

Learning Document Structure for Retrieval and Classification

Jayant Kumar, Peng Ye, David Doermann

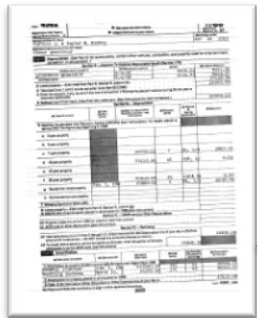
University of Maryland College Park

{jayant, pengye, doermann}@umiacs.umd.edu

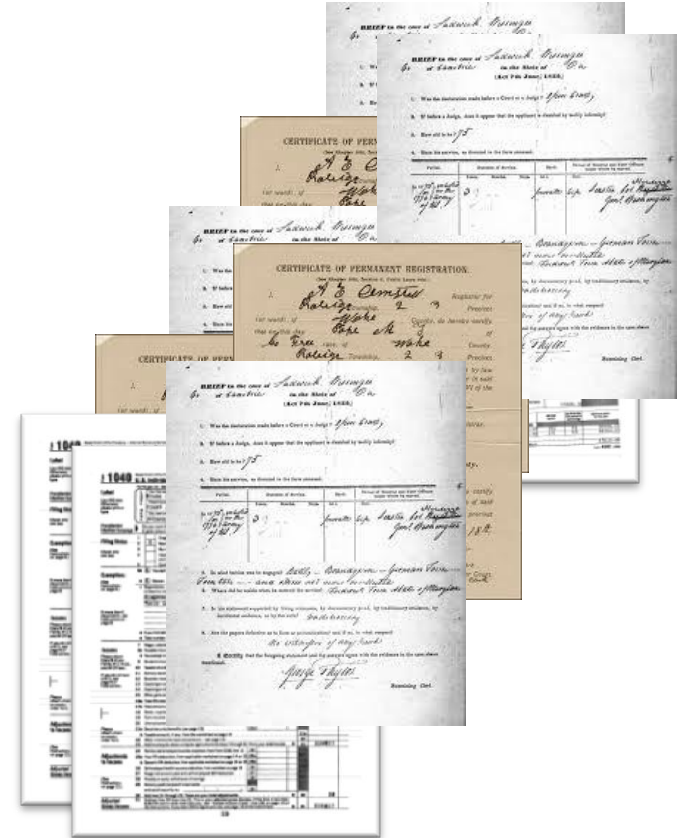
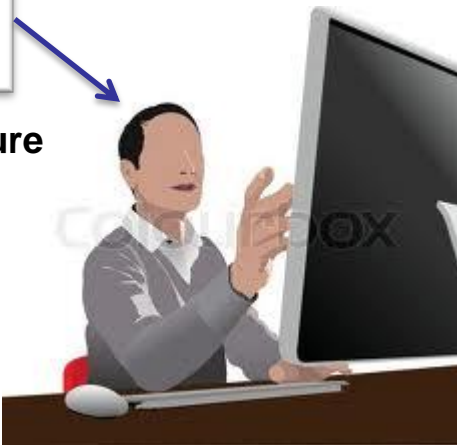


Motivation

- Large Scale Document Image Search
 - By Genre – Preexisting layouts
 - By Example – Similar to “what we have”
 - By User Defined Characteristics



Layout /structure
of interest



Large heterogeneous collection

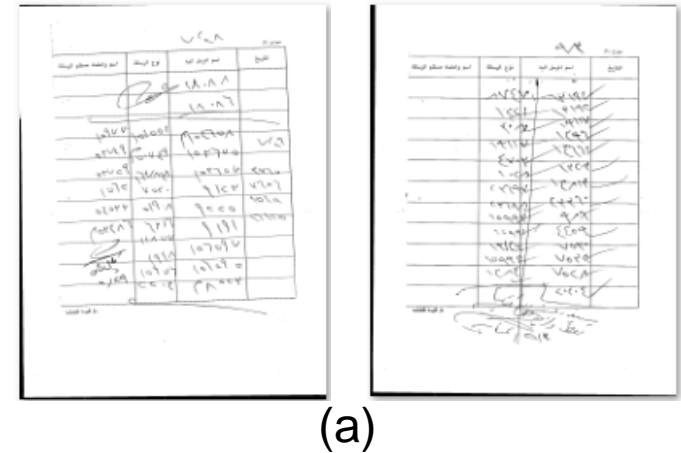
Structural Similarity based Retrieval

Problem:

- Retrieve “similar” documents from a large heterogeneous collection of document images.

Challenges:

- Inconsistent layout
 - Exhibit only similar high-level structure.
- Imbalanced data
 - The number of relevant documents for training may be limited



Relevant Documents

Degree of Structural Similarity

- **Exact Match:** Same underlying structure with some rotation/translation,
 - e.g., Tax forms
- **Approximate Match:** Global structure looks similar with local variations,
 - e.g. handwritten drawn tables cell properties vary, but table looks similar
- **Conceptual Match:** Only at a very abstract level can documents be described as similar.
 - e.g., forms with machine printed headers and handwritten answers

Limitations of Previous Methods

- Layout-specific or Content-specific [Business letters (Dengel 1993)]
- Strict assumptions on layout [Layout similarity (Shin 2001, Diligenti 2003)]
- Disregard spatial relationships [Bag-of-words model (Qiu 2002)]
- Poor performance with limited training data
- Highly sensitive to noise and degradations
- Incapable of finding *important* regions in relevant images

Observations

- Structure relevant at many levels, requiring strong local features
 - Use codebook of local structural patterns and SURF features
- Document elements typically have a horizontal or vertical bias
 - Pool features locally to capture structure
- Relevant properties can be local and a minority in the document
 - Learn which partitions are important!

المستوى	مشروع القرش	التبوع	الاشترالك	المجموع
الحزبي	دينار	فلس	دينار	فلس

(a) Table heading

١٩٠
المسكرة وآخر رجعة محل فيها : (بجانب)

(b) Text lines



(c) Borders design

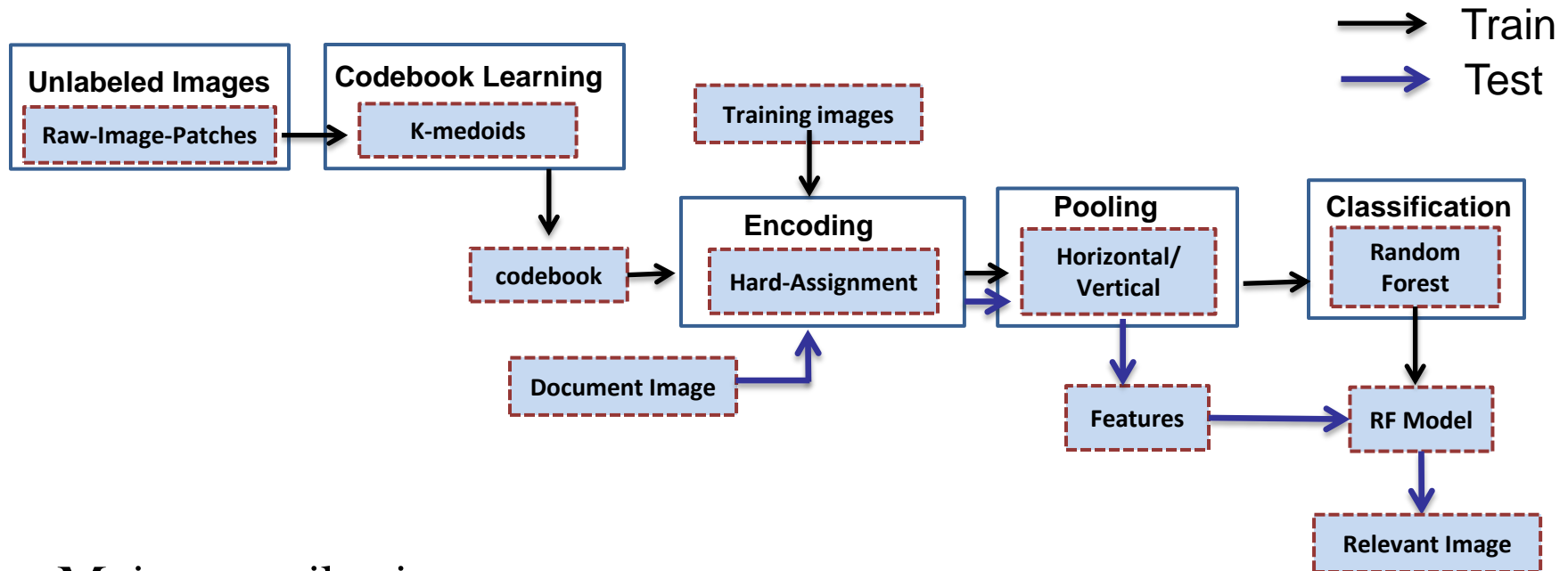
Student's Name	Birth Date	Sex	School	Grade Level / ID#
1516 Washington, Urbana, 01/01/1988	11/23/88	F	Edison	72
1516 Washington, Urbana, 01/01/1988	11/23/88	F	Edison	72

(d) Form with Rule-Lines

Horizontal-vertical
pooling



Proposed Method



Main contributions:

- Recursive horizontal-vertical partitioning for structural-similarity feature computations
- Random-forest based variable importance measures for important structural pattern finding

Codebook based Features

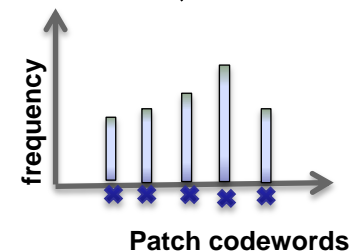
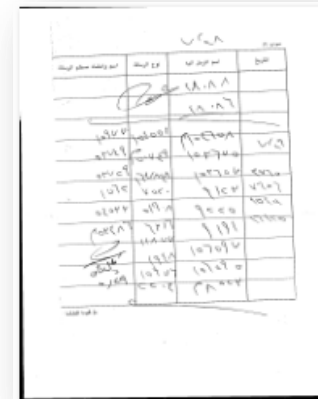


Unlabeled Images

K-medoids

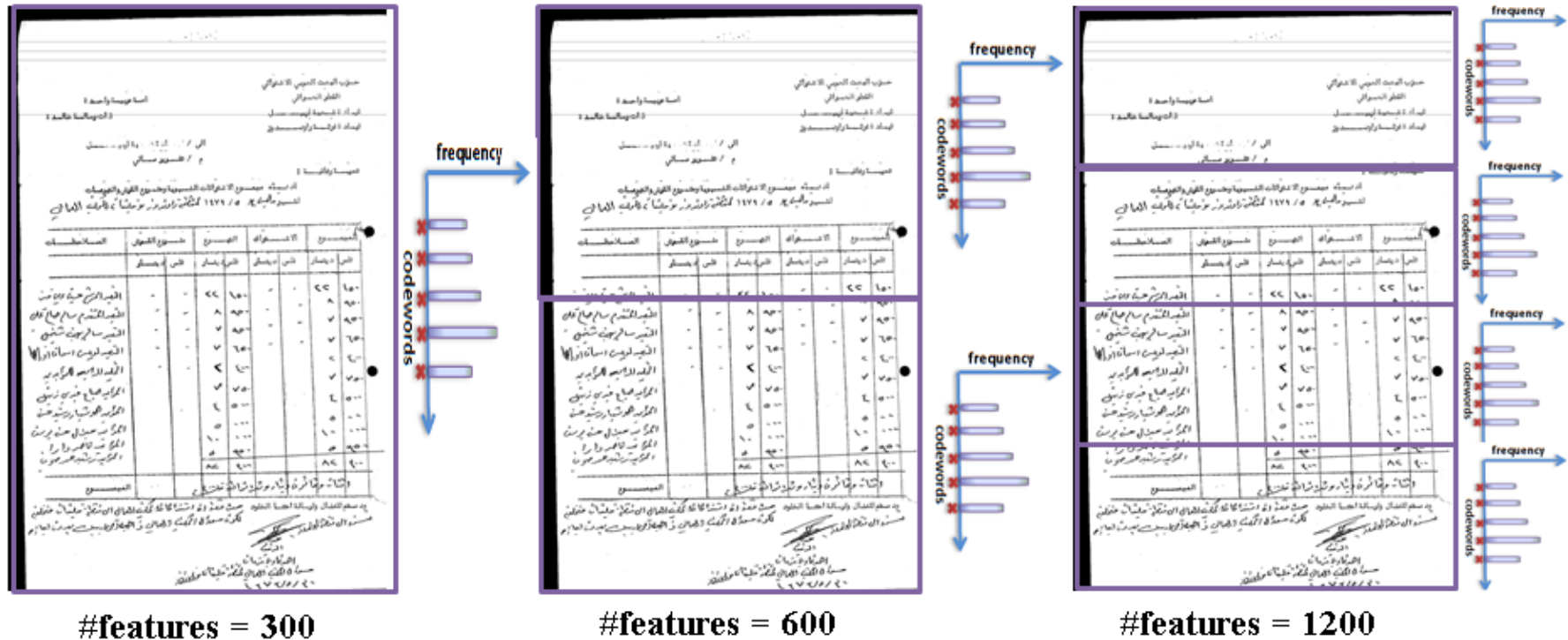


Raw-image-patch codebook based features



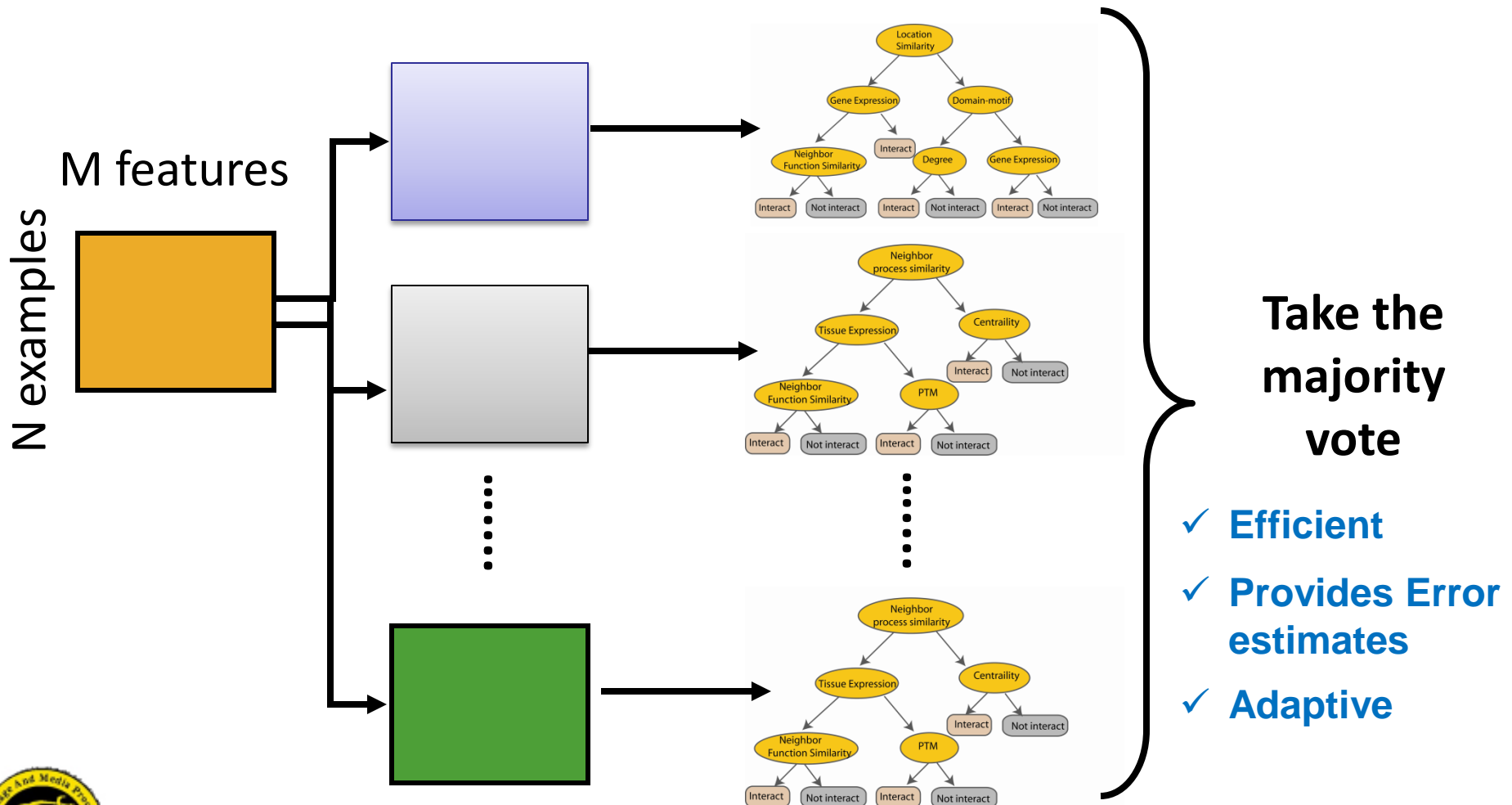
- Very efficient
- Captures local structures

Horizontal-Vertical Pooling of Features



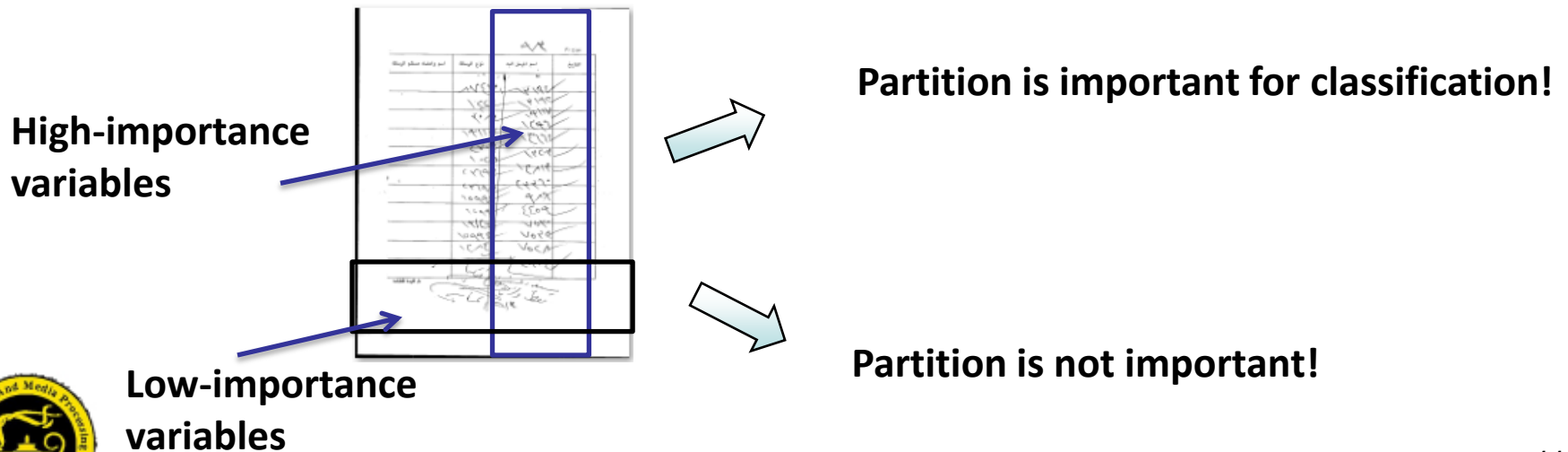
- ✓ Each local descriptor (histogram) characterizes local structure statistics
- ✓ Local histograms are concatenated to form final feature vector for each image

Random Forest Classifier



Adaptiveness Property

- Values of a particular variable are permuted in OOB sample, accuracy is again calculated.
- Decrease in accuracy is averaged over all trees
 - Used as measure of importance of variable in random forest.



Experimental Protocol

➤ **Three datasets:**

- **Retrieval of hand-drawn/printed table images** (Approximate match)
- **Retrieval of handwritten mixed-forms** (Conceptual Match)
- **Grouping of NIST-tax forms** (Exact Match)
 - No training data

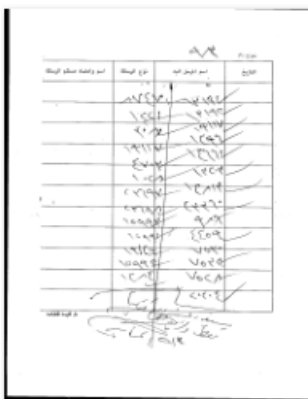
➤ **Evaluation:**

- **F1-Score** based on precision and recall of relevant documents
- **Purity** of clusters for grouping



Datasets

	Training	Testing
Table Dataset¹	150 tables/250 non-tables	132
Mixed-form Dataset²	240 form/320 non-forms	230
NIST Tax Forms	--	20 Classes 5590 Images



رقم	اسم العميل	تاريخ الميلاد	اسم العائلة	اسم العائلة
1
2
3
4
5
6
7
8
9
10

Sample from table dataset



1- الاسم الكامل: ...

2- التاريخ: ...

3- العنوان: ...

4- الهاتف: ...

5- البريد الإلكتروني: ...

6- الملاحظات: ...

7- التوقيع: ...

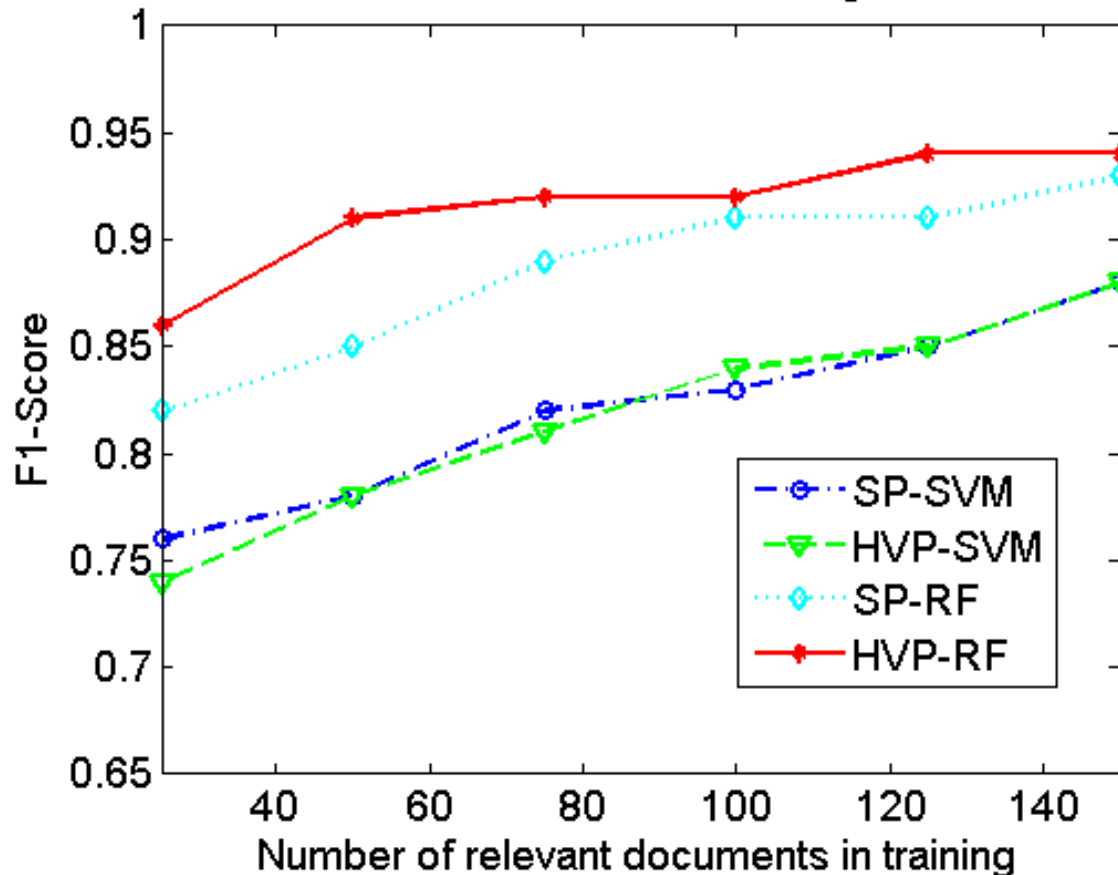
Sample from mixed-form dataset

1,2 Dataset available at:

<http://lampsrv02.umiacs.umd.edu/projdb/project.php?projType=1>

Results – Table images

Performance on table images



#Patches per image: 3000

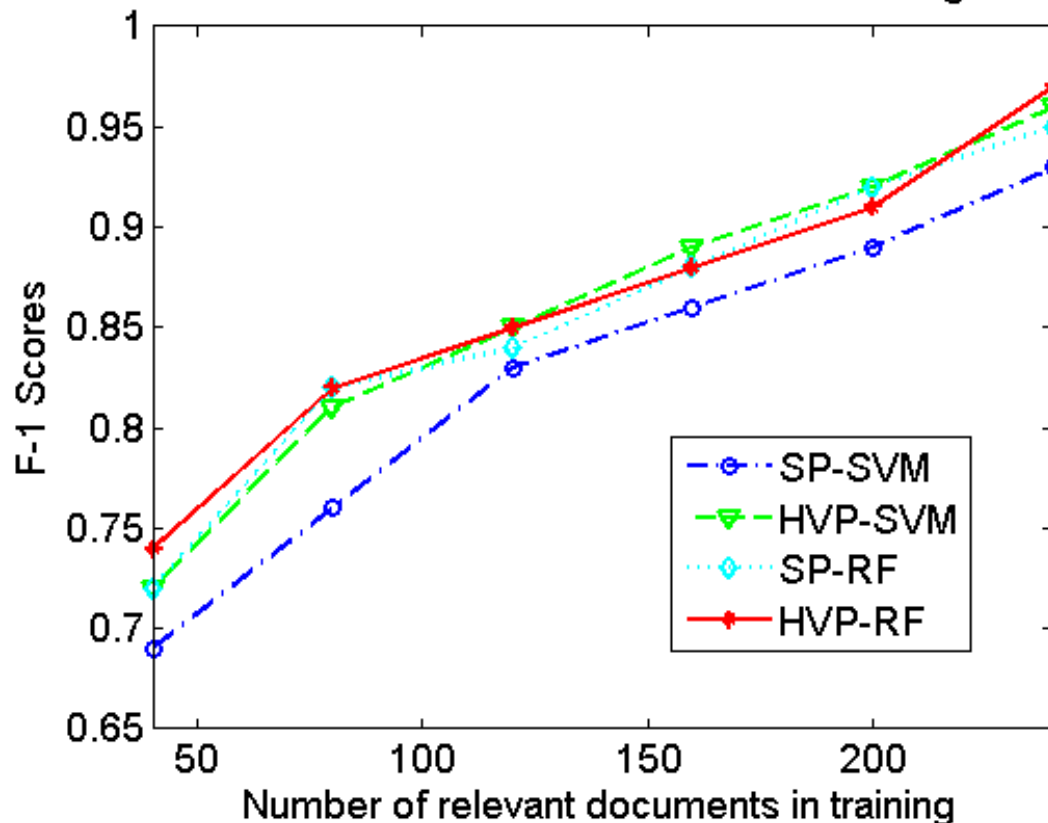
#nTrees: 1000

#mTry: $\sqrt{\text{\#attributes}}$

Accuracy on Balanced data: 97.8%

Results– Mixed-Form Images

Performance on mixed-form document images



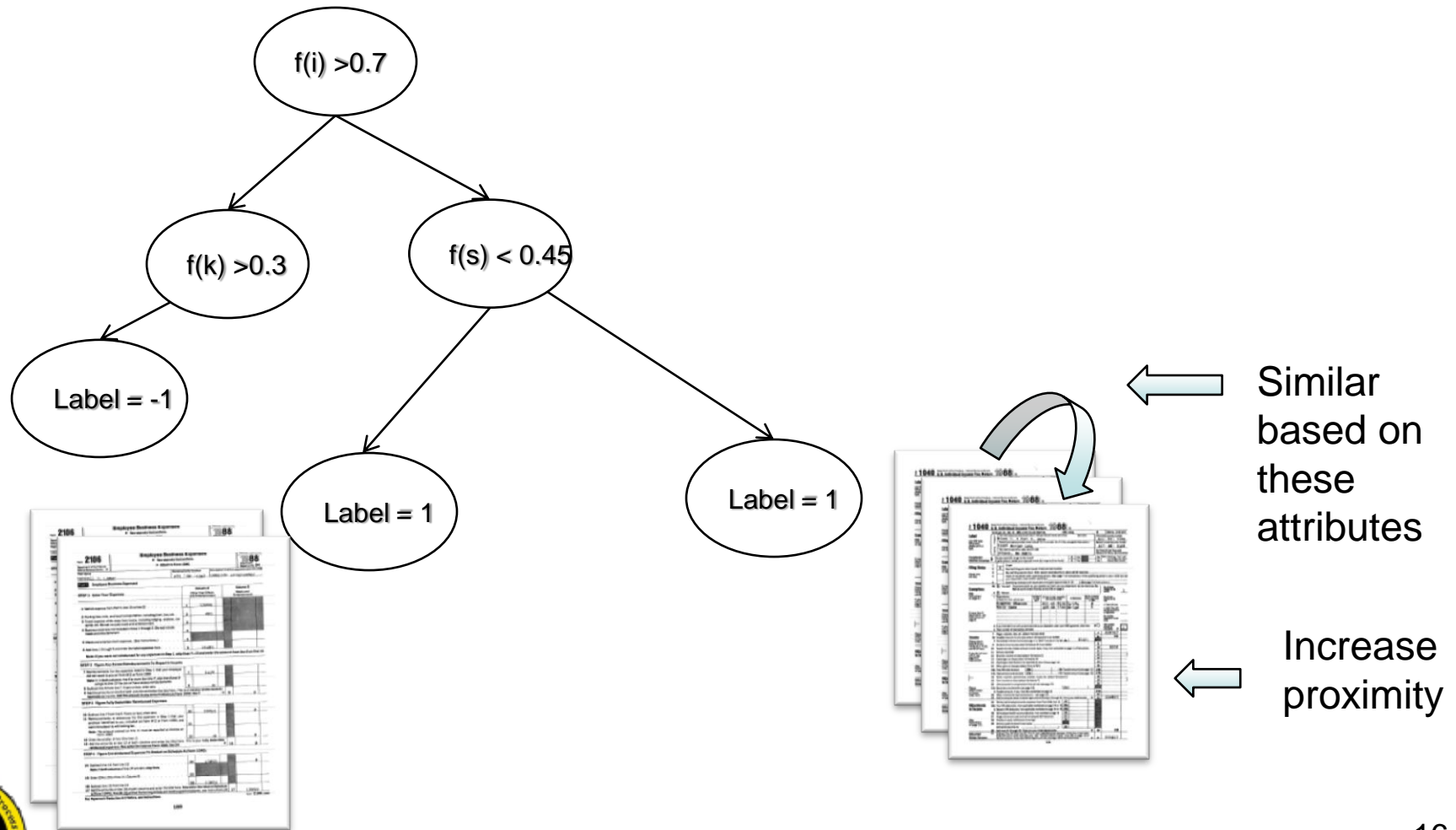
#Patches per image: 3000

#nTrees: 1000

#mTry: $\sqrt{\text{\#attributes}}$

Accuracy on Balanced data: 98.9%

Computing Proximities using Random Forest



Results on NIST Tax-form Images

- 20 different types of tax forms (1040_1, 1040_2, 2106_1, 2106_2, 4562_1, 6251 etc.)
- **Purity = 1.0 using Normalized-cuts** (Shi and Malik 2001)

This is a scan of a 2008 U.S. Individual Income Tax Return (Form 1040). The form is filled out with various numbers and text, including the taxpayer's name, address, and income details. It features a standard layout with numbered lines for different sections of the return.This is a scan of a 2008 Employee Business Expenses form (Form 2106). The form is filled out with details about the taxpayer's business expenses, including a list of expenses and their corresponding amounts. It includes a section for the taxpayer's signature and the preparer's information.This is a scan of a 2008 Depreciation and Amortization form (Form 4562). The form is filled out with details about the taxpayer's depreciation and amortization expenses, including a table of assets and their respective values. It includes a section for the taxpayer's signature and the preparer's information.This is a scan of a 2008 Tax on Short-Term Capital Gains and Dividends form (Form 6251). The form is filled out with details about the taxpayer's short-term capital gains and dividends, including a table of assets and their respective values. It includes a section for the taxpayer's signature and the preparer's information.

Summary and On-going Work

- Horizontal-Vertical pooling is an effective way to capture local structure statistics of document images
- Random Forest classifier is a good candidate for structural similarity based retrieval
- Approach is efficient and scalable
- Extensions possible to un-supervised and semi-supervised grouping of document images

Thank You!

