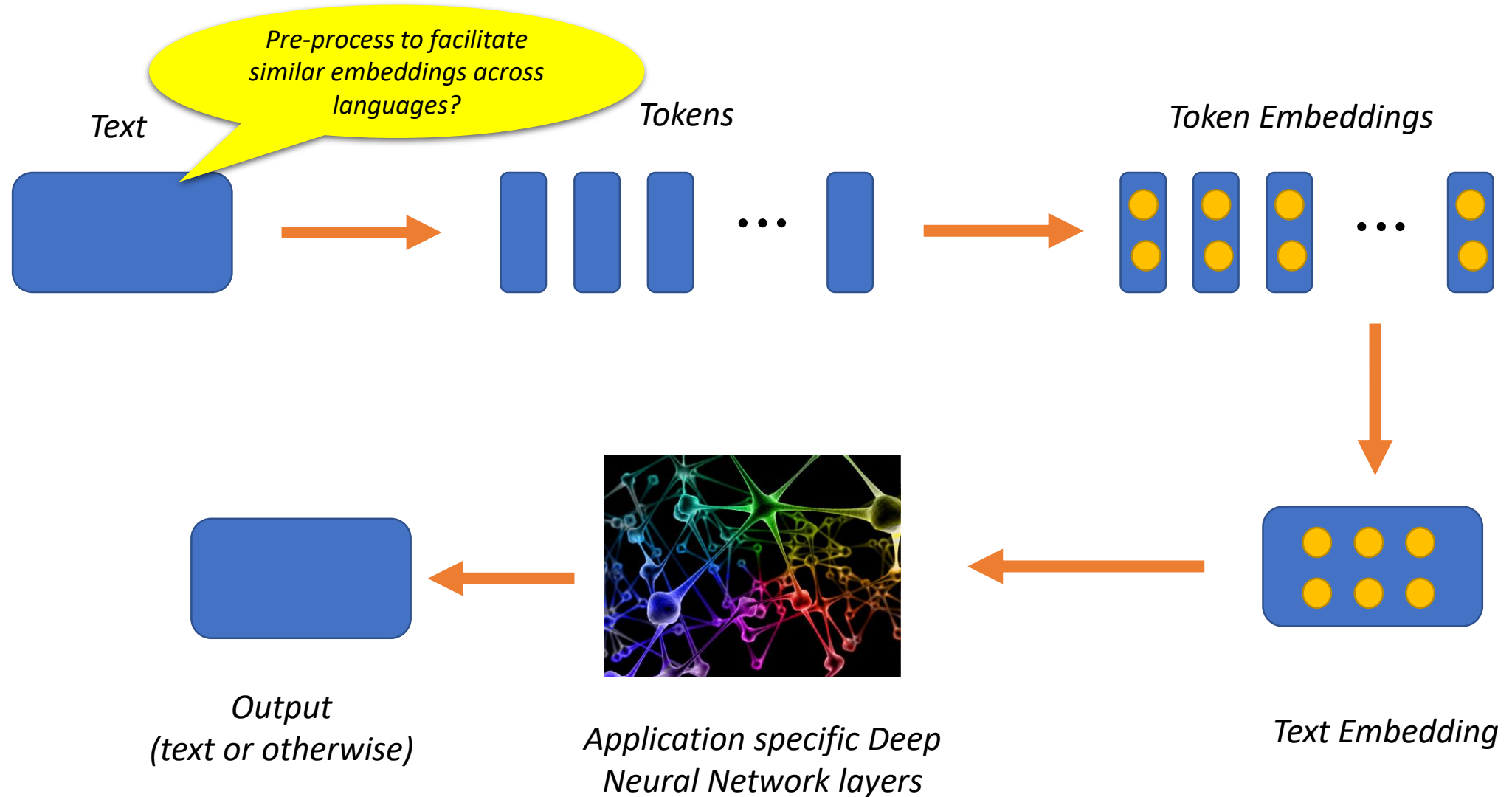
# A Typical Multilingual NLP Pipeline



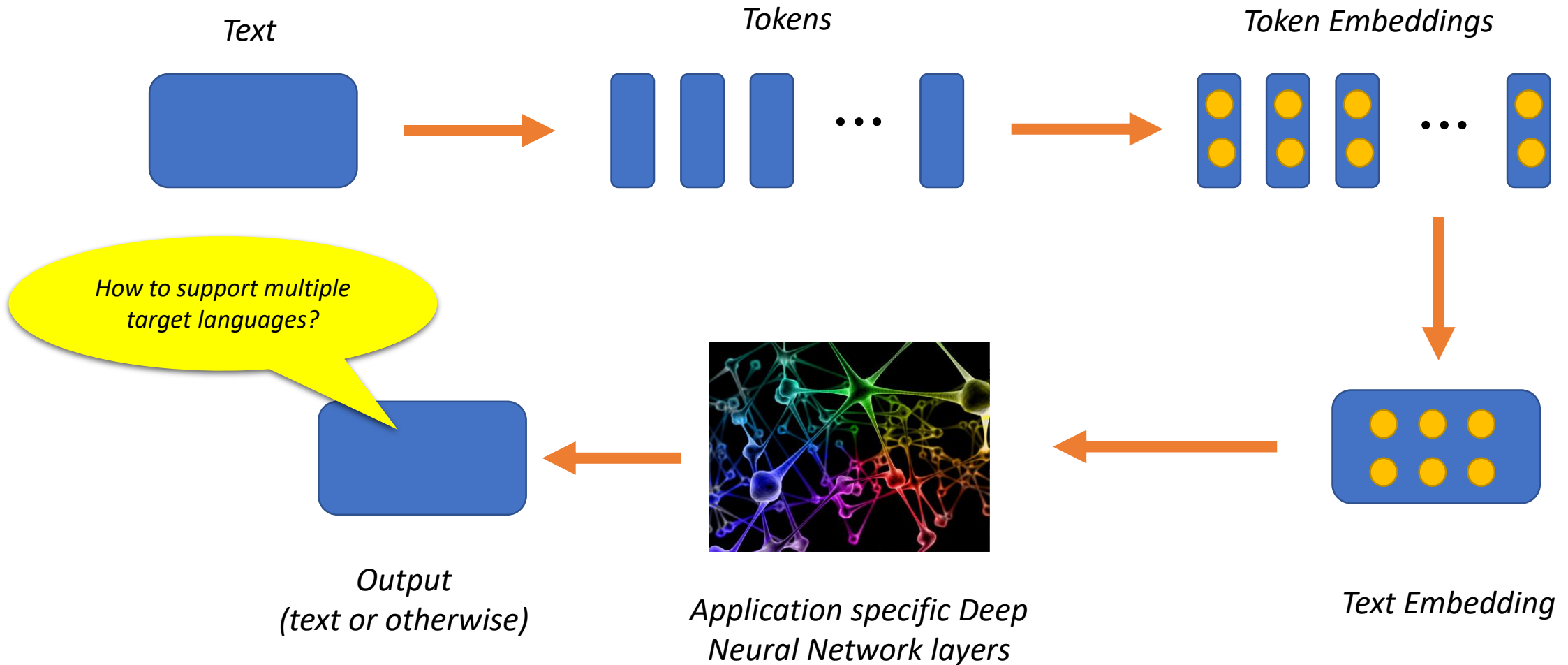
# A Typical Multilingual NLP Pipeline



# A Typical Multilingual NLP Pipeline



# A Typical Multilingual NLP Pipeline



# More Reading Material

This was a small introduction, you can find more elaborate presentations and further references to explore below:

## SMT Tutorials

- *Machine Learning for Machine Translation (An Introduction to Statistical Machine Translation)*. **Tutorial at ICON 2013** with Prof. Pushpak Bhattacharyya, Piyush Dungarwal and Shubham Gautam. [\[slides\]](#) [\[handouts\]](#)
- *Machine Translation: Basics and Phrase-based SMT*. **Talk at the Ninth IIIT-H Advanced Summer School on NLP (IASNLP 2018), IIIT Hyderabad**. [\[pdf\]](#) [\[pptx\]](#)

## NMT Tutorial

### Machine Translation for Related Languages

- *Statistical Machine Translation between related languages*. Tutorial at NAACL 2016 with Prof. Pushpak Bhattacharyya and Mitesh Khapra. [\[abstract\]](#) [\[slides\]](#)
- *Machine Translation for related languages*. Tech Talk at AXLE 2018 (Microsoft Academic Accelerator). [\[pdf\]](#) [\[pptx\]](#)
- *Translation and Transliteration between related languages*. Tutorial at ICON 2015 with Mitesh Khapra. [\[abstract\]](#) [\[slides\]](#) [\[handouts\]](#)

### Multilingual Training

- *Multilingual Learning*. Invited Talk at IIIT Hyderabad Machine Learning Summer School (Advances in Modern AI) 2018. [\[slides\]](#)

# Thank you!

[anoop.kunchukuttan@gmail.com](mailto:anoop.kunchukuttan@gmail.com)