

# CS626: Speech, Natural Language Processing and the Web

## *Introduction*

Pushpak Bhattacharyya  
Computer Science and Engineering  
Department  
IIT Bombay  
*28<sup>th</sup> July, 2022*

# Sequel Course

- CS772- Deep Learning for Natural Language Processing
- NLP concepts covered in CS626 will see their realization on Deep Neural Nets
- These two courses create a solid platform for launching substantial research and development in NLP

# Ode to *Scientists and Engineers*

Scientists ask WHY

Engineers ask WHY NOT

Scientists wonder at WHAT-IS

Engineers wonder WHAT-COULD-BE

World couldn't do without either.

Scientists STUDY

Engineers MAKE

And ever the twain shall meet.

# What is “Language”

## Oxford English Dictionary

1. the principal method of human ***communication***, consisting of words used in a structured and conventional way and conveyed by speech, writing, or gesture.

"a study of the way children learn Language"

2. a system of communication used by a particular country or community.

"the book was translated into twenty-five  
languages"

# General point: Properties of Human Languages

(George Yule, *“Study of Language”*, 1998)

- **Displacement** (Indicators that change with time and place: I saw him yesterday at the market; I will see him tomorrow in the school)
- **Arbitrariness** (name → Meaning; water, chair)
- **Productivity/creativity** (potentially infinite no. of sentences)
- **Cultural Transmission** (child acquires parent's language)
- **Discreteness** (sound and meaning units separated)
- **Duality** (Surface structure, deep structure)

# What is “Linguistics”

- Scientific study of language, its underlying and governing rules
- **Descriptive:** describe the language objects, language phenomena AS THEY ARE
- **Prescriptive:** prescribe what is allowed and not allowed, e.g. disallow double negative- “/ *did not see nobody in the hall*”; control language behaviour

# What is NLP

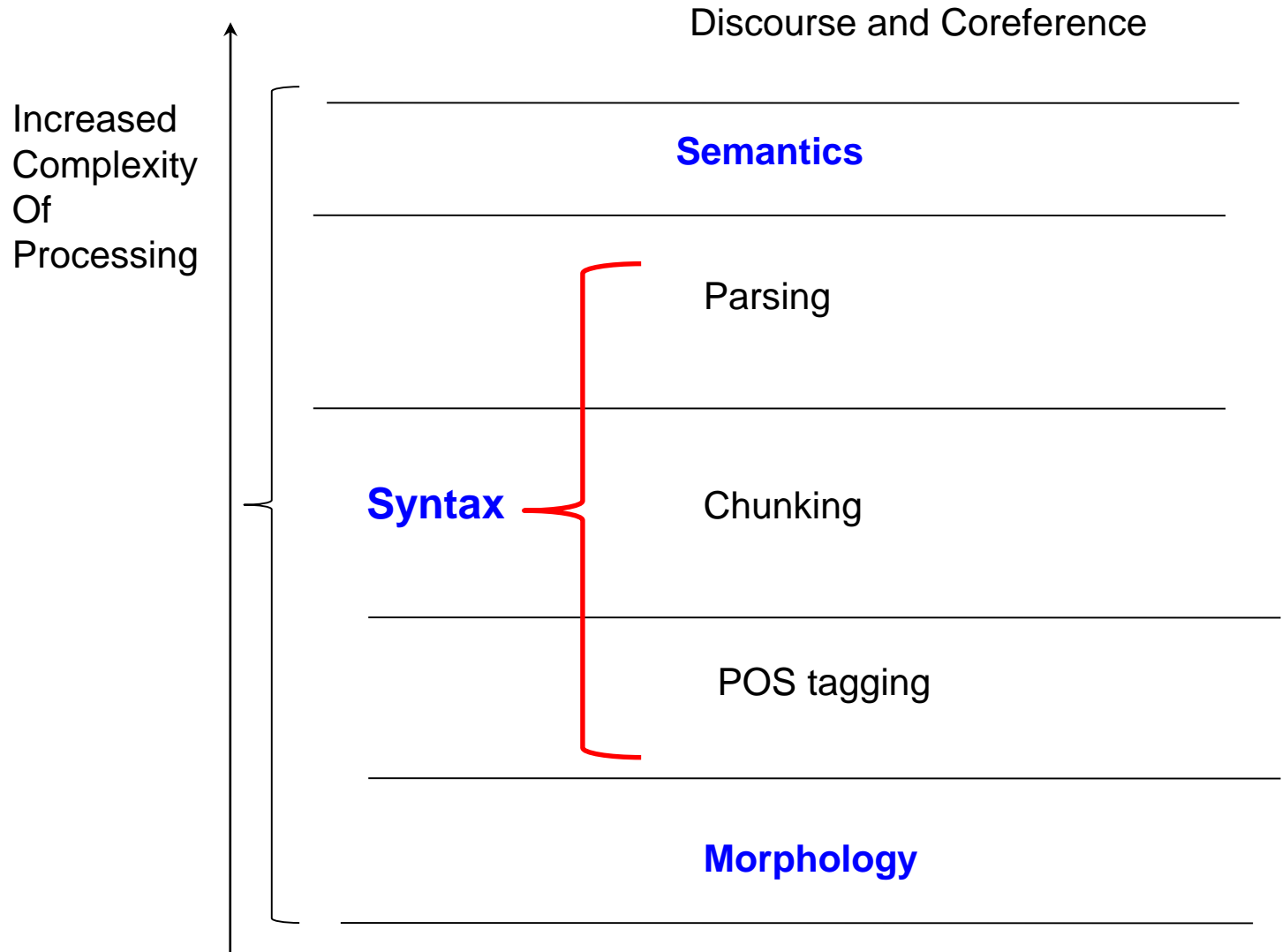
NLP= Language + Computation



*(due to ML)*

= Linguistics + Probability

# NLP Layers





# Linguistic Strata vs. Languages

	Hindi	Swahili	Tamil	<i>etc.</i>
<b>Sound:</b> Phonetics, Phonology				
<b>Structure:</b> Morphology, Syntax				
<b>Meaning:</b> Semantic, Pragmatics				

# Speech-NLP Stack vs. Languages

	Hindi	Swahili	Tamil	<i>etc.</i>
Sound: ASR, TTS				
Structure: MA, POS, NE, Chunker, Parser				
Meaning: SRL, Knowledge Nets, SA-EA-OM, QA, Summarizer				

# Language Typology

Classification according to structural features

# Proto-Language (*Wikipedia*)

Meaning:	Sanskrit	Latin:
"three"	<i>trayas</i>	<i>tres</i>
"seven"	<i>sapta</i>	<i>septem</i>
"eight"	<i>ashta</i>	<i>octo</i>
"nine"	<i>nava</i>	<i>novem</i>
"snake"	<i>sarpa</i>	<i>serpens</i>
"king"	<i>raja</i>	<i>regem</i>
"god"	<i>devas</i>	<i>divus</i> ("divine")

One of the indications that languages descended from a single source


# Word order based

- Object–subject–verb (OSV)
- Object–verb–subject (OVS)
- Subject–verb–object (SVO): English
- Subject–object–verb (SOV): Most Indian Languages
- Verb–subject–object (VSO)
- Verb–object–subject (VOS)

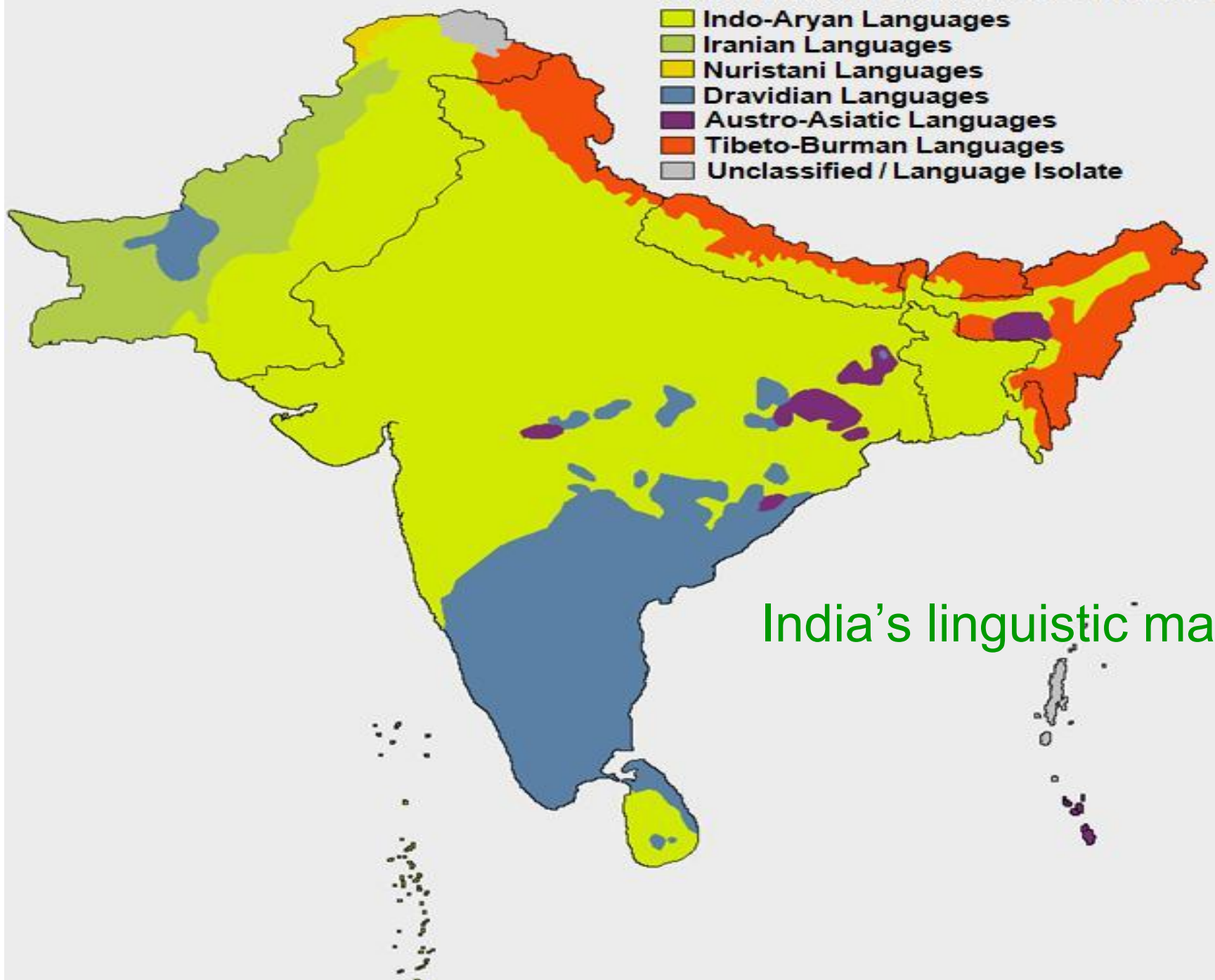
# Dominant Word Order Distribution Across Languages (Wikipedia)

Type	Languages	%	Families	%
<b>SOV (Hindi)</b>	<b>2,275</b>	<b>43,3%</b>	<b>239</b>	<b>65.3%</b>
<b>SVO (English)</b>	<b>2,117</b>	<b>40.3%</b>	<b>55</b>	<b>15%</b>
<b>VSO (tagalog in phillipines)</b>	<b>503</b>	<b>9.5%</b>	<b>27</b>	<b>7.4%</b>
<b>VOS (Malagasy in Madagascar)</b>	<b>174</b>	<b>3.3%</b>	<b>15</b>	<b>4.1%</b>
<b>NODOM (Sanskrit)</b>	<b>124</b>	<b>2.3%</b>	<b>26</b>	<b>7.1%</b>
<b>OVS (Korean and Japanese, many times)</b>	<b>40</b>	<b>0.7%</b>	<b>3</b>	<b>0.8%</b>
<b>OSV (Warao in Venezuela)</b>	<b>19</b>	<b>0.3%</b>	<b>1</b>	<b>0.3%</b>

# Some interesting cases

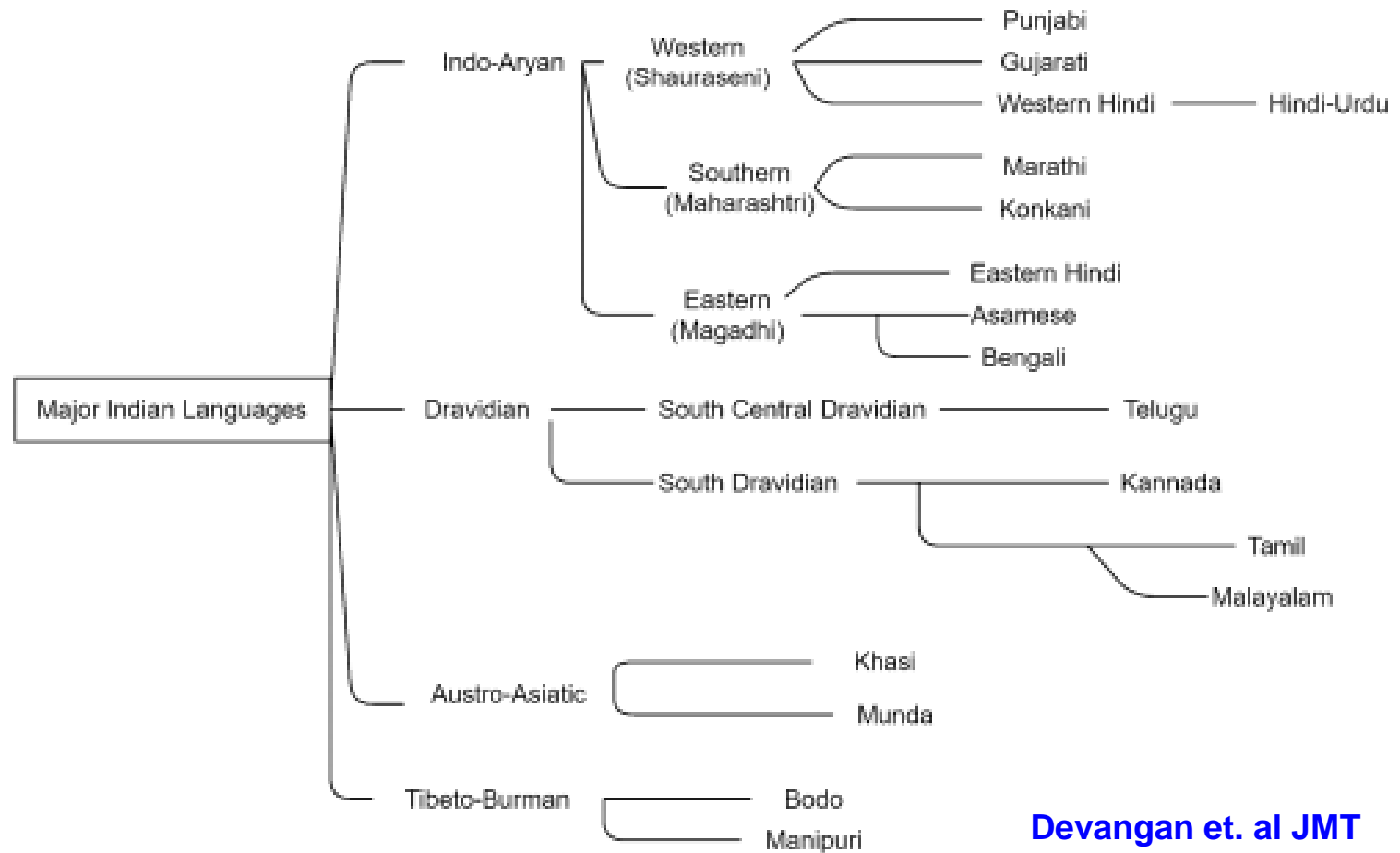
- German word order depends on the position of the main verb (MV)
  - *I know the boy who lives in Berlin*  

  - *Ich kenne den Jungen, der in Berlin lebt*
- Hindi:
  - *Mei us ladke ko jaantaa hu jo barlin me rahataa haai*
  - *Mei us ladke ko jo barlin me rahataa haai jaantaa hu*

## SOUTH ASIAN LANGUAGE FAMILIES



India's linguistic map





Devangan et. al JMT  
2021

**Fig. 1** Tree diagram to illustrate the language closeness of major Indian languages

Main Challenge of NLP:  
**AMBIGUITY**

# An interesting whatsapp conversation (English and Bengali)

Lady A: Yesterday you told me about shop that sells artificial jewellery

<bn>ki naam jeno?</bn> (what did you say was the name?)

Lady B: nykaa

Lady A (offended): What do you mean Madam? Is this the way to talk?

Lady B: <bn>kena ki holo?</bn> (why what happened?)

*Lady A did not reply: she was angry!!!*

# Root cause of the problem: Ambiguity!

- NE-non NE ambiguity (proper noun-common noun)
- Aggravated by code mixing
- “Nykaa”: name of the shop
- Sounds similar to “ন্যাকা” (nyaakaa), meaning somebody “who feigns ignorance/innocence” in a derogatory sense
- An offensive word

# NYKAA Fashion

Inbox (173) x (281) What: x CS772-202: x Google Cal: x https://www: x NEW Th: x Courses | D: x Jewellery O: x + -

nykaafashion.com/jewellery/c/77?root=nav\_3&ptype=listing%2Cjewellery%2Ccategories%2C1%2Cshop-all-jewellery&utm\_content=ads&utm\_s...

tra ₹300 off. Use code: NFAPP300 App Download Help

**NYKAA**  
FASHION

All Brands **Women** Men Kids Home Tech More

Search for products, styles, brands

0

What's New ▾ Indian Wear ▾ Western Wear ▾ Bags ▾ Footwear ▾ Jewellery ▾ Lingerie ▾ Sportswear ▾ Sleep & Lou ▾

Silver 6582 ☐

Rhodium 2102 ☐

18K Gold 1781 ☐

22k Gold 1777 ☐


Rose Gold 1404 ☐

MORE FILTERS

!


Select only 1 category to view more filters

SELECT 1 CATEGORY




BESTSELLER HIDDEN GEMS

**Odette**  
Multi-Color Stone Enticing Long Onyx Neckl...  
₹1,957 ₹5,150 **62% Off**



OFFER

**Fabula**  
Green Meenakari Red Beads & Kundan Ethni...  
₹722 ₹5,553 **87% Off**




LATEST SEASON

**Twenty Dresses by Nykaa Fashion**  
I Am Trending Ear Cuff  
₹487 ₹695 **30% Off**

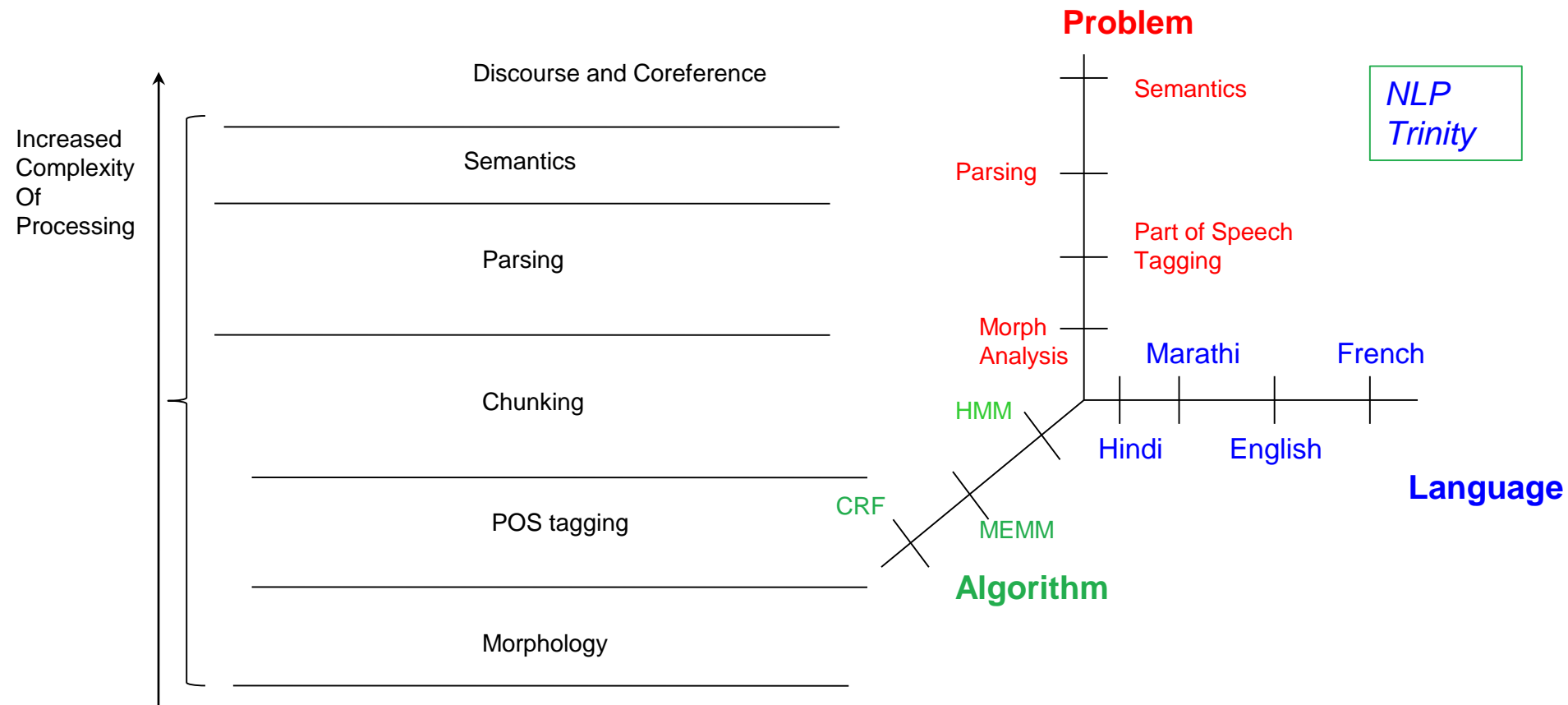
pos-labels-Hin-m....doc ^

Show all x



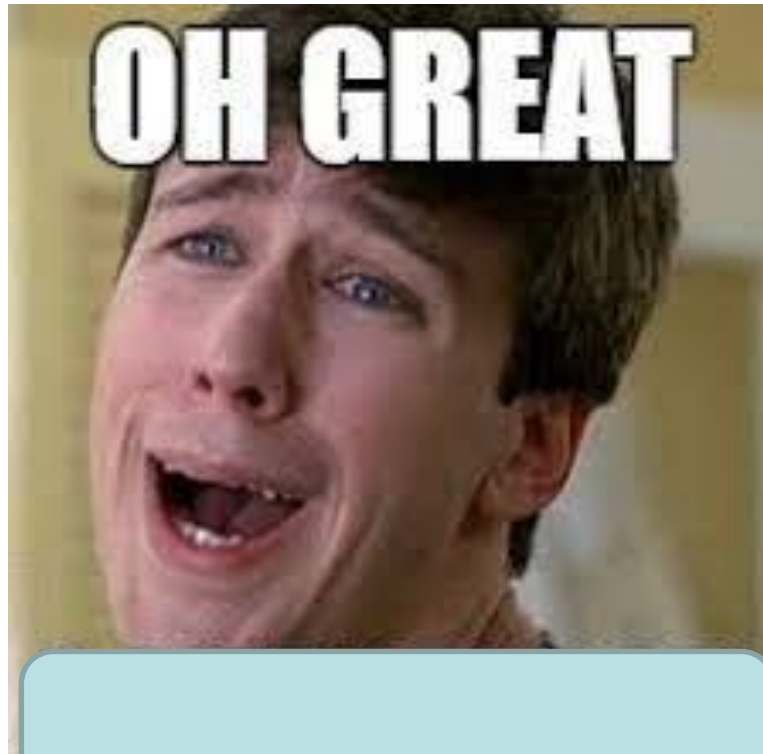
ENG IN 15:07 05-01-2022 18

# Ambiguity at every layer, for every language, for every mode



# Role of Multimodality

- Signals from other modes
- E.g., Sarcasm



**Frequent Observation:  
Data + Classifier > Human  
decision maker !!**

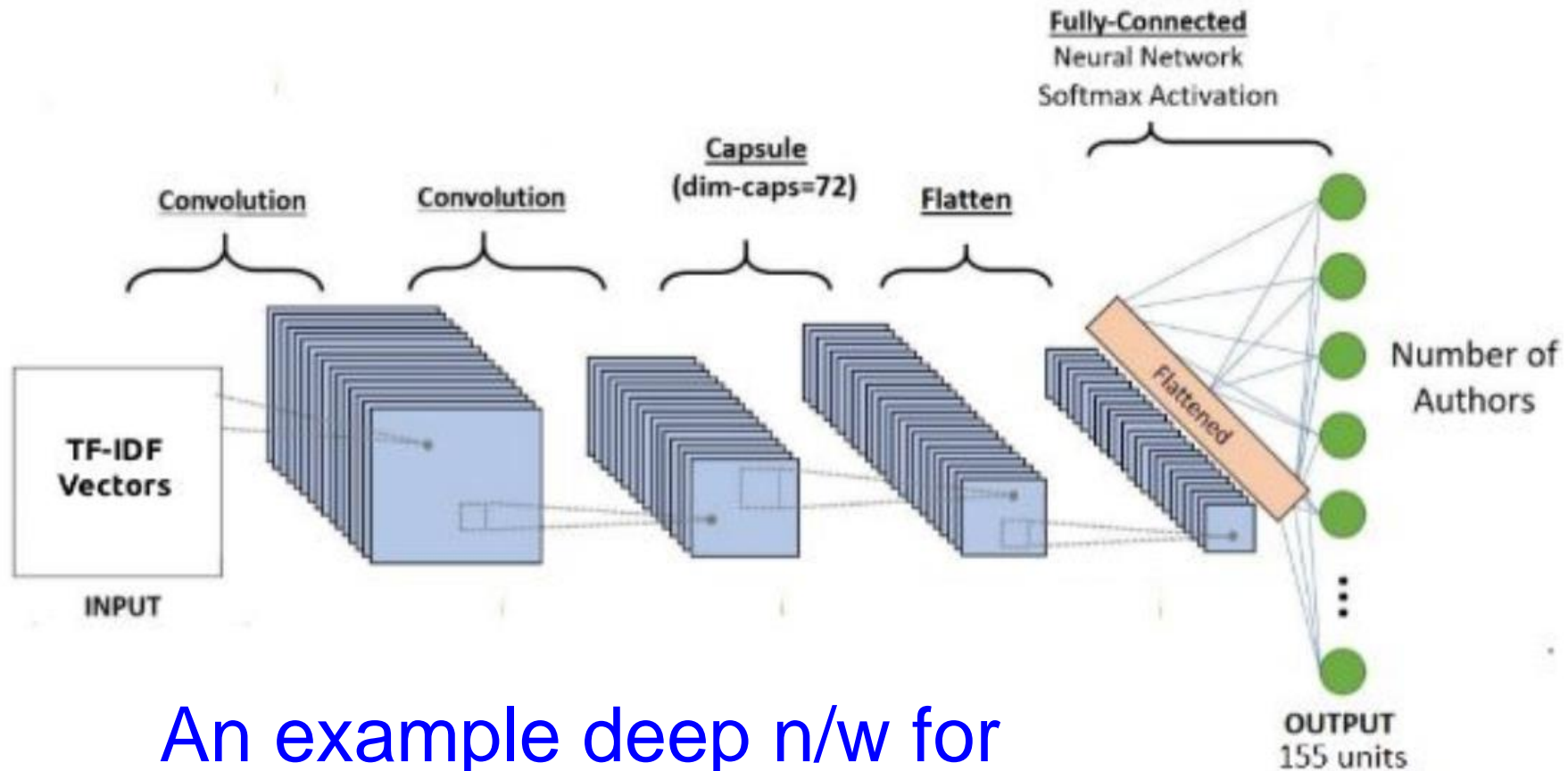
**Case for ML-NLP**



# LEARN from Data with Probability Based Scoring

- With LOTS of data, learn with
  - High precision (small possibility of error of commission)
  - High recall (small possibility of error of omission)
- But depends on human engineered features, i.e., capturing essential properties

# Modern Modus Operandi: End to End DL-NLP



An example deep n/w for author identification

# Problem Knowledge and Deep Learning

- Large number of parameter in DL-NLP:  
Why?
- Fixing large number parameter values need large amounts of data (text for NLP).
- If we **know underlying distribution** then we can make predictions.

**IMP:** The number of needed parameters can be reduced by using knowledge.

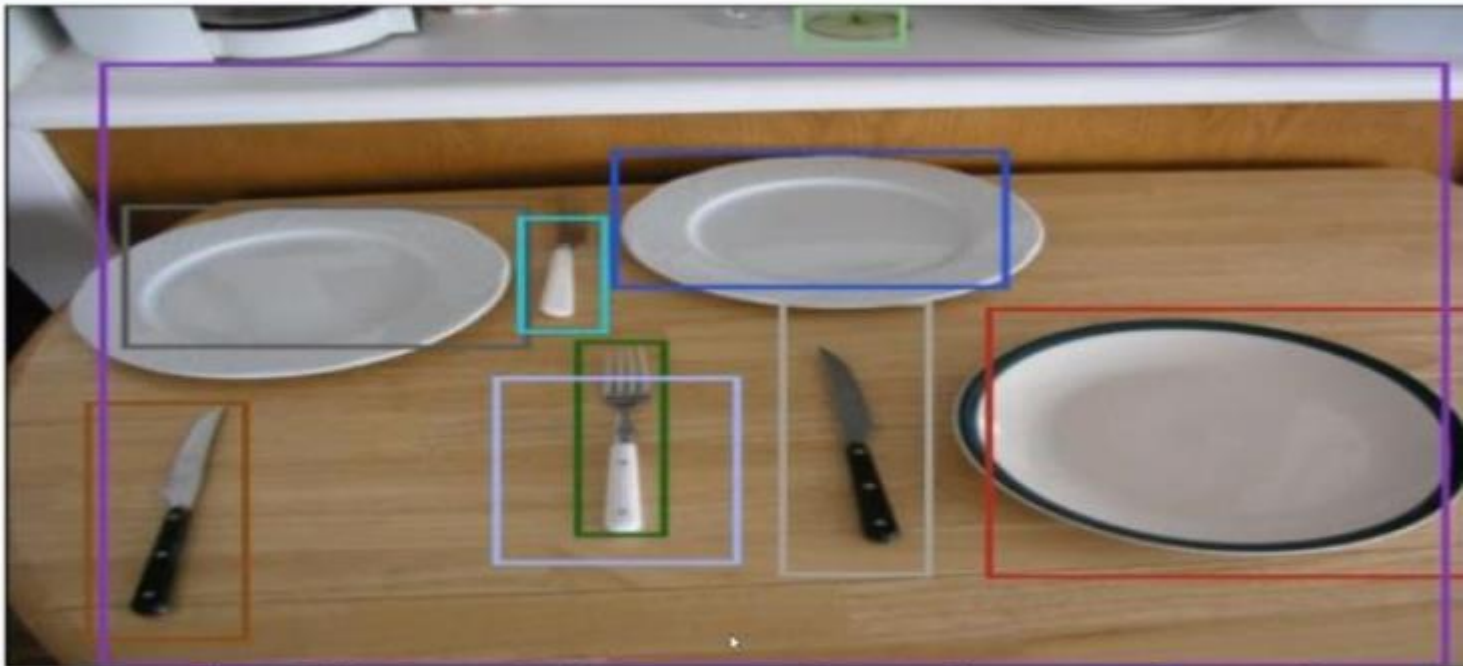
# NLP is Important

Cutting edge applications

# Large Applications to reduce the problem of scale

- (A) Machine Translation (demo)
- (B) Information Extraction
- (C) Sentiment and Emotion Analysis
- Complexity and applicability increases by requirement and introduction of Multilinguality, Multimodality

# Dense Image Captioning



सफेद और नीले रंग की मेज पर. सफेद प्लेट पर सफेद प्लेट।. सफेद प्लेट पर सफेद प्लेट।. सफेद और चांदी के बर्तन।. काला और काला चाकू।. एक लकड़ी की मेज पर है।. काला और काला चाकू।. मैं हरा और हरा <unk>. सफेद और चांदी के साथ एक चाक।. सफेद और सफेद रंग का होता है।

# OCR-MT-TTS

- Input image:
- English transcription: Take the risk or loose the chance
- Hindi Translation: जोखिम लें या मौका गंवा दें।
- Hindi speech



Take the risk  
or  
lose the chance

Demos



# Course Logistics

# Venue

- First Introductory Lecture: F.C.Kohli Auditorium, Thursday 28<sup>th</sup> July, 10.30AM
- After that in new CSE building, 101 (first floor)

# Course Logistics- MS teams

- Join the MS Teams using the code **jun37pk**
  - Login to MS teams using the LDAP credentials
  - Select **Join or create a team**
  - Choose **Join a team with a code**
  - Join using the code above
  - You will be added to the team CS626-2022
    - General channel: For notifications
    - Live lectures channel: For live lectures

# Course website

Website:

[https://www.cse.iitb.ac.in/~cs626/2022/  
#about](https://www.cse.iitb.ac.in/~cs626/2022/#about)

Visit the above website for course related information.

For more information about the research in NLP visit CFILT website.

<https://www.cfilt.iitb.ac.in/>

# Class Schedule

- Monday: 8:30 AM to 9:25 AM
- Tuesday: 9:30 AM to 10:25 AM
- Thursday: 10:30 PM to 11:25 AM

# Moodle

Login to Moodle with LDAP credentials.

- Select the course CS626

All course related notifications will be notified via Moodle also.

# Course TAs

- Nihar Ranjan Sahoo, PhD CSE  
(nihar@cse.iitb.ac.in)
- Sandeep Singamsetty, M.Tech CSE)  
(213050064@iitb.ac.in)
- NVS Abhishek, M.Tech CSE  
(213050019@iitb.ac.in)

# Evaluation Scheme (tentative)

- 50%: Reading, Thinking, Comprehending
  - Quizzes (15)
  - Midsem (15)
  - Endsem (20)
- 50%: Doing things, Hands on
  - Assignments (25%)
  - Reading ONE paper and doing a preliminary implementation of the same (25%)
- Quiz every last Thursday of the month



# Course Content: Task vs. Technique Matrix

Task (row) vs. Technique (col) Matrix	Rules Based/Kn owledge- Based	Classical ML				Deep Learning		
		Perceptron	Logistic Regression	SVM	Graphical Models (HMM, MEMM, CRF)	Dense FF with BP and softmax	RNN- LSTM	CNN
Morphology								
POS								
Chunking								
Parsing								
NER, MWE								
Coref								
WSD								
Machine Translation								
Semantic Role Labeling								
Sentiment								
Question Answering								

# Books

1. Dan Jurafsky and James Martin, Speech and Language Processing, 3<sup>rd</sup> Edition, 2019.
2. Christopher Manning and Heinrich Schutze, Foundations of Statistical Natural Language Processing, MIT Press, 1999.
3. Pushpak Bhattacharyya, Machine Translation, CRC Press, 2017.
4. Ian Goodfellow, Yoshua Bengio and Aaron Courville, Deep Learning, MIT Press, 2016.

# Journals and Conferences

- Journals: Computational Linguistics, Natural Language Engineering, Journal of Machine Learning Research (JMLR), Neural Computation, IEEE Transactions on Neural Networks
- Conferences: ACL, EMNLP, NAACL, EACL, AACL, NeurIPS, ICML

# Useful NLP, ML, DL libraries

- NLTK
- Scikit-Learn
- Pytorch
- Tensorflow (Keras)
- Huggingface
- Spacy
- Stanford Core NLP