# OCR-Based Document Classification

## Gillian Pan
*Mark Greene, Dr. Yash Shrivastava*
School of Electrical and Information Engineering
FACULTY OF ENGINEERING & INFORMATION TECHNOLOGIES

## DOCUMENT CLASSIFICATION

- Document classification is the sorting of documents into predefined categories.
- The trend of the paperless office in the recent years prompted a large growth in more efficient document processing and retrieval, driven largely by the incentive of reducing labor costs.
- Document classification has been studied previously to support the classification of wartime archival documents, medical journals, identity documents, forms and financial documents
- Toshiba hopes to integrate document classification into their e-BRIDGE Re-Rite product to enable workflow automation.

## OPTICAL CHARACTER RECOGNITION (OCR)

- Extracts text , blocks from scanned images
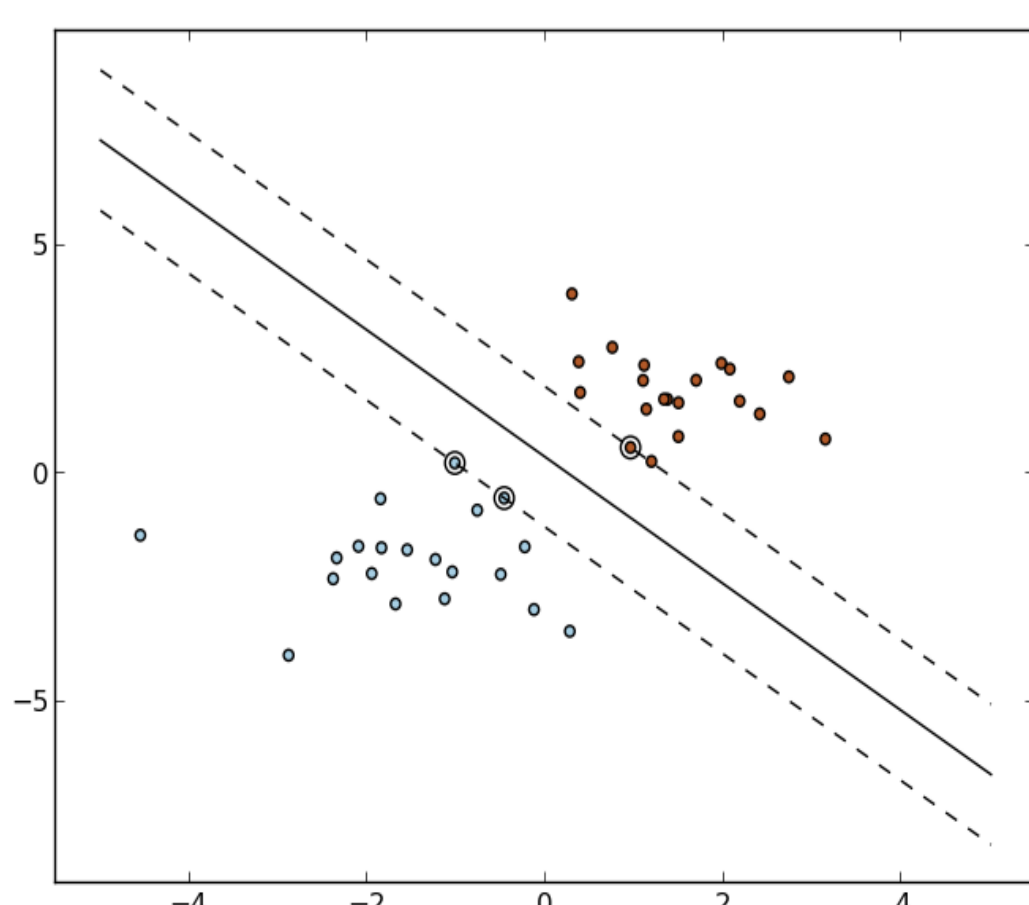- Segments the Text in region blocks

## MACHINE LEARNING

- A branch of Artificial Intelligence, Machine Learning is good at making predictions with large volumes of data.
- Also used for speech recognition, handwriting recognition and spam filtering

### Support Vector Machine (SVM)

- A support vector machine constructs a hyper-plane or set of hyper-planes in a high or infinite dimensional space and finds a separating hyperplane with the maximum margin i.e. the largest distance to the nearest training data points of any class

## CLASSIFIER FEATURES

### Preprocessing
- Levenshtein Distance, Thesaurus

### Semantic Features
- Keyword Exist, Number Words, Number Numeric Words

### Visual Features
- Number of blocks of each type, Total area taken by each block type, Block Widths
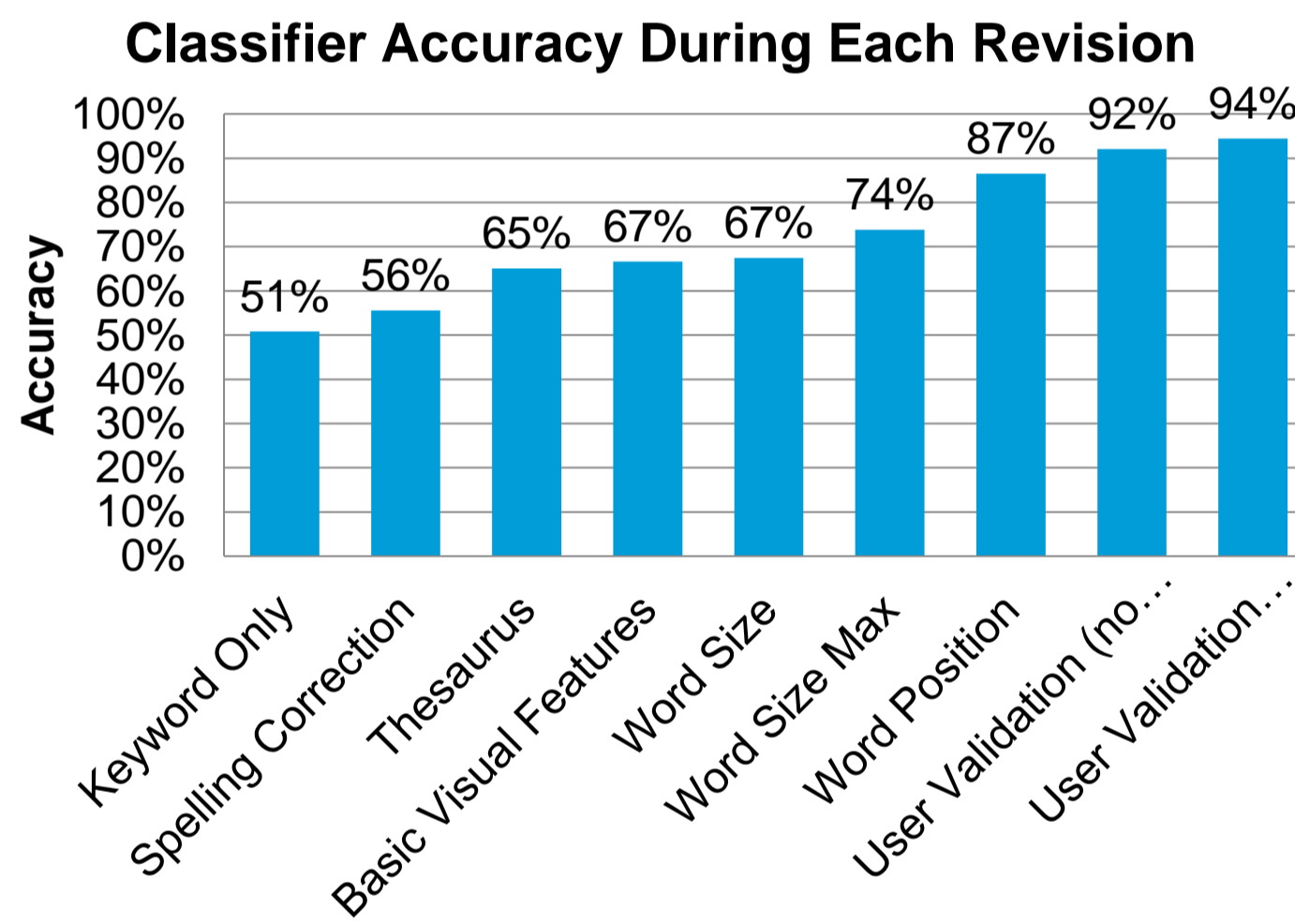
### Combination Semantic and Visual Features
- Word Size, Word Position

### User Validation
- Keyword Search , Prediction Probabilities

## ACCURACY

- Below shows the graph showing the amount of improvement gained at each revision of the classifier. Each revision trained the model containing an additional feature
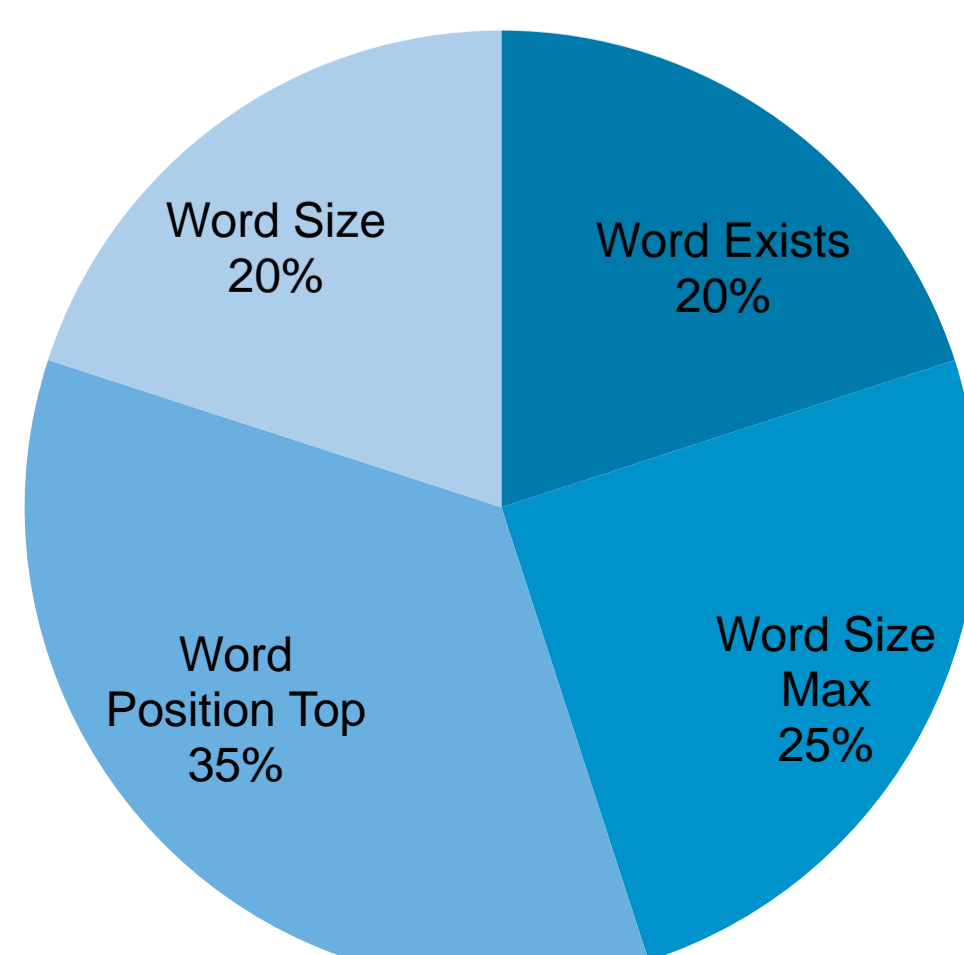
**Classifier Accuracy During Each Revision**

| Revision | Accuracy |
| --- | --- |
| Keyword Only | 51% |
| Spelling Correction | 56% |
| Thesaurus | 65% |
| Basic Visual Features | 67% |
| Word Size | 67% |
| Word Size Max | 74% |
| Word Position | 87% |
| User Validation (no…) | 92% |
| User Validation… | 94% |

### Accuracy Improvement

| Feature Name | % Improvement |
| --- | --- |
| Word Size | 0.79% |
| Basic Visual Features | 1.59% |
| User Validation (probabilities) | 2.38% |
| Spelling Correction | 4.76% |
| User Validation (no keyword) | 5.56% |
| Word Size Max | 6.35% |
| Thesaurus | 9.52% |
| Word Position | 12.70% |

## FEATURE IMPORTANCE
**20 Most important Features**

Word Size 20%
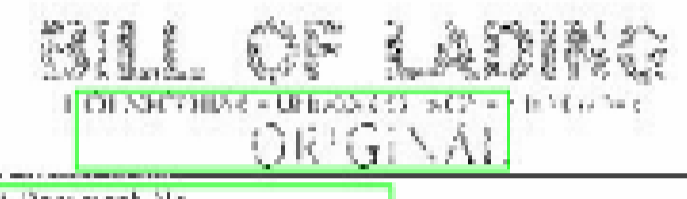Word Exists 20%
Word Position Top 35%
Word Size Max 25%

## THINGS TO LOOK OUT FOR
**Multiple words in document.**

### OCR Errors
- Blocks are not always identified by the OCR

- Below words were identified as an image.

**Purchase Order**

### More OCR Errors (Spelling)

| Words (Levenshtein Distance < 3) | |
| --- | --- |
| Incorrect spelled instances of the word 'invoice' | Nvoiceno, iioic, invuiife, voce, iioic, rvoices, nvoi, Invoc, invoceb, invocf, invoco, jlnvoice, jlnvoice, nwoice, tnvotce, Invoicc, Invoice, mnvoice, unvoice, nvoice |
| Total | 20 |

### Form Errors
- Forms that do not follow the normal conventions of the form type. E.g. Packing slips does not always contain the words 'packing slip' in it.

## GENERALIZATION
- Tries to find out how well the model might perform in practice and attempts to determine if there is any bias caused by any unlucky splits of the available data.

### K-Folds Cross Validation

| Test ID | Predictions Correct | Errors | Accuracy |
| --- | --- | --- | --- |
| 1 | 121 | 4 | 96.80% |
| 2 | 122 | 3 | 97.60% |
| 3 | 121 | 4 | 97.58% |
| 4 | 119 | 6 | 95.20% |
| 5 | 116 | 9 | 92.80% |
| Average | 119.8 | 5.2 | 96.00% |

## CONCLUSIONS
- Accuracy of the classification and generalization results were consistently high.
- The features that the classifier considered important were reflective of how humans classify documents.
- A large accuracy improvement was seen from fixing OCR errors.
- The amount of user interaction needed has to be considered to integrate the technologies described here into Toshiba's e-BRIDGE Re-Rite.