

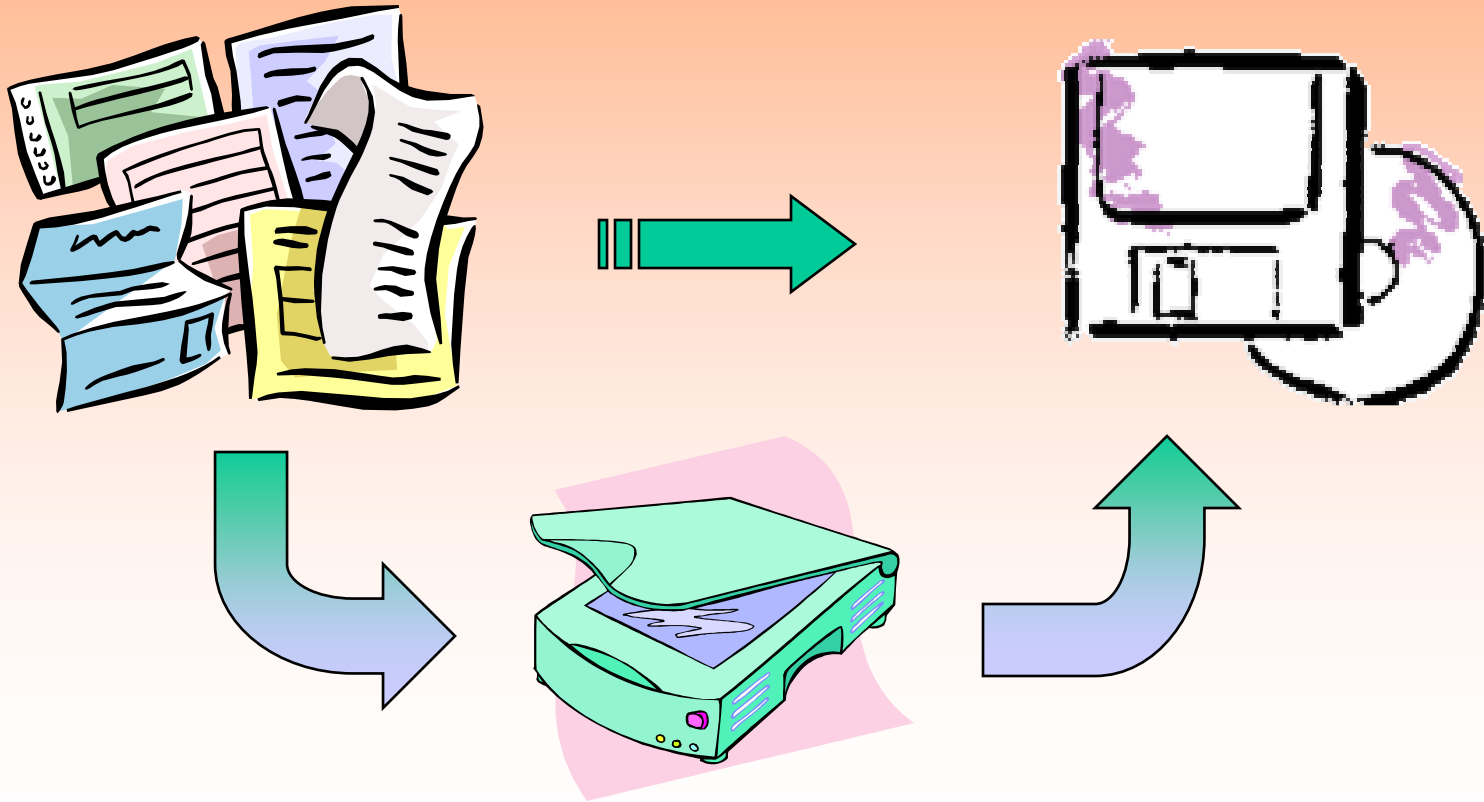
***A Voting System for
Automatic Correction of OCR
Output***

Tomi Klein Miri Kopel

Bar Ilan University

Introduction

OCR = Optical Character Recognition



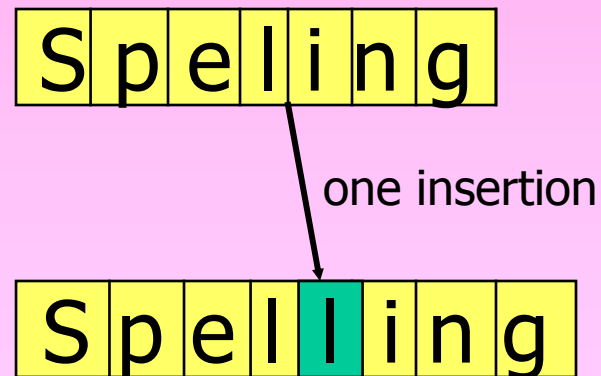
Problems:

- **Even the most accurate OCR devices don't give 100% accuracy.**
- **Manual correction is too expensive**

Known Techniques for Spelling Correction

Edit Distance:

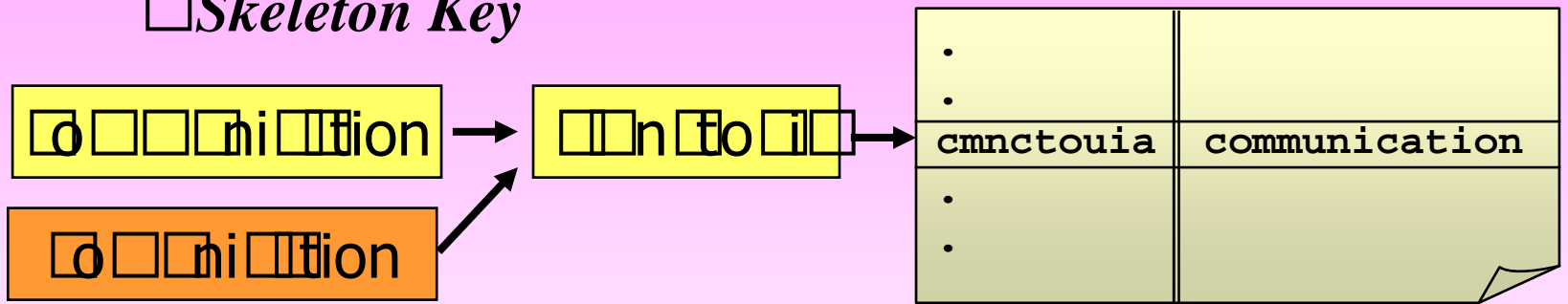
The minimum number of editing operations (i.e., insertion, deletion and substitution of letters) required to transform one string to another.



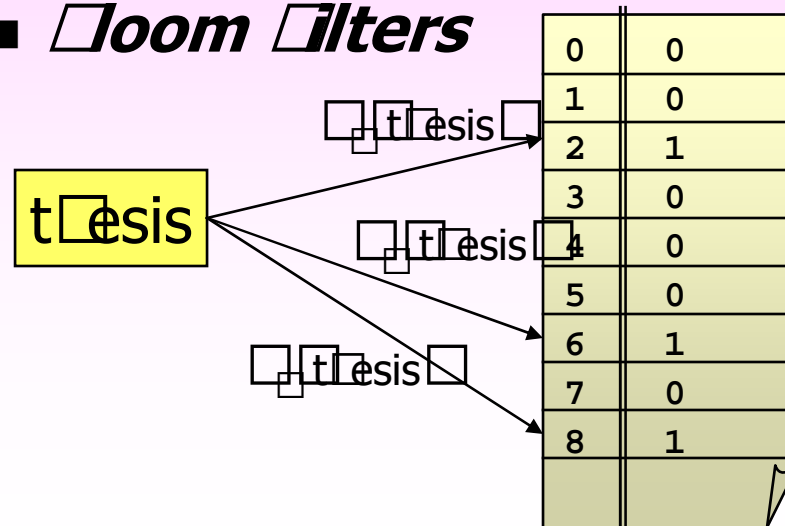
Known Techniques for Spelling Correction (cont...)

- Hashing:

- *Skeleton Key*



- *loom filters*

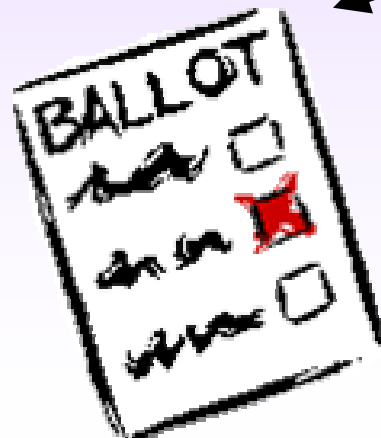





So, what is the problem ?

- **Most of the techniques are relevant only for typing errors.**
- **Designed for isolated words.**
- **We are interested in a fully automatic system.**



The Algorithm



- 
its other typical errors occur
- 
For example, the confusion matrix
- 
consider the following

The Algorithm (cont...)

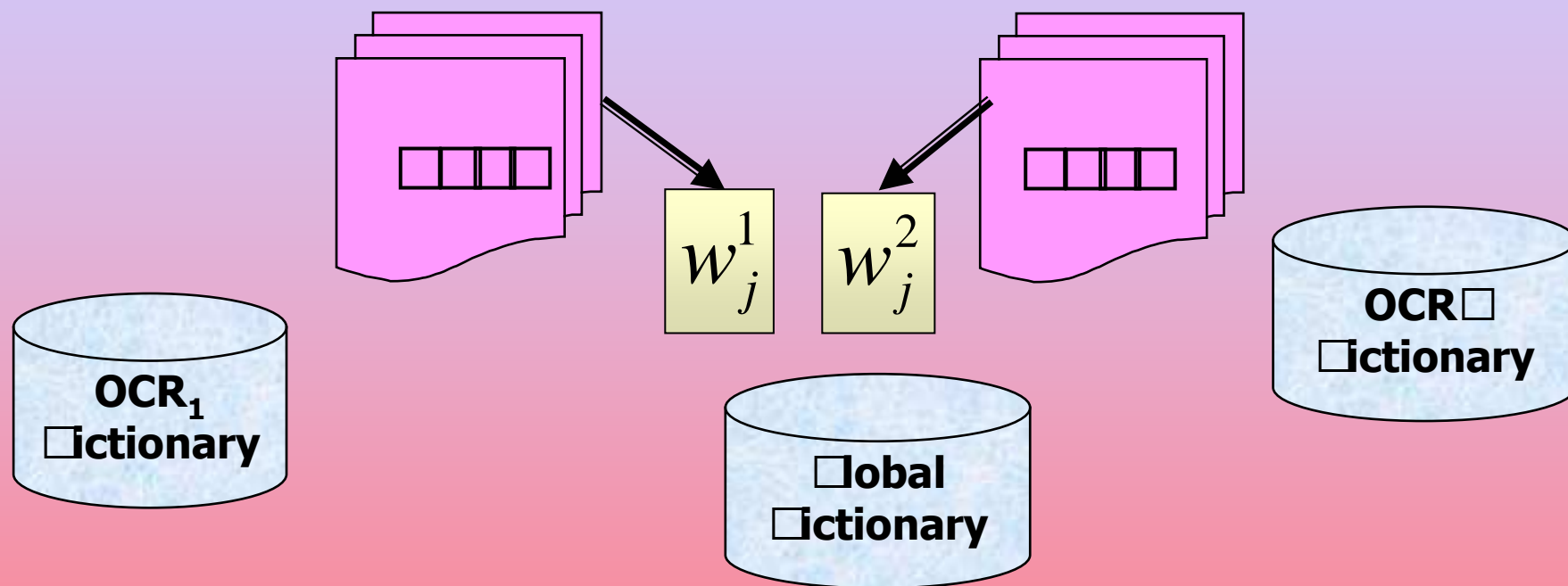
Initialization Phase

- Build the dictionary.
- Create the confusion matrices.
- Construct the n -gram and word-gram tables.
- Build the documents sets.



The Algorithm (cont...)

Detection and Correction Phase



Word 1 Word 2 Word 3 Word 4

Word 1 Word 2 Word 3 Word 4 Word 5 Word 6 Word 7 Word 8 Word 9 Word 10

Word 1 Word 2 Word 3 Word 4 Word 5 Word 6 Word 7 Word 8 Word 9 Word 10

Word 1 Word 2 Word 3 Word 4 Word 5 Word 6 Word 7 Word 8 Word 9 Word 10

The Algorithm (cont...)

- If words are identical:
Accept the word.

$$w_j^1 = w_j^2$$

- **It's important**

$$w_j^1 \neq w_j^2$$

- **ictionaries**

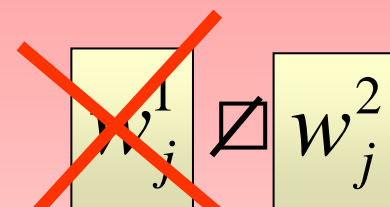
- **is a letter**

- **on'sion**

- **onte**

The Algorithm (cont...)

If the words mismatch:



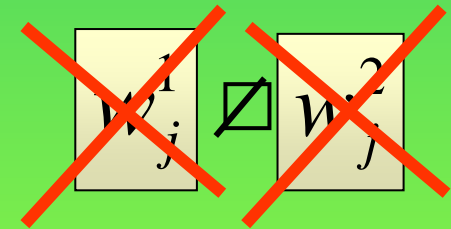
- If only one of them is valid Then:
Accept the valid one.

$$\text{dictionary}(w_j^i) = 0.6 \cdot \text{freq}(w_j^i, \text{local_dictionary}) + 0.4 \cdot \text{freq}(w_j^i, \text{global_dictionary})$$

The Algorithm (cont...)

- Not re in en

Generate candidates or not
in en except the best one



$$\text{Candidates} = \text{Generate_Candidates}(w_j^1) \cup \text{Generate_Candidates}(w_j^2)$$

The Algorithm (cont...)

Parameters for running the Markov chain

Transition probabilities

Initial state

Order of Markov chain

Transition matrix

$$is_close = \begin{cases} 1 & \text{if } edit_distance(w_j^i, w_{j_k}^i) = 1. \\ 0 & \text{otherwise.} \end{cases}$$

$$mark(w_{j_k}^i) = 0.6 \cdot (error_freq(w_j^i, w_{j_k}^i, OCR^i)) + \\ 0.4 \cdot (word_gram(w_{j-1}^i, w_{j_k}^i, w_{j+1}^i) + \\ dictionary(w_{j_k}^i) + \\ is_close(w_j^i, w_{j_k}^i))$$

The Algorithm (cont...)

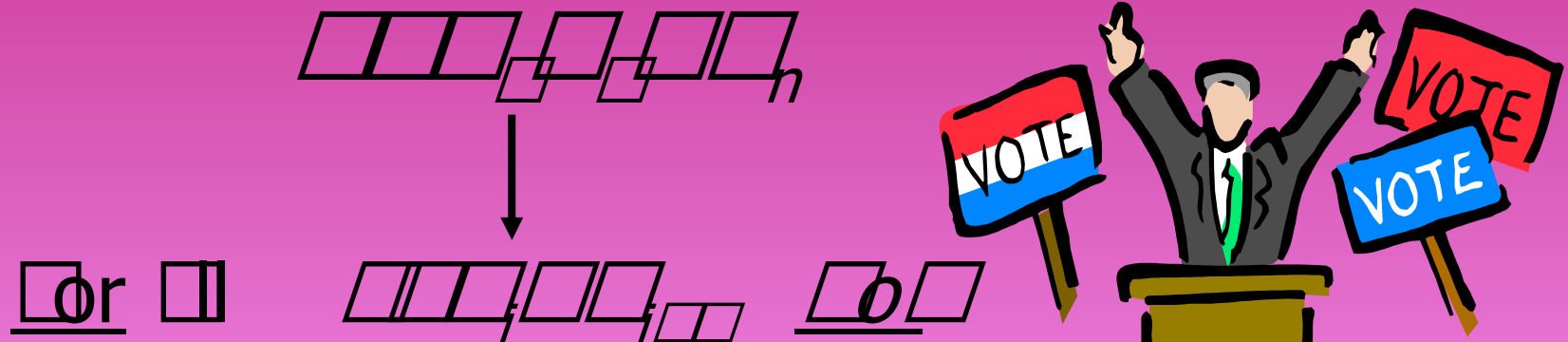
- **If both words are valid:**

$$\text{context}(w_j^i) = \text{freq}(w_j^i, \text{local_dictionary}) + \text{word_gram}(w_{j-1}^i, w_j^i, w_{j+1}^i).$$

$$\text{mark}(w_j^i) = \text{accuracy}(\text{OCR}^i) \cdot [0.6 \cdot \text{context}(w_j^i) + 0.4 \cdot \text{freq}(w_j^i, \text{global_dictionary})]$$

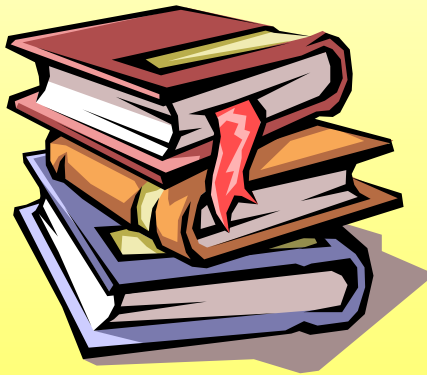
The Algorithm (cont...)

Generate Candidates(w)



- The condition $w_i \neq w_{i+1}$
- There are no rules applied
- The generated word is added to the candidates set

The Experiments



□□si□□□□□rde□□ai□re□er

□□□□an□et□□□□C

The Environment of the Experiments

The confusion matrix for English text

| <i>Source String</i> | <i>Error String</i> | <i>OCR 1</i> | <i>OCR 2</i> |
|----------------------|---------------------|--------------|--------------|
| i | l | 1.58% | 0.48% |
| 0 | o | 3.54% | 0.77% |
| da | d | 0.75% | 0.10% |
| 2 | z | 0.30% | 0.72% |
| f | l | 0.15% | 0.72% |
| i | l | 1.58% | 0.48% |
| e | g | 1.18% | 0.10% |

Examples of successful corrections

English and French texts

| <i>Original Word</i> | <i>OCR 1</i> | <i>OCR 2</i> | <i>Accepted Word</i> |
|----------------------|--------------|--------------|----------------------|
| going | going | goring | going |
| survivors | surveyors | survivors | survivors |
| we're | we'te | we'e | we're |
| thankfully | thankfvly | thankfiily | thankfully |
| school | school | sciiool | school |
| neighborhood | ne~hborhood | nei&iborho~ | neighborhood |
| precisent | pr,cisent | preci~ent | precisent |
| details | d,tails | detaHs | details |

Examples of erroneous decisions

Errors can be generated for several reasons:

1. We accept an incorrect word directly from the OCR for e.g. when both OCRs generate the same mistake.
 - We choose the wrong word between two valid alternatives.
 - Only one of the words is in the dictionary but it is not the correct one.
 - We corrected the input word which was correct to an incorrect one.
 - We couldn't find a proper candidate for the misspelled input word.

| <i>Original Word</i> | <i>OCR 1</i> | <i>OCR 2</i> | <i>Accepted Word</i> | <i>Error type</i> |
|----------------------|--------------|--------------|----------------------|-------------------|
| 5,600 | 5,60o | 5,60o | 5,60o | 1 |
| Ian | lan | lan | lan | 1 |
| class | dass | dass | dass | 1 |
| 100 | 100 | too | too | 2 |
| true | true | toe | toe | 2 |
| little | lithe | lisle | lithe | 2 |
| We | the | We | the | 2 |
| circles | circ.les | circus | circus | 3 |
| says | savs | sas | san | 4 |
| 2001 | 2ooi | zooi | zoo | 4 |

Analysis of the Results

Model 1:
OCR1, dictionary
with no post
processing.



Model 2:
OCR2, dictionary
with no post
processing.



Model 3:
OCR1
post processing



Model 4:
OCR2
post
processing



Analysis of the Results (cont...)

Model 1 :

OCR 1 : OCR 2:
simple dictionary
loop.



Model 2 :

OCR 1 : OCR 2:
dictionary loop
with
proabilities.



Model 3 :

OCR 1 : OCR 2:
post
processing.



Analysis of the Results (cont...)

| | | Error Rate |
|---|---|------------|
| 1 | OCR 1 – no post-processing | 14.0% |
| 2 | OCR 2 – no post-processing | 12.1% |
| 3 | OCR 1 – dictionary + candidates generation | 8.5% |
| 4 | OCR 2 – dictionary + candidates generation | 8.0% |
| 5 | Both – comparison + simple dictionary lookup | 7.2% |
| 6 | Both – comparison + dictionary's frequencies | 5.1% |
| 7 | Both – full correction system | 3.6% |

Analysis of the Results (cont...)

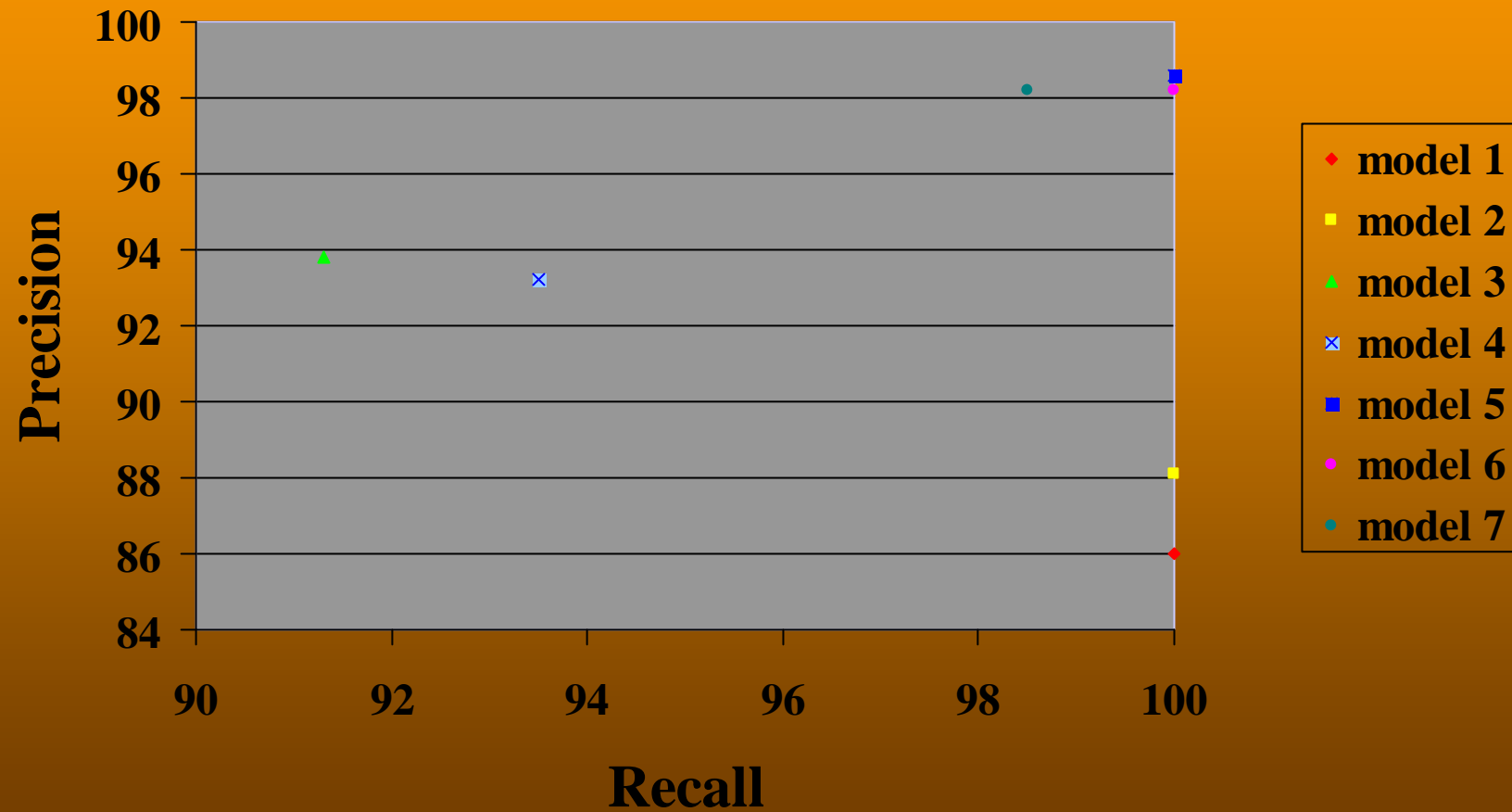
| | Accept words from OCR as is | Do not accept the OCR word, and try to suggest candidates |
|----------------------------------|-----------------------------|---|
| Correct word written to output | (A) true positive | (B) true negative |
| Incorrect word written to output | (C) false positive | (D) false negative |

$$\text{Recall} = \frac{A}{A + B} = \frac{\text{\# correct words accepted directly from OCR}}{\text{\# correct words accepted}}$$

$$\text{Precision} = \frac{A}{A + C} = \frac{\text{\# correct words accepted directly from OCR}}{\text{\# total words accepted directly from OCR}}$$

Analysis of the Results (cont...)

Recall vs. Precision



Further Work

- More OCR devices.
- Context: NLP techniques, word class
- Specifications for certain language.