

Any 6

Speaker: Purushottam Bhattacharyya, IIT Bombay

Topic: Machine Translation

Morphemes → The smallest part in a word having a meaning

Eg Inconceivable

In conceive able

↓ ↓ ↓
Morphemes

(1 word → 1 meaning)

Isolating

Languages

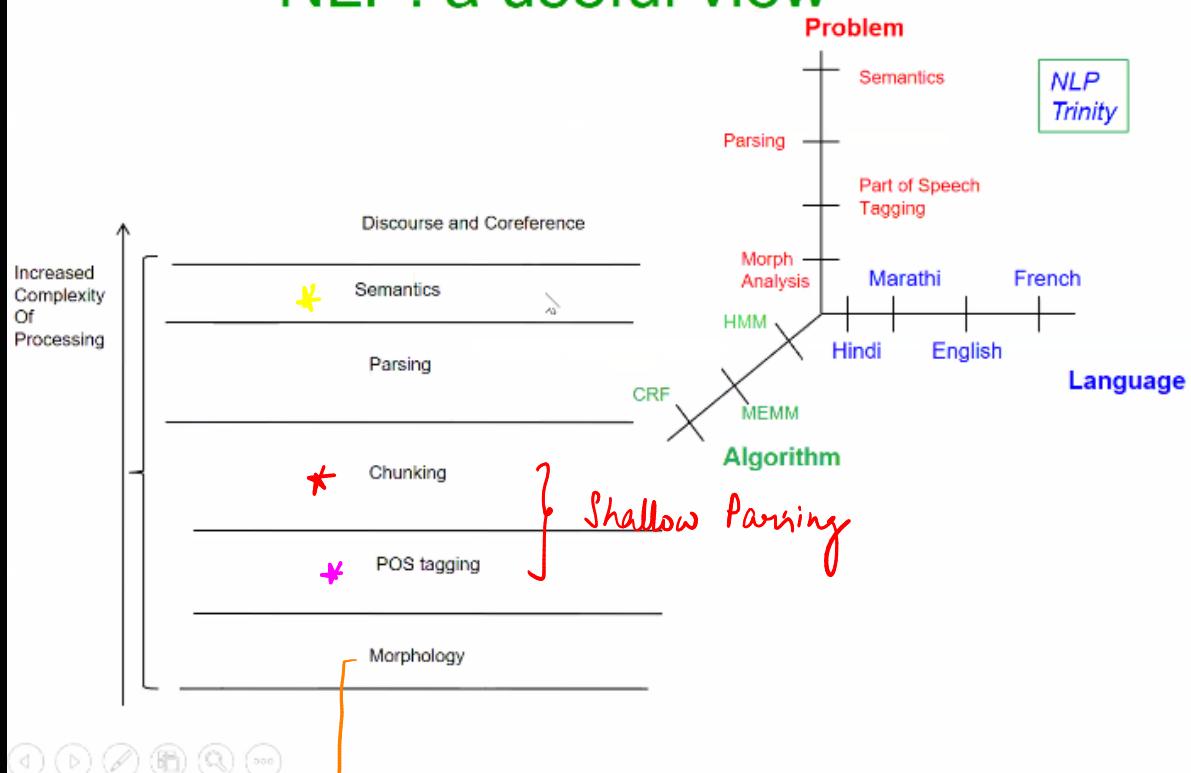
Analytic

(glue many words together)

Agglutinative

Synthetic
(1 word many meanings)

NLP: a useful view



→ Breaking a word into its parts. → One category

↓
Agglutination

(Stitching many words together)

Extreme Agglutination → Manipuri etc

* Part of Speech (POS Tagging) → Classification Problem.

For each part in a word or sentence → category is assigned

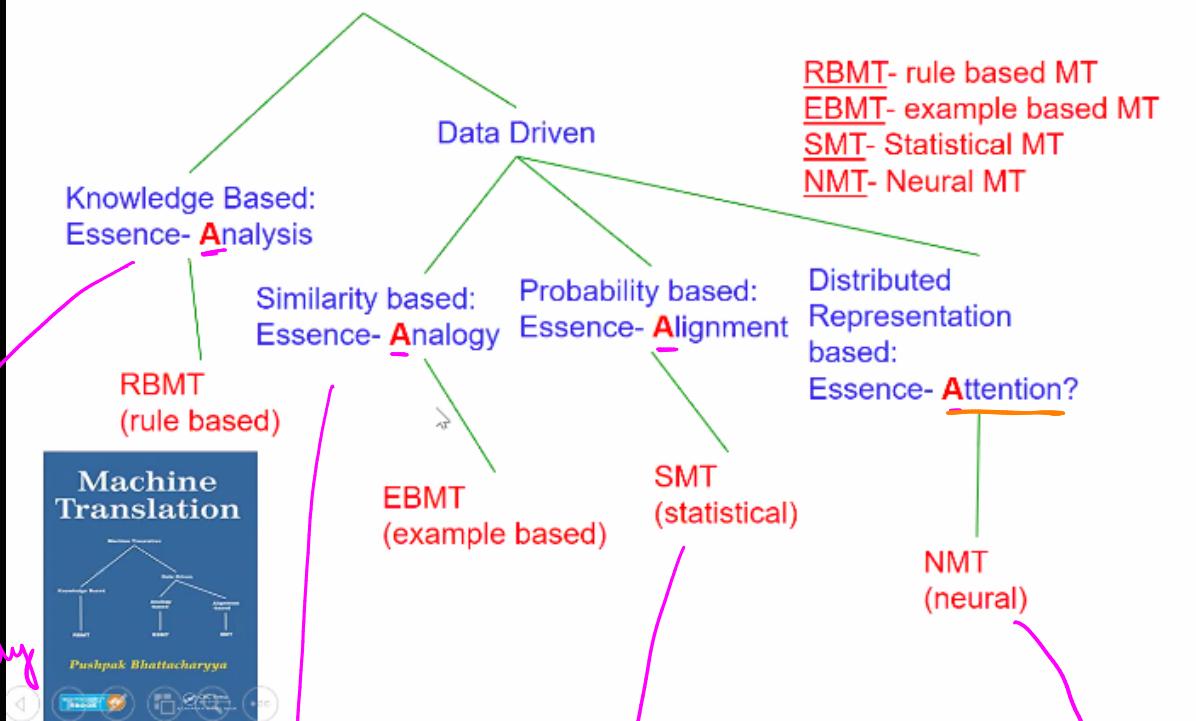
* Chunking (in case of small phrases)

* Semantics → Semantic Role Identification (Place, Time, etc)
Sense of word

* Rule Based MT ~ Still Important.

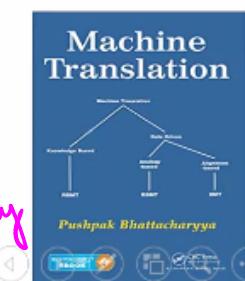
Machine Translation: Translating from one language to another by computer

Machine Translation



Source Sentence is analysed & converted to a form to directly generate target sentence

~ 1960's/1970's NLP Stack is born
 Analysis needs to be good



Can reuse

(translation memory)

Analogy

IBM → 1990's

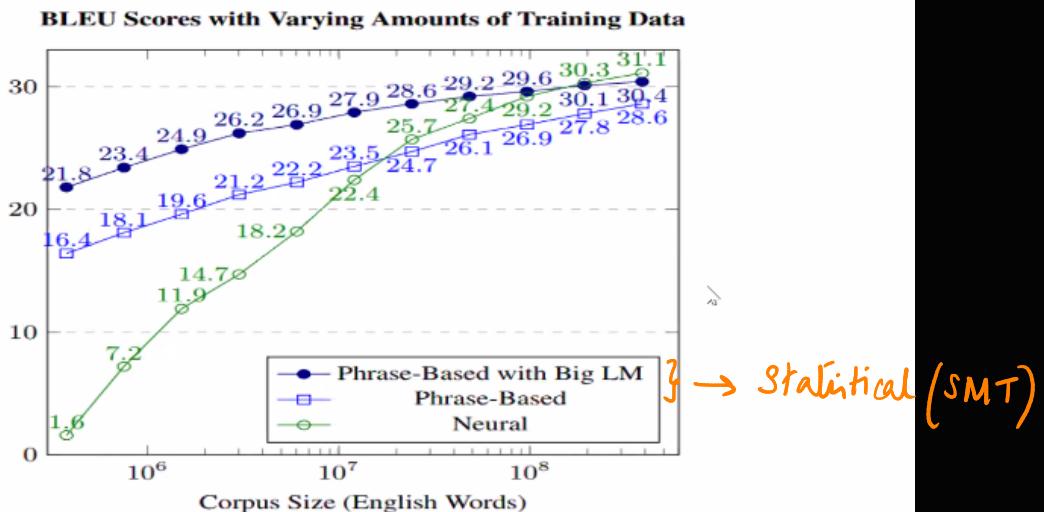
Expectation Maximization

{ Makes alignment b/w words in a pair of || sentences }

Source & Target Sentences

Encoder Decoder
{ Transformer }

Today's Ruling Paradigm: NMT which is data intensive



Philipp Koehn and Rebecca Knowles. 2017. *Six Challenges for Neural Machine Translation*. In Proceedings of the First Workshop on Neural Machine Translation, pages 28–39.

As data ↑, Neural N/W method (Transformers) beats all other methods
But!! * Neural → Lack of interpretability/explainability SMT etc.

Challenge of MT: Language Divergence

- Languages have different ways of expressing meaning
 - * –Lexico-Semantic Divergence
 - * –Structural Divergence

Our work on English-IL Language Divergence with illustrations from Hindi (Dave, Parikh, Bhattacharya, Journal of MT, 2002)

Understanding Divergences:

Different ways of expressing meaning

Indo-European
Tibeto-Burman

English:
Hindi: *yana kambal bahut naram hai*
Bangla: *ei kambal ti khub naram <null>*
Marathi: *haa kambal khup naram aahe*
Manipuri: *kampor blanket asi mon mon laui*

Diagram illustrating language relationships:

```

graph TD
    English --- Hindi
    English --- Marathi
    English --- Bengali
    English --- Manipuri
    Hindi --- Marathi
    Hindi --- Bengali
    Marathi --- Bengali
    Manipuri --- English
  
```

Definitive Pronoun (DP) This blanket is very soft [intensifier relationship] Meaning graph.

Definitive relation Stage Relationship

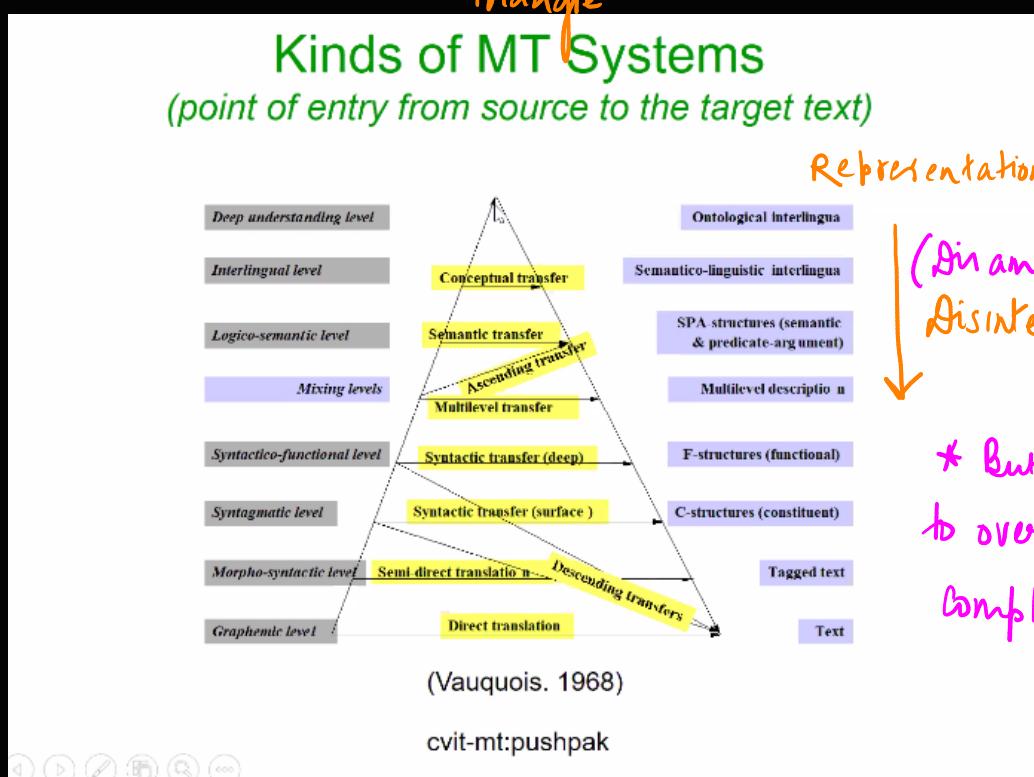
English → DP Subject Verb Object / Hindi → DP Subject Object Verb Copular
 Bangla → DP Subject Object Verb Copula ?? / Marathi = " (same) /
 Manipuri → Object (Thin) Object² Verb Copular (No intensifier (very))

- * Hindi, Bangla & Marathi → almost 1 to 1 map
- * Bangla classifier → definitive Pronouns Rules *
- * Verb going to the end from English → Hindi, Marathi, Manipuri - is called Structural Divergences
- * Very soft changing to soft soft (Reduplication) in Manipuri is called Lexical Divergences

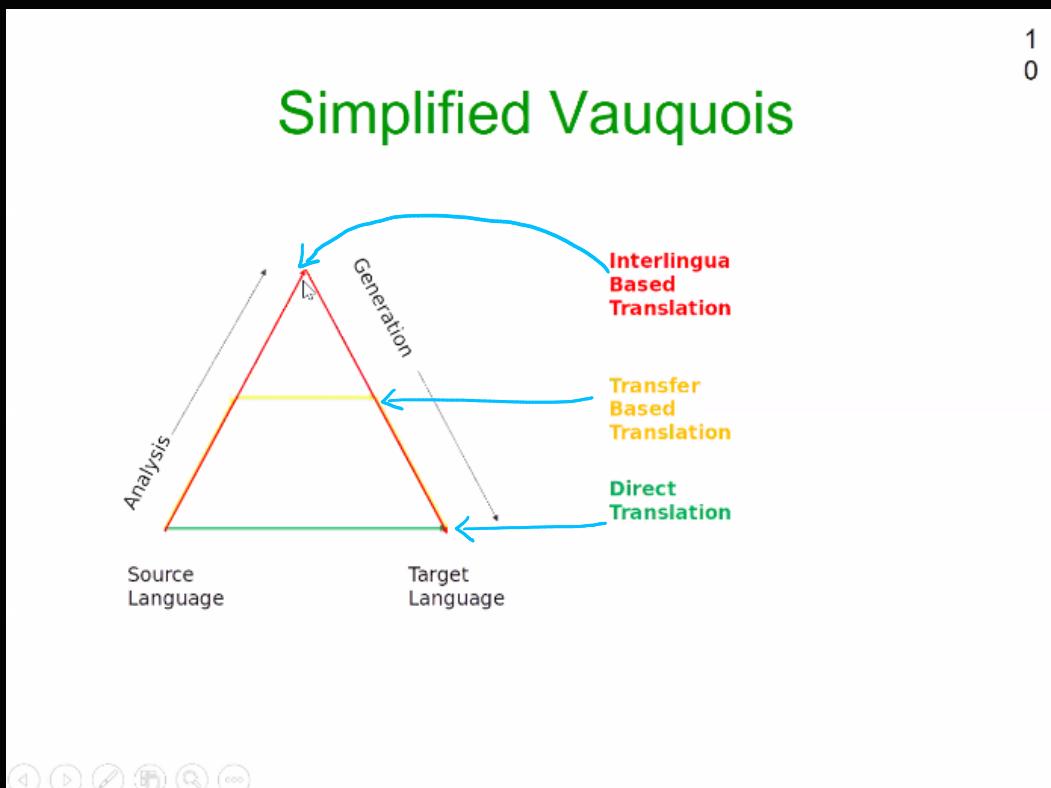
In ML algorithms will pick up this patterns automatically (NN's)

* NN's are data hungry & all Rule Based systems have the problem of false +ve's & false -ve's.

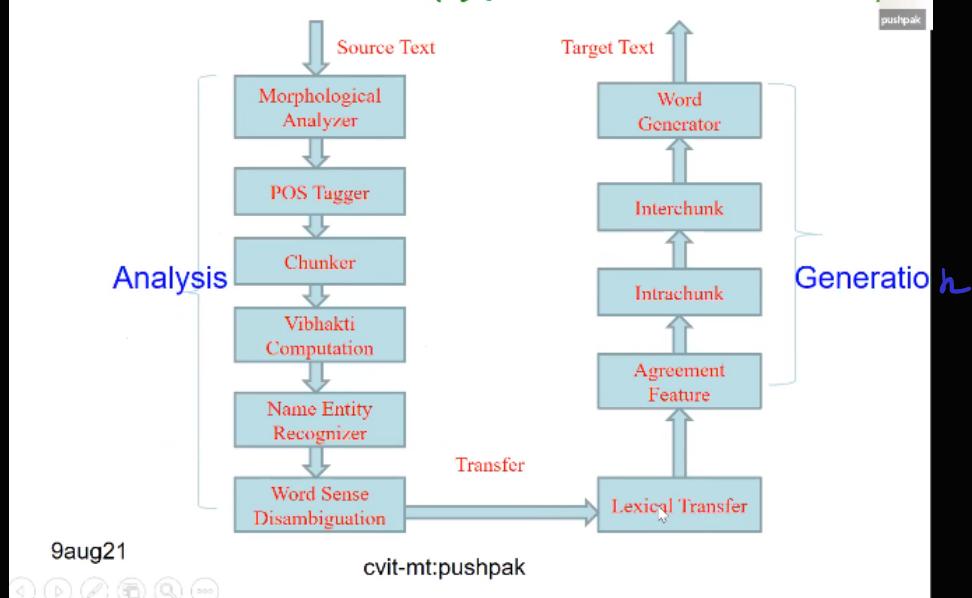
Vauquois
Triangle



* We can retain ambiguity & still do a successful target sentence generation : { Different languages can work with different levels of ambiguity }



Rule based MT (typical architecture)



Key Research Areas

Machine Translation

Sentiment Analysis

Information Retrieval

Lexical Semantics

Information Extraction

Cognitive NLP

*Linguistics is the eye and computation
the body!*

Challenges of IL Computing (1/2)

- **Scale and Diversity:** 22 major languages in India, written in 13 different scripts, with over 720 dialects
- **Code Mixing** ("kyo ye hesitation?"); **Gerundification** ("gaadi chalaoing")
- **Absence of basic NLP tools and resources:** ref nlp pipeline
- **Absence of linguistic tradition for many languages**

↳ Mixing Language

↳ Very Rich

Pushpak Bhattacharyya, Hema Murthy, Surangika Ranathunga and Ranjiva Munasinghe, *Indic Language Computing*, CACM, V 62(11), November 2019.

ILT Challenges (2/2)

- Script complexity and non-standard input mechanism: InScript Non-optimal
- Non-standard transliteration (“mango” → ‘am”, “aam”, Am”)
- Non-standard storage: proprietary fonts
- Challenging language phenomena: Compound verbs (“has padaa”), morph stacking (“gharaasamorchyaanii”)
- Resource Scarcity

↳ Qwerty keyboard issues with Indian typing

→ Media have their own fonts not following Unicode Standards
→ tokens stacked together → word

Indian Language SMT (2014)

	hi	ur	pa	bn	gu	mr	kK	ta	te	ml	en
hi	61.28	68.21	34.96	51.31	39.12	37.81	14.43	21.38	10.98	29.23	
ur	61.42		52.02	29.59	39.00	27.57	28.29	11.95	16.61	8.65	22.46
pa	73.31	56.00		29.89	43.85	30.87	30.72	10.75	18.81	9.11	23.83
bn	37.69	32.08	31.38		28.14	22.09	23.47	10.94	13.40	8.10	18.76
gu	55.66	44.12	45.14	28.50		32.06	30.48	12.57	17.22	8.01	19.78
mr	45.11	32.60	33.28	23.73	32.42		27.81	10.74	12.89	7.65	17.62
kK	41.92	34.00	34.31	24.59	31.07	27.52		10.36	14.80	7.89	17.07
ta	20.48	18.12	15.57	13.21	16.53	11.60	11.87		8.48	6.31	11.79
te	28.88	25.07	25.56	16.57	20.96	14.94	17.27	8.68		6.68	12.34
ml	14.74	13.39	12.97	10.67	9.76	8.39	9.18	5.90	5.94		8.61
en	28.94	22.96	22.33	15.33	15.44	12.11	13.66	6.43	6.55	4.65	

Baseline PBSMT - % BLEU scores (S1)

- Clear partitioning of translation pairs by language family pairs, based on translation accuracy.

scores fell due to Dravidian influence

Extreme agglutination



* MT will scarce data in a problem → Challenge (*)

Subword Based MT came into picture after this results (specially with low Dravidian family)

Mitigating the Resource Problem

Three ways (1/2)

(1) Artificially boost the resource



– Subword based NLP

- Characters, Syllables, Orthographic Syllables, Byte Pair Encoding
 - Given, “khaa+uMgaa → will+eat” AND “jaa+rahaa_hE → is+going”
 - Produce “khaa+rahaa_hE → is+eatin”

→ doesn't appear in training but

Subword Training helps

↓
needs to be a scoring to prevent non linguistic constructs.

Three ways (2/2)

(2) Take help from another language

– Cooperative NLP

(3) Use “higher level language properties”

e.g., Part of Speech, Sense ID etc.

But there is a pitfall- NLP's “Law of Trade off”

- Trade Off:

- Precision vs. Recall*

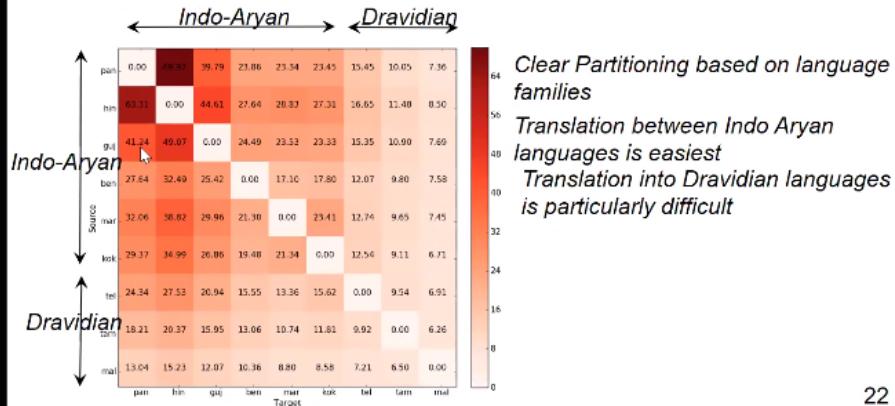
- Sparsity vs. Ambiguity*

- Information_Injection vs. Topic_Drift*



Word level translation (BLEU score)

pushpak



22

9aug21

cvit-mt:pushpak



Subwords (for “jaauMgaa”)

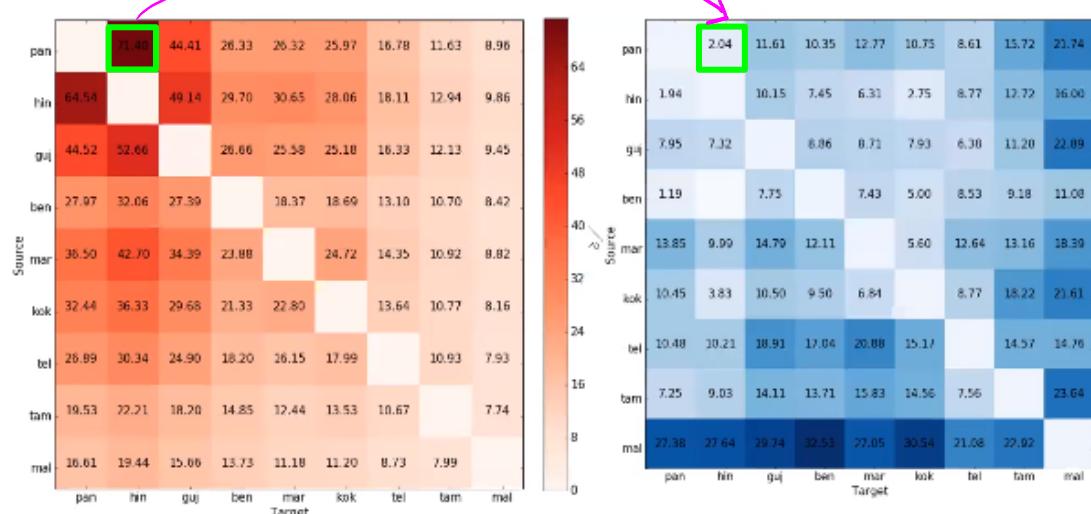
- Characters: “j+aa+u+M+g+aa”
- Morphemes: “jaa+”uMgaa”
- Syllables: “jaa+”uM”+gaa”
- Orthographic syllables: “jaau”+”Mgaa”
- BPE (depends on corpora, statistically frequent patterns): both “jaa” and “uMgaa” are likely

* Breaking the words into syllables in a complex task in itself



Morph level translation

pushpak



BLEU
scores

9aug21



% improvement over word level
scores

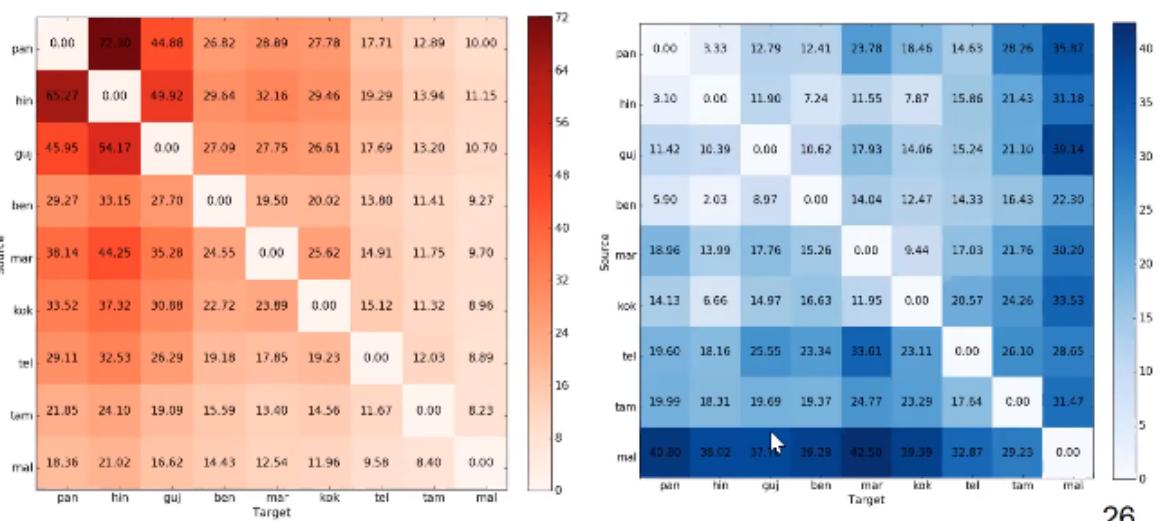
cvit-mt:pushpak

25

↑ Subword Based Translation.
∴ Subwords can isolate the morphemes

BPE level translation

pushpak



more ↑ %. seen ..

Factor based SMT

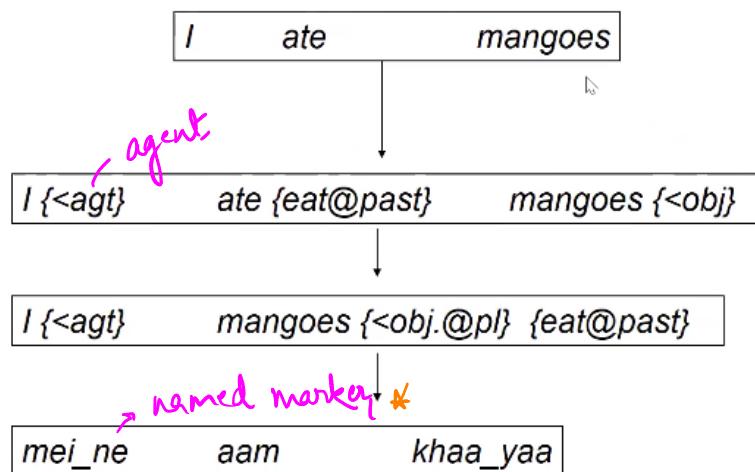
Ananthakrishnan Ramanathan, Hansraj Choudhary, Avishek Ghosh and Pushpak Bhattacharyya, Case markers and Morphology: Addressing the crux of the fluency problem in English-Hindi SMT, ACL-IJCNLP 2009, Singapore, August, 2009.

9aug21

cvit-mt:pushpak

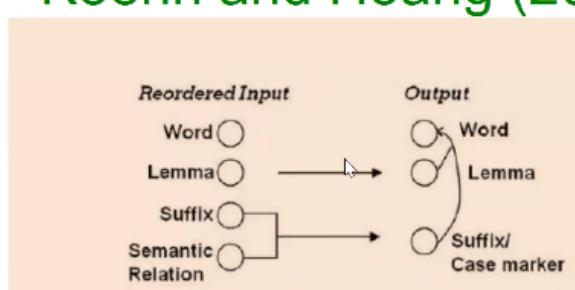
27

Semantic relations+Suffixes→Case Markers+inflections



Hindi →

Our Factorization based on Koehn and Hoang (2007)



1. a lemma to lemma translation factor (boy → लड़का (ladak))
2. a suffix + semantic relation to suffix/case marker factor (-s + subj → ए (e))
3. a lemma + suffix to surface form generation factor (लड़का + ए (ladak + e) → लड़के (ladake))

9aug21

cvit-mt:pushpak

29

} Rule
Operated

* as the verb
ate is
transitive.

* Read
←

Experiment: Corpus Statistics

	#sentences	#words
Training	12868	316508
Tuning	600	15279
Test	400	8557

Very small corpus

Results: The impact of suffix and semantic factors

Model	BLEU	NIST
Baseline (surface)	24.32	5.85
lemma + suffix	25.16	5.87
lemma + suffix + unl	27.79	6.05
lemma + suffix + stanford	28.21	5.99

Results: The impact of reordering and semantic relations

Model	Reordering	BLEU	NIST
surface	distortion	24.42	5.85
surface	lexicalized	28.75	6.19
surface	syntactic	31.57	6.40
lemma + suffix + stanford	syntactic	31.49	6.34

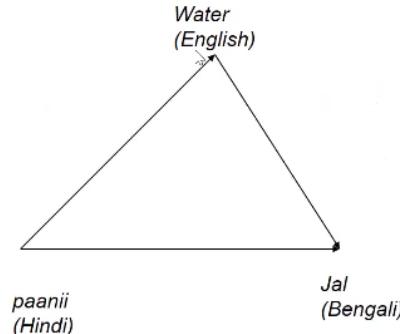
Subjective Evaluation: The impact of reordering and semantic relations

Model	Reordering	Fluency	Adequacy	#errors
surface	lexicalized	2.14	2.26	2.16
surface	syntactic	2.6	2.71	1.79
lemma + suffix + stanford	syntactic	2.88	2.82	1.44

Cooperative NLP: Pivot Based MT

Raj Dabre, Fabien Cromiere, Sadao Kurohashi and
Pushpak Bhattacharyya, Leveraging Small Multilingual
Corpora for SMT Using Many Pivot Languages,
NAACL 2015, Denver, Colorado, USA, May 31 - June
5, 2015.

Triangulation



L1 → bridge → L2 (Wu and Wang 2009)

- Resource rich and resource poor language pairs
- Question-1: How about translating through a 'bridge'?
- Question-2: how to choose the bridge?

Mathematical preliminaries

$$e_{\text{best}} = \arg \max_e p(e|f) \\ = \arg \max_e p(f|e)p_{\text{LM}}(e)$$

Where $p(f|e)$ is given by:

$$p(f|e) = p(\bar{f}^I | \bar{e}^I) = \prod_{i=1}^I \phi(\bar{f}_i | \bar{e}_i) d(a_i - b_{i-1}) p_W(\bar{f}_i | \bar{e}_i, a)^{\gamma}$$

▷

$$\phi(\bar{f}_i | \bar{e}_i) = \sum_{\bar{p}_i} \phi(\bar{f}_i | \bar{p}_i) \phi(\bar{p}_i | \bar{e}_i)$$

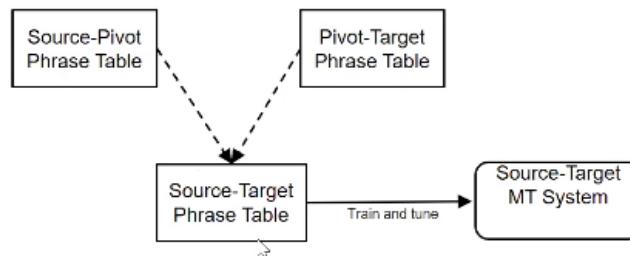
f = target language .
 e = source sentence

e = source language .
 e = source sentence

$$p_W(\bar{f}_i | \bar{e}_i, a) = \prod_{l=1}^n \frac{1}{|m/(l,m) \in a|} \sum_{V(l,m) \in a} w(f_l | e_l)$$

* Introduction of
pivot in scarce
resource problem
is helpful.
↓
few shot ML

Triangulation approach



- Important to induce language dependent components such as phrase translation probability and lexical weight

Oct 19, 2014

FAN, Pushpak Bhattacharyya

48

Mauritian Creole (MCR) → French (FR) → English (E)

- MCR and FR share vocabulary and structure

Vocabulary matching

French	Creole	English
avion	Avion	aeroplane
bon	Bon	good
gaz	Gaz	gas
bref	bref	brief
pion	pion	pawn

Source → pivot → target

9aug21

cvit-mt:pushpak

39

Experiment on MCR→FR→E

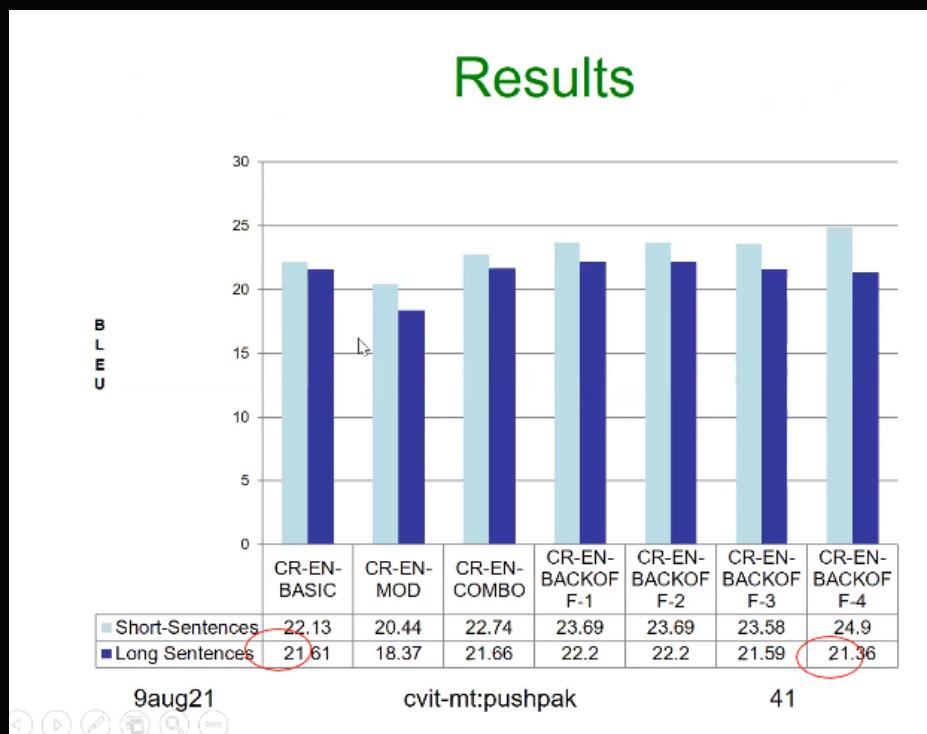
Language pair	#Sentences	#unique words (L1-L2)
En-Fr	2000000	127405- 147812
En-Cr (train + tune)	25010	16294-17389
En-Cr (test)	284 (142 short + 142 long)	1168-1070 + 3562-3326
Fr-Cr	18354	13769-13725

9aug21

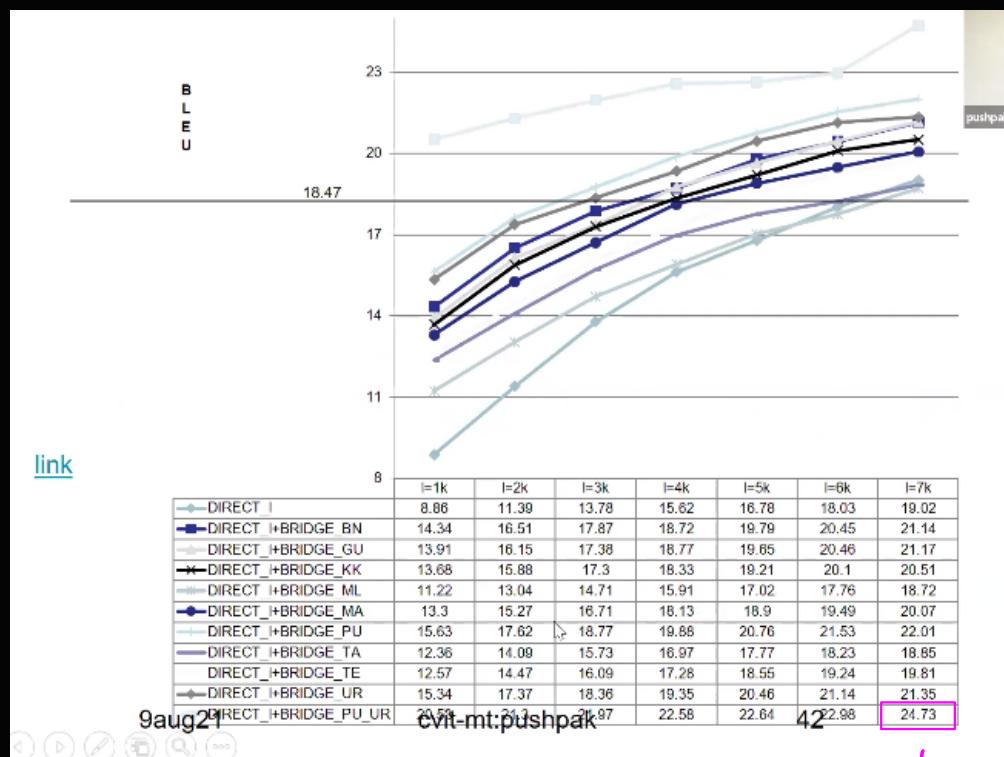
cvit-mt:pushpak

40

Results



Long sentences showed better results

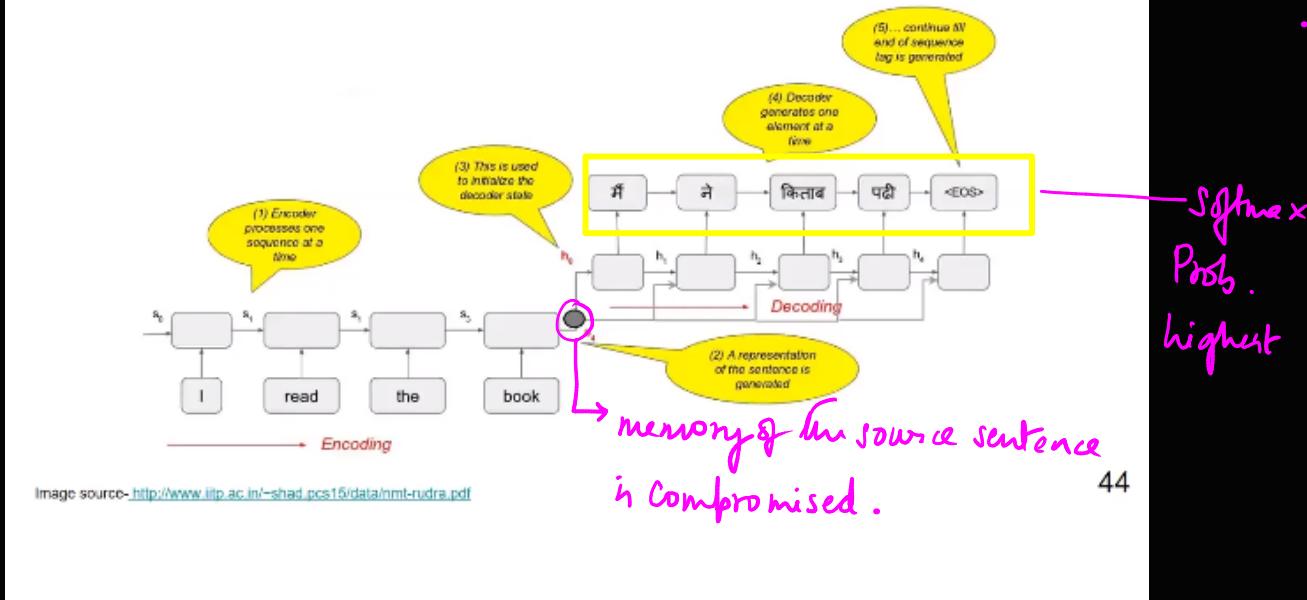


Punjabi & Urdu as pivots. → good score

Neural Machine Translation

pushpak

Encoder-Decoder model



44

↓

Attention Networks (Transformers) came into picture.

Some representative accuracy figure for Indian Language NMT (highly constrained domain)

Language pair	BLEU score
Hi - Mr	31.25
Hi - Pa	63.38
Pa - Hi	68.31
Hi - Gu	49.98
Gu - Hi	53.22 (↑ from 53.09 from SMT)

Unsupervised Neural Machine Translation

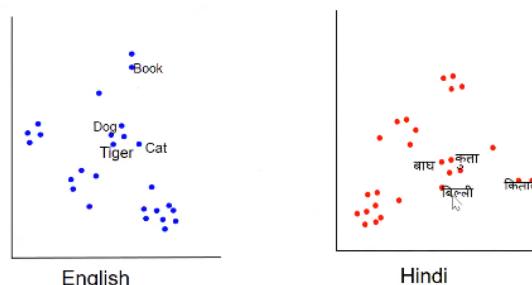
Sukanta Sen, Kamal Kumar Gupta, Asif Ekbal and Pushpak Bhattacharyya. Multilingual Unsupervised NMT using Shared Encoder and Language-Specific Decoders. 57th Annual Meeting of the Association for Computational Linguistics (ACL 2019), Florence, Italy. 2019

Background: Unsupervised NMT (Artetxe et al., 2018)

1. Consists of one shared Encoder and multiple Decoders
2. Fixed cross-lingual embedding at encoder side
3. Denoising Auto-encoding
4. Back-translation

Mikel Artetxe, Gorka Labaka, Eneko Agirre, and Kyunghyun Cho. 2018. Unsupervised Neural Machine Translation. ICLR

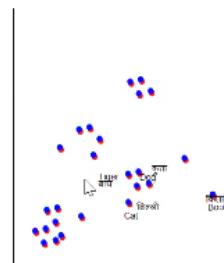
Background: Geometric Structure of Words (use iso-metricity)



After Cross-lingual Mapping

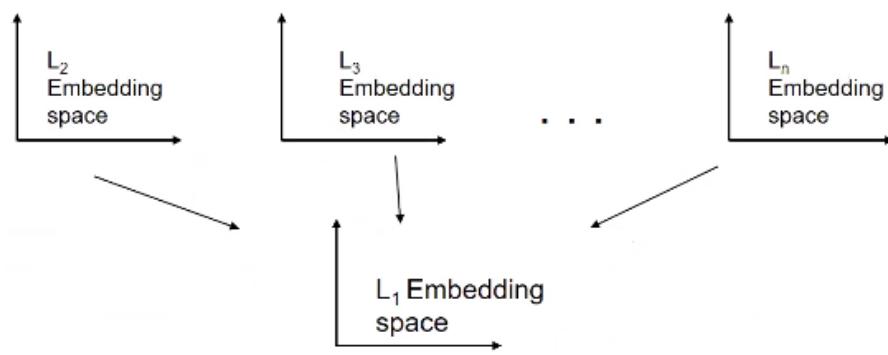
This involves strong assumption that embedding spaces across languages are isomorphic, which is not true specifically for distance languages (Søgaard et al. 2018). However, without this assumption unsupervised NMT is not possible.

Søgaard, Anders, Sebastian Ruder, and Ivan Vulić. 2018. On the limitations of unsupervised bilingual dictionary induction. ACL.



Multilingual Embedding

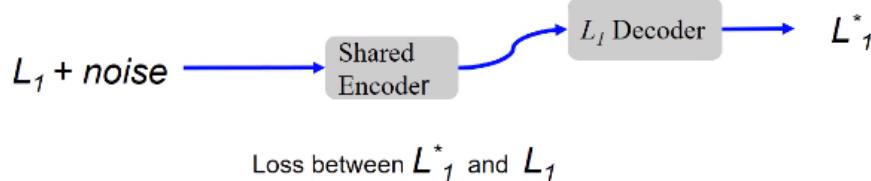
$L_1 \rightarrow$ English.
 $L_2, L_3, L_4 \rightarrow$
 Non English.



We pairwise map all non-English embedding spaces into the English embedding space using Conneau et al., 2018



Denoising Auto-encoding for L_1



Noise is introduced through swapping/deletion of words. For example,

Original: **An investment in knowledge pays the best interest .**

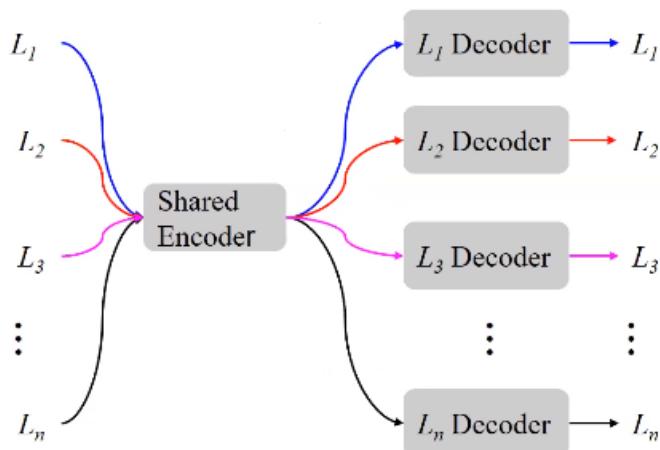
After noising: **investment An in knowledge pays best the interest .**

Model tries to predict original from noisy sentence. We do it for sentences in all four languages.



Denoising and Back-Translation

* Takes
cross path
via the
shared
Encoder *



Results

System	newstest2013			newstest2014		
	Base	Multi	▲	Base	Multi	▲
Fr→En	13.81	14.47	+0.66	14.98	15.76	+0.78
Es→En	13.97	15.45	+1.48	-	-	-
En→Fr	13.28	13.71	+0.43	14.57	14.69	+0.12
En→Es	14.01	14.82	+0.81	-	-	-
De→En	11.30	11.94	+0.64	10.48	11.21	+0.73
En→De	7.24	8.09	+0.85	6.24	6.77	+0.53

Table 1: BLEU scores on *newstest2013* and *newstest2014*. ▲ shows improvements over bilingual models. Spanish (Es) is not part of the *newstest2014* test set. **Base:** Baseline. **Multi:** Multilingual

Sample Outputs (French→ English)

Source	Reference	Bilingual (Baseline)	Multilingual (Ours)
La préparation à gérer une classe dans un contexte nord-américain, québécois.	Preparation to manage a class in a North-American and Quebec context.	The build-up to manage a class in a Australian, Australian.	The preparation to handle a class in a Latin American context.
Il va y avoir du changement dans la façon dont nous payons ces taxes.	There is going to be a change in how we <u>pay</u> these <u>taxes</u> .	There will be the change in the course of whom we <u>owe</u> these <u>bills</u> .	There will be the change in the way we <u>pay</u> these <u>taxes</u> .

Sample Outputs (Spanish→ English)

Source	Reference	Bilingual (Baseline)	Multilingual (Ours)
Los estudiantes, por su parte, aseguran que el curso es uno de los más <u>interesantes</u> .	Students, meanwhile, say the course is one of the most <u>interesting</u> around.	The students, by their part, say the practice is one of the most <u>intriguing</u> .	The students, by their part, say the course is one of the most <u>interesting</u> .
No duda en contestar que nunca aceptaría una <u>solicitud</u> de una persona desconocida.	He does not hesitate to reply that he would never accept a <u>request</u> from an unknown person.	No doubt ever answering doubt it would never accept an <u>argument</u> an unknown person.	No doubt in answer that he would never accept a <u>request</u> of a unknown person.

Sample Outputs (German→ English)

Source	Reference	Bilingual (Baseline)	Multilingual (Ours)
Auch diese Frage soll letztlich Aufschluss darüber geben, welche Voraussetzungen es für die Entstehung von Leben gibt.	This question should also provide information regarding the preconditions for the <u>origins</u> of life.	This question will also ultimately give clues about what there are for the <u>evolution</u> of life.	This question will ultimately give clues to how there are conditions for the <u>emergence</u> of life.
Ihm werde weiterhin vorgeworfen, unerlaubt geheime Informationen weitergegeben zu haben.	He is still accused of passing on secret information without authorisation.	Him will continue to be accused of stealing unlawful information.	Him would continue to be accused of illegally leaking secret information.

Results (Zero-shot Translation in Unsupervised NMT)

→	Es	Fr	De
Es	-	13.92	4.78
Fr	13.87	-	4.59
De	7.40	6.78	-

Table 2: BLEU scores of translation between non-English languages on *newstest2013*. Consider rows are source and columns are target. The network is not trained for these language pairs and still it is possible to translate between these pairs by using the shared encoder and language specific decoders.

Sample Outputs (Unseen pairs)

Source	Reference	Multilingual
Les dirigeants républicains justifient leur politique par la nécessité de lutter contre la fraude électorale.	French→Spanish Los dirigentes republicanos justificaron su política por la necesidad de luchar contra el fraude electoral.	Los dirigentes republicanos <OOV> su política por la necesidad de luchar contra la fraude electoral.
Chacun sait que son livre fait partie de cet édifice.	French→German Jeder weiß, dass sein Buch Teil dieses Gebäudes ist.	Jeder weiß, dass sein Buch Teil seines Gebäudes machte.
Seine Zahlen auf Ebene der internationalen Turniere sind beeindruckend.	German→Spanish Sus números a nivel de torneos internacionales son impresionantes.	Sus cifras sobre el nivel de torneos internacionales son impresionantes.
Diese Einschränkungen sind nicht ohne Folgen.	German→French Ces restrictions ne sont pas sans conséquence.	Ces restrictions ne sont pas sans conséquences.
Tomemos por caso la elección directa del presidente , que ha sido un logro de la presión pública.	Spanish→German Nehmen Sie nur einmal die direkte Wahl des Präsidenten, die ein Verdienst des öffentlichen Drucks war.	Nehmen Sie über die direkte Wahl des Präsidenten, hat dies ein Erfolg ein der öffentlichen Druck.
Las inversiones en la materia superan los 1.5 billones de dólares.	Spanish→French Les investissements dans ce domaine dépassent les 1,5 milliards de dollars.	Les investissements dans la matière dépassent les 1,5 milliards de dollars.

Unsupervised NMT with filtering on Back Translation

Jyotsana Khatri and Pushpak Bhattacharyya, [Filtering Back-Translated Data in Unsupervised Neural Machine Translation](#),
28th Int'l Conf on Computational Linguistics (COLING20),
Online Conference, December 8-13, 2020.

Central Idea

- Filter back-translated data by giving more weight to good pseudo parallel sentence pairs.
- Quality of pseudo parallel sentence pairs measured using round trip BLEU

The screenshot shows the HEMAT Sentence Translation interface. At the top, there is a navigation bar with links for Sentence Translation, Document Translation, Feedback, Home, and a user profile for Pushpak Bhattacharyya. Below the navigation bar, the main title is "Sentence Translation". A sub-header indicates it is for the "Hindi-English Machine Aided Translation For Judicial Domain".

The interface features two main input fields:

- Source Text (Left):** "The toy train from Kalka to Shimla is considered as the most beautiful rail line in India."
- Target Language Selection (Top):** "Select Translation Direction" dropdown set to "English → Hindi".
- Target Text (Right):** "कालका से शिमला तक की खिलौना रेल को भारत की सबसे सुंदर रेल-रेखा माना जाता है।"

Below this, another section shows:

- Source Text (Left):** "Select Target Langauge" dropdown set to "Hindi".
- Target Language Selection (Top):** "Select Target Langauge" dropdown set to "Marathi".
- Target Text (Right):** "कालका पासून सिमल्यापर्यंत खेळणी रेल्वेला भारतातील सर्वात सुंदर रेल्वे रेषा मानले जाते।"
- Translation in target language:** "Translation in target language" dropdown set to "Marathi".
- Target Text (Bottom):** "कालकापासून सिमल्यापर्यंत खेळणी रेल्वेला भारतातील सर्वात सुंदर रेल्वे रेषा मानले जाते।"

Sentence Translation

Hindi-English Machine Aided Translation For Judicial Domain

Select Translation Direction	
English → Hindi	
The criminal was apprehended.	अपराधी को पकड़ लिया गया था।

Select Source Langauge	Enter text in source language:
Hindi	अपराधी को पकड़ लिया गया था।

Select Target Langauge	Translation in target language:
Marathi	गुन्हेगार पकडला गेला .

Select Source Langauge	Enter text in source language:
Marathi	गुन्हेगार पकडला गेला .

Select Target Langauge	Translation in target language:
English	the culprit was arrested .

* Meaning is serviced but this method of evaluation by completing a circle for a pair of sentence in a stringent form of evaluation. { or else errors can accumulate to give bad target output }

Another Example ; when o/p is disturbed { post editing is required }

- * Work Error Rate
- * Cognitive Based Load
- * Eye Tracking based Cognitive Load.

Sentence Translation

Hindi-English Machine Aided Translation For Judicial Domain

Select Translation Direction

English → Hindi ▾

Most of the court cases of India are traffic violation cases.

भारत के अधिकांश न्यायालय मामले यातायात उल्लंघन के मामले हैं। I

Select Source Langauge

Hindi

Enter text in source language:

भारत के अधिकांश न्यायालय मामले यातायात उल्लंघन के मामले हैं।

Select Target Langauge

Marathi

Translation in target language:

भारतातील बहुतांश न्यायालये वाहतूक नियमांचे उल्लंघन करत आहेत .

I

Select Source Langauge

Marathi

Enter text in source language:

भारतातील बहुतांश न्यायालये वाहतूक नियमांचे उल्लंघन करत आहेत .

Select Target Langauge

English

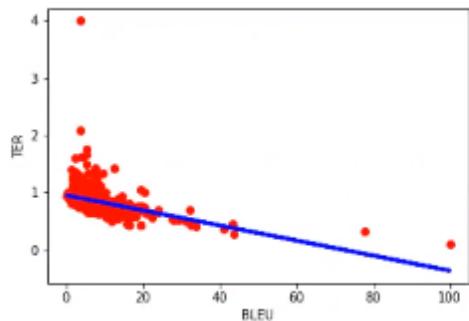
Translation in target language:

most of the courts in india are violating traffic rules .

I

* Postediting: Most of the court cases in India are for violating traffic rules

Bible



$$\text{Equation of Line : } y = -0.013x + 0.955$$

$$\text{Slope}(m) = -0.013$$

$$\text{Intercept}(c) = 0.955$$

Lang-pair	# sentences in source language	# sentences in target language	Training dataset
en-fr	5M	5M	WMT News Crawl
en-de	5M	5M	WMT News Crawl
en-ro	5M	2.28M	WMT News Crawl

Table 1: Details of training data

Lang-pair	Validation data	Test data
en-fr	Newstest 2013	Newstest 2014
en-de	Newstest 2013	Newstest 2016
en-ro	Newsdev 2016	Newstest 2016

Table 2: Details of test data

Results

Method	en-fr	fr-en	en-ro	ro-en	en-de	de-en
Song et al., 2019	26.59	25.42	25.53	24.8	17.62	24.78
Song et al., 2019 + Filtering	26.37	25.5	25.29	24.64	17.51	24.87
Lample and Conneau, 2019	27.95	27.02	26.26	25.8	18.26	25.22
Lample and Conneau, 2019 + Filtering	28.4*	27.69*	26.96*	26.24*	17.35	25.91*

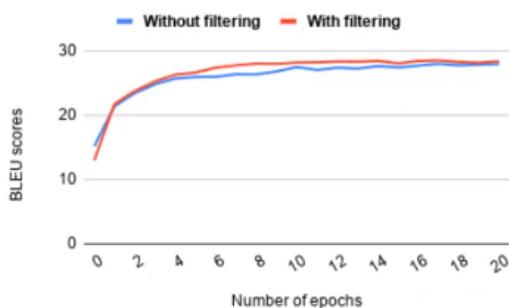


Figure 1: Progress with iterations: BLEU scores for each epoch for en-fr

Neural Indian Language MT

Comprehensive ILNMT, including some points from the bahubhashak project

Shubham Dewangan, Shreya Alva, Nitish Joshi and Pushpak Bhattacharyya, [Experience of Neural Machine Translation between Indian Languages](#), Journal of Machine Translation (JMT), Springer, Accepted. [Doi](#)



Recent Publications

- Shubham Dewangan, Shreya Alva, Nitish Joshi and Pushpak Bhattacharyya, *Experience of neural machine translation between Indian languages*, **Machine Translation volume 35, 2021**.
- Akash Banerjee, Aditya Jain, Shivam Mhaskar, Sourabh Dattatray Deoghare, Aman Sehgal, Pushpak Bhattacharya, *Neural Machine Translation in Low-Resource Setting: a Case Study in English-Marathi Pair*, **MT Summit 2021**.
- Tamali Banerjee, Rudra V Murthy, Pushpak Bhattacharya, *Crosslingual Embeddings are Essential in UNMT for distant languages: An English to IndoAryan Case Study*, **MT Summit 2021**.
- Kamal Gupta, Dhanvanth Boppana, Rejwanul Haque, Asif Ekbal, Pushpak Bhattacharyya, *Investigating Active Learning in Interactive Neural Machine Translation*, **MT Summit 2021**.
- Tamali Banerjee, Rudra V Murthy, Pushpak Bhattacharya, Scrambled Translation Problem: A Problem of Denoising UNMT, **MT Summit 2021**.



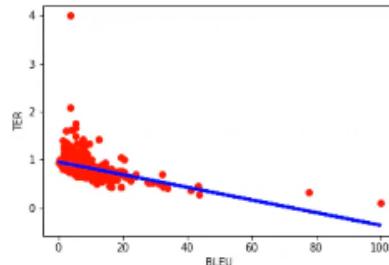
*

TER → Transmission Error Rate

When

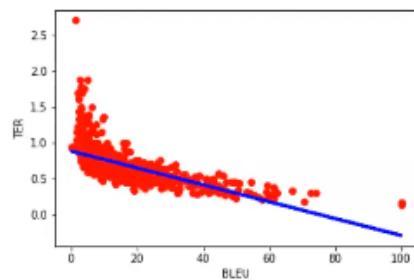
TER = 1, BLEU = 0
= 0, = 1

Bible



Equation of Line : $y = -0.013x + 0.955$
Slope(m) = -0.013
Intercept(c) = 0.955

ILCI

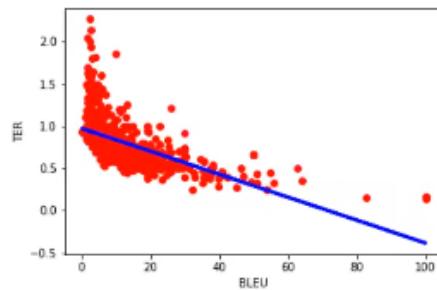


Equation of Line : $y = -0.011x + 0.880$

Slope(m) = -0.011

Intercept(c) = 0.880

PIB

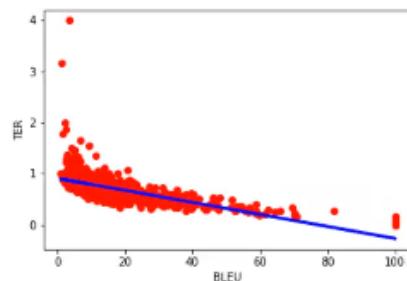


Equation of Line : $y = -0.013x + 0.974$

Slope(m) = -0.013

Intercept(c) = 0.974

PMI

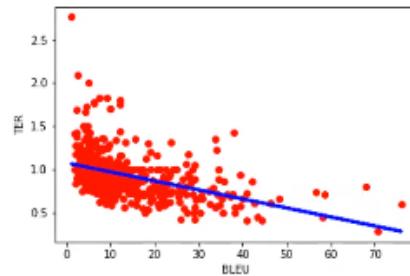


Equation of Line : $y = -0.011x + 0.912$

Slope(m) = -0.011

Intercept(c) = 0.912

Tourism



Equation of Line : $y = -0.010x + 1.075$

Slope(m) = -0.010

Intercept(c) = 1.075

Summary

- MT Paradigms
- Data Driven MT: SMT and NMT
- Tricks of Resource Mitigation
- Unsupervised NMT
- Experience of IL-NMT

Summary on resource mitigation tricks

- Several techniques explored and demonstrated their efficacy.
 - Phrase Table Injection, has great potential to boost BLEU scores, particularly when Dravidian languages are involved.
 - Harnessing monolingual data with back translation, forward translation is advantageous.
 - Enhancements like morph and word feature injection

Conclusions and Future Work

1. Data needed for quick system development, robustness
2. Linguistic insight needed for data curation, annotation and insightful error analysis
3. NMT ensured fluency- syntactic goodness
4. No way out other than various tricks for ameliorating data hunger

Future lies in:

Cross Lingual, Multilingual Word embedding

Final Message

“NLP is a task in Trade Off”

e.g., Not too much of subwords or cooperation

**(beware of ‘ambiguity insertion’),
not too little**

(beware of ‘sparsity’) !!