

---

# Language Models for Text Recognition: An Overview

---

Shravya Shetty and Sargur Srihari

CEDAR

University at Buffalo

State University of New York, USA

---

# Plan of Presentation

1. Language models
    - N-gram probabilities
    - Generative model for text recognition
    - Improvements to n-gram models
  2. Text Recognition
  3. Sentence Level Language Models
  4. Word Level Language Models
  5. Conclusion
  6. References
-

---

# 1. Language Models

- Language Models are probabilistic models that capture **language regularities**
  - N-gram models:
    - unigram, bigram and trigram
    - **word models**: capture word dependencies
    - **part of speech models**: capture word class dependences
  - Shown to be effective/improve performance in language technologies
-

---

# Applications of Language Models

- Speech Recognition
  - Machine Translation
  - Information Extraction
  - Text recognition (OCR, Handwriting)
    - Post-processing of recognition results
  - Spelling correction
  - MINDS
    - ML(statistical), IR, NLP, DAR, ASR
-

# Language Model for Text Recognition

Minimum error rule of recognizing text image:

$$\underset{\text{word sequences}}{\operatorname{arg\,max}} \{ P(\text{word sequence} \mid \text{text image}) \}$$

- Task is to determine conditional probability for each word sequence
- Alternatively, by Bayes rule

$$\underset{\text{word sequences}}{\operatorname{arg\,max}} \{ P(\text{text image} \mid \text{word sequence}) \times P(\text{word sequence}) \}$$

---

# Generative Model for text recognition

$\underset{\text{word sequences}}{\arg \max} \{ P(\text{text image} | \text{word sequence}) \times P(\text{word sequence}) \}$

- Basically a generative model where joint probability is computed
  - First Term a conditional probability, computed by say Naïve Bayes
  - Second term calculated using the n-gram model (and chain rule)

$$P(\text{word sequence}_1^n) = P(\text{word}_1) \times P(\text{word}_2 | \text{word}_1) \times P(\text{word}_3 | \text{word}_2, \text{word}_1) \times \\ P(\text{word}_4 | \text{word}_3, \text{word}_2, \text{word}_1) \dots P(\text{word}_n | \text{word}_{n-1} \dots \text{word}_1)$$

---

# N-gram Language Model

- Probability of word sequence 1..N is

$$p(w_1^N) = \prod p(w_i/w_{i-1}) \quad \textit{bigram}$$

or

$$\prod p(w_i/w_{i-1}, w_{i-2}) \quad \textit{trigram}$$

- Shown to be effective in natural language applications
-

---

# Limitations of n-gram models

- Sequences not occurring in training data have zero probability
  - Many n-grams occur too few/many times due to corpus
  - Long-term dependencies are there
  - Improvements to n-gram models possible
-

---

# Techniques for Improving n-gram Models

## ■ Smoothing

- Probability mass of n-grams with frequency greater than a threshold is redistributed across all the n-grams

## ■ Clustering

- Predict cluster of similar words instead of single word.

## ■ Caching

- Recently observed words are likely to occur again. Combine with more general models

## ■ Skipping

- words not directly adjacent to target word contain useful information

## ■ Sentence-mixture models

- Model different kinds of sentences separately
-

---

## 2. Text Recognition as a Document Processing Task

---

OCR

ICR

OHR

ILR

# Processing Steps

Scan

Line Segmentation

Word Segmentation



The image displays three sequential screenshots of the CEDAR-FOX software interface, illustrating the processing steps for a scanned document. Each screenshot shows a window titled "CEDAR-FOX - [sample\_hindi.PNG-- Questioned Document]" with a menu bar (File, View, Features, Search, Tools, Recognition, Window, Help) and a toolbar. The main content area shows handwritten Hindi text.

**First Screenshot (Scan):** Shows the original scanned image of the text. The text is: स्तुति उत्तर स्वयंभू दित्य आ, हिन्दू आगन्तुक राष्ट्रीय व्य, प्रश्नन्न र्शैत्र्य शून्य -द्व, स्त्री स्कूल अच्चाईयाँ अ.

**Second Screenshot (Line Segmentation):** Shows the text with horizontal lines drawn across it, indicating the segmentation of the document into lines. The text is: 1 स्तुति उत्तर स्वयंभू दित्य आत्मा, 2 हिन्दू आगन्तुक राष्ट्रीय व्यक्ति, 3 प्रश्नन्न र्शैत्र्य शून्य -द्वीय, 4 स्त्री स्कूल अच्चाईयाँ आश्व

**Third Screenshot (Word Segmentation):** Shows the text with individual words highlighted in different colors, indicating the segmentation of the document into words. The text is: स्तुति उत्तर स्वयंभू दित्य आत्मा, हिन्दू आगन्तुक राष्ट्रीय व्यक्ति बुल्लू, प्रश्नन्न र्शैत्र्य शून्य -द्वीय पर्याप्त, स्त्री स्कूल अच्चाईयाँ आश्वान भर्ष

# Isolated Word Recognition

## Task: Convert Word Image to its Textual Form

- **Analytic Recognition**
  - ❑ Segments of word image matched to characters of words in lexicon
- **Holistic Recognition** →
  - ❑ Word shape matched to prototypes or features of words in lexicon.
- **Combining Results**
  - ❑ Neural Network

The Query Image

पातु

Search In

Current Document

The Other Document

Preprocessed Documents

Search only in ROI

Searching results

Rank	Word Image	Distance	Document ID	
No.1	पातु	0	5_1	Highlight
No.2	पातु	0	5_1	Highlight
No.3	पातु	0	5_1	Highlight
No.4	पातु	0	6_2	Highlight

Go

Goto Top

Go Up

---

# Language Models in Text Recognition

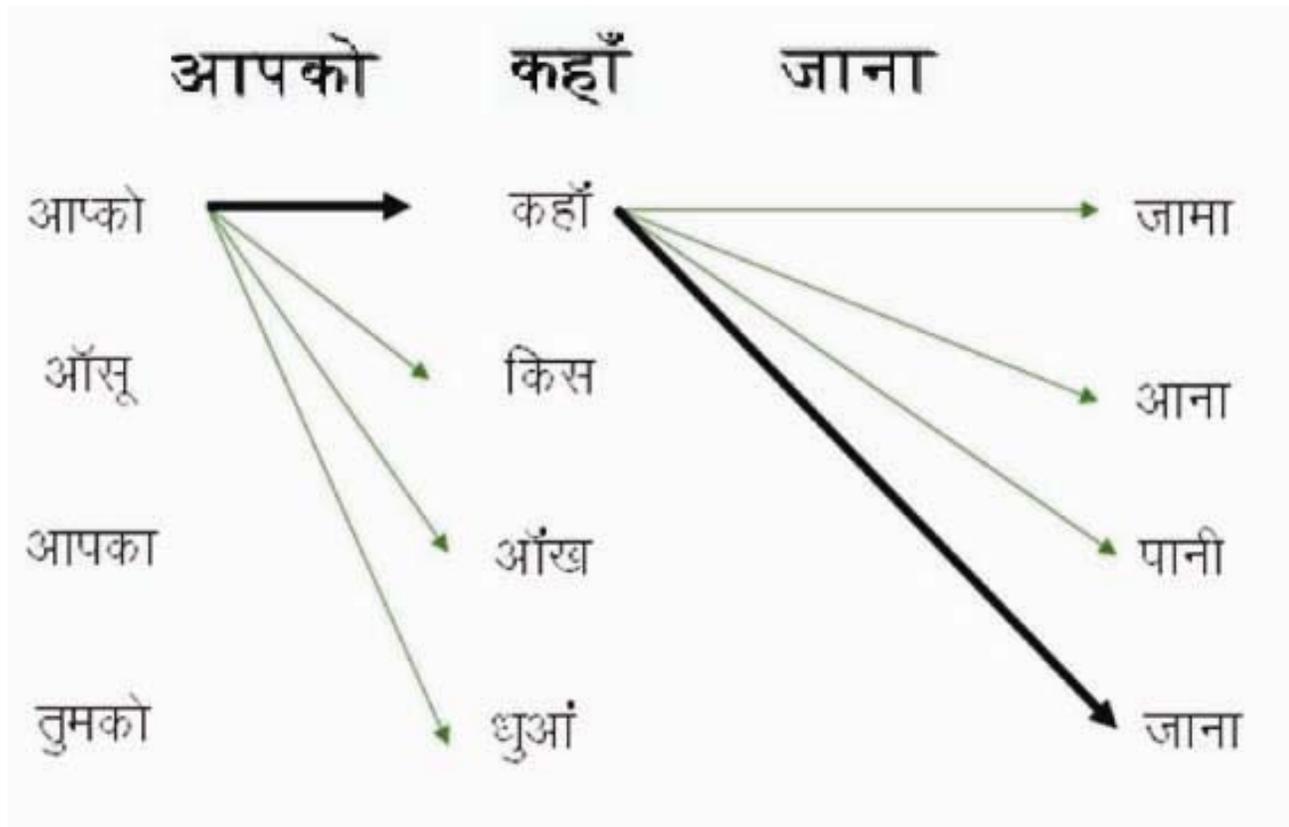
- Language models can be at
    - inter-word level (sentence) or
    - inter-character (word) level
  - Sentence level
    - language model learnt from a corpus consisting of sentences
  - Word level
    - language model learnt using dictionaries
-

---

# 3. Sentence Level Language Models

---

# Recognition Post-processing: Finding most likely word sequence using word level models



Word n-grams  
Word-class n-grams  
(POS, NE)

To make  
recognition  
choices  
or to limit choices

# Language model for correcting recognition results

- Can improve performance when sentence follows average statistics of the language
- Trigram language model
- Smoothing using Interpolated Kneyser-Ney
  - Modified backoff distribution based on no of contexts
  - Higher- and lower-order distributions are combined

- An implementation for handwritten essays:
  - N best list of word recognition results are used
  - Second order HMM is used to incorporate trigram model
  - Find most likely sequence of hidden states given a sequence of observed paths in a second order HMM– Viterbi Path

$$P_{\text{IKN}}(w_i | w_{i-1}) = \frac{C(w_{i-1}w_i) - D}{C(w_{i-1})} + \lambda(w_{i-1}) \frac{|\{v | C(vw_i) > 0\}|}{\sum_w |\{v | C(vw) > 0\}|}$$

# Results with Language Modeling (Handwritten Essays)

## ORIGINAL TEXT

Lady Washington role was hostess for the nation. It's different because Lady Washington was speaking for the nation and Anna Roosevelt was only speaking **for the people** she ran into **on wer travet** to see the president.

Four  
Essays

204  
(100%)

## WORD RECOGNITION

Lady washingtons role was hostess for the nation first to different because lady washingtons was speeches for for martha and taylor roosevelt was only meetings **for did people** first vote polio **on her because** to see the president

124  
(62%)

## LANGUAGE MODELING

Lady washingtons role was hostess for the nation but is different because george washingtons was different for the nation and eleanor roosevelt was only everything **for the people** first ladies late **on her travel** to see the president

145  
(72%)

---

# Correction using Word Class

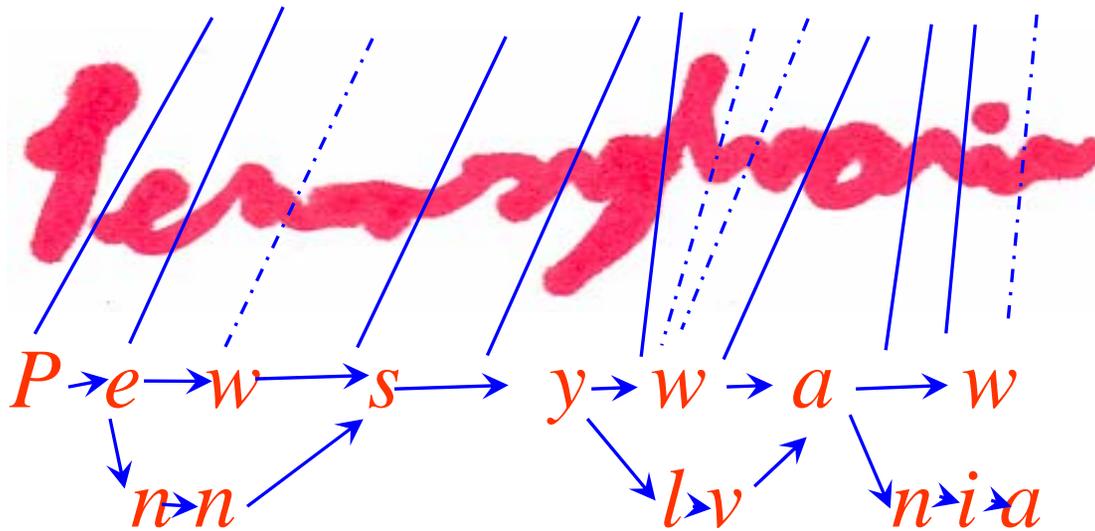
- Words in language model are replaced by corresponding POS and NE tags
  - Bigram probabilities for word classes learnt from corpus
  - Word correction performed using Viterbi decoding
  - Slightly improved performance
-

---

# 4. *Word Level Language Models*

---

# Word level models for Latin alphabet



- For each character segment, train recognizer (NN, SVM, etc.) to recognize letter
- Learn  $P(\text{letter} \mid \text{character image})$
- Include the language model
- $P(\text{letter} \mid \text{character image}) \times P(\text{letter} \mid \text{previous letters})$

- Recognition and language model are tightly integrated
- Word recognition uses continuous density HMMs
- Search space is a network of word HMM's and word transition modeled by a language model

# Character and Alphabet Models in Indic Languages

- n-grams in a character can be at the character level or the alphabet level
- Example of **character** level bigram

$$p(x \text{ पूछती } x) = p(\text{पू} | x) \cdot p(\text{छ} | \text{पू}) \dots$$

- Example of **alphabet** level bigram

$$p(x \text{ पूछती } x) = p(\text{प} | x) \cdot p(\text{८} | \text{प}) \cdot p(\text{छ} | \text{८}) \dots$$

# Character/Alphabet N-grams in Devanagari

- 126 alphabets and 5538 characters
- Bigram and trigram frequency counts can be obtained for them. E.g., the Emille corpus

Tokens	Number of classes	Perplexity ( $Pp$ )	
		Bigram	Trigram
Alphabet	126	13.64	6.23
Character	5538	18.51	3.53

- N-grams at character level give lower perplexity (less uncertainty)

# Word Level Language Model for Devanagari



र → मा → य → ण

श → गा → रा ↗

- Determine character path
- Recognize characters
- Use Language model

# Syntactic Rules for Devanagari Words

- Formation of characters determined by rules
- Can be used to reject character sequences
- Example rules –
  - Half-consonant cannot have a vowel modifier

Illegal ----- ङॆ

- Only one vowel modifier is allowed on a consonant

झ । ॆ ।      छ । ॆ ।  
First vowel modifier      Second vowel modifier

---

# 5. Conclusion

- Statistical language models found useful in all language technologies
    - Can reduce error rate by  $1/3$  or  $1/2$
  - Tight coupling of HMMs for text recognition and n-gram language models possible
  - Indic Languages have special challenges
    - Character and Alphabet sets
  - More structure in Indic Languages
    - POS models (vibhakti) should work well for Sanskrit
    - Hybrid models are useful
  - References
-

---

# References

- i. **J. Hull**, *A Computational Theory of Visual Word Recog.*, PhD dissert, CEDAR 1988.
  - ii. **R. Srihari, C. Ng, C. Baltus**, *Language Models in On-line Handwriting Recognition*, 3<sup>rd</sup> IWFHR, 1993
  - iii. **G. Kim**, *Handwritten Word and Phrase Recognition*, PhD Dissert, CEDAR1996.
  - iv. **A. Vinciarelli, S. Bengio H. Bunke**, *Off-line recognition of unconstrained Handwritten sentences using HMM and Statistical Language Models*, 2003.
  - v. **S. Srihari, H. Srinivasan, C. Huang, S. Shetty**, *Spotting words in Latin Devanagari and Arabic Scripts*, Vivek, 2006.
  - vi. **S. Srihari, R. Srihari, H. Srinivasan, S. Shetty**, “On scoring handwritten essays”, IJCAI 2007.
  - vii. **S. Kompalli**, *A Stochastic Framework for Font-Independent Devanagari OCR*, PhD dissert, CEDAR 2007.
-