

Visual Segmentation for Information Extraction from Heterogeneous Visually Rich documents

Ritesh Sarkhel*, Arnab Nandi
Department of Computer Science and Engineering

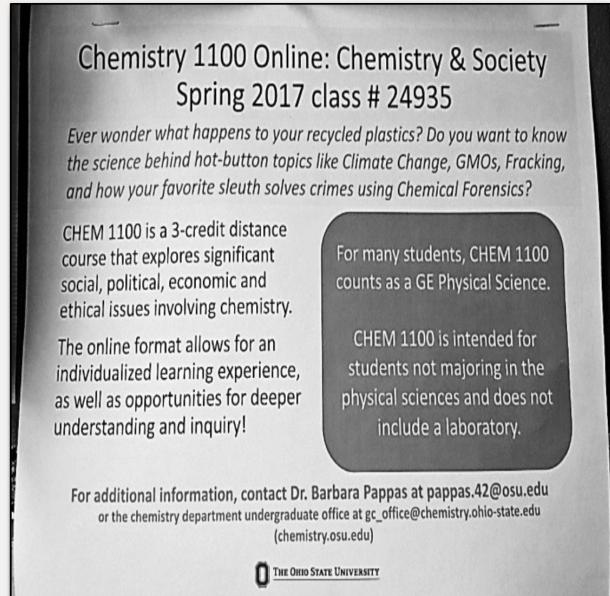


Outline

- ❖ Visually Rich Documents : Lots of useful information
- ❖ Challenges and observations
- ❖ Our contributions
 - VS2-Segment
 - VS2-Select
- ❖ Experimental results
 - Event information extraction
 - Outperforms text-based baseline by an average F1-score of 5.07%



Visually Rich Document



Chemistry 1100 Online: Chemistry & Society
Spring 2017 class # 24935

Ever wonder what happens to your recycled plastics? Do you want to know the science behind hot-button topics like Climate Change, GMOs, Fracking, and how your favorite sleuth solves crimes using Chemical Forensics?

CHEM 1100 is a 3-credit distance course that explores significant social, political, economic and ethical issues involving chemistry.

The online format allows for an individualized learning experience, as well as opportunities for deeper understanding and inquiry!

For additional information, contact Dr. Barbara Pappas at pappas.42@osu.edu or the chemistry department undergraduate office at gc_office@chemistry.ohio-state.edu (chemistry.osu.edu)

THE OHIO STATE UNIVERSITY

For many students, CHEM 1100 counts as a GE Physical Science.

CHEM 1100 is intended for students not majoring in the physical sciences and does not include a laboratory.



OPEN HOUSE
Friday, June 2nd 8AM - 6PM
320 Flyerco Way
Anaheim, CA

About The Property
To edit this text, click it once and edit the text in the menu to the right. Flyerco is easy to use and made for someone Realtors like yourself. If you want to replace an image you can click the image once and click replace image from the menu to the right as well.

Features

- Easy To Use
- Download Anywhere
- Print Anywhere
- Shareable Online
- Unlimited Downloads

Denise Amara
Realtor
714-350-8290
support@flyerco.com
www.flyerco.com

FlyerCo

- Posters
- Leaflets
- Banners etc.

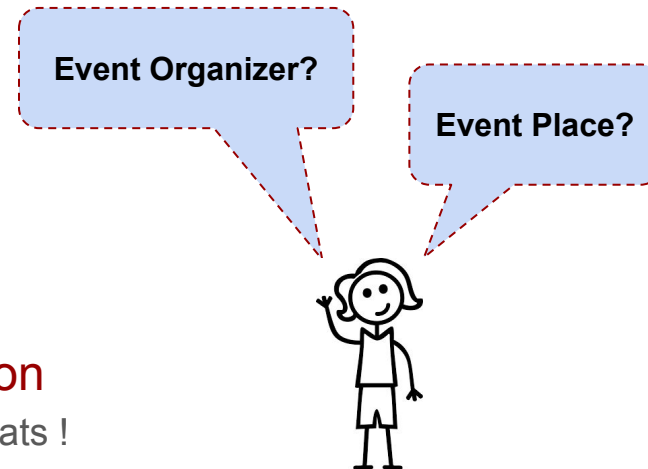
Properties of Visually Rich Documents

- Visual cues to highlight distinct semantic entities
 - Color
 - Font size
 - Text orientation and positioning
 - Negative-space
 - Symmetry etc.
- May be sparsely worded



Motivation

- **Rich source of ad-hoc content**
 - Contains a lot of useful information
 - Not easily available in an indexed database
- **Motivating example:**
 - Event information extraction
- **Objective: Automated Information Extraction**
 - **Recall:** Robust towards different document formats !



Outline

- ❖ Visually Rich Documents : Lots of useful information
- ❖ **Challenges and observations**
- ❖ Our contributions
 - VS2-Segment
 - VS2-Select
- ❖ Experimental results
 - Event information extraction
 - Outperforms text-based baseline by an average F1-score of 5.07%



Possible Solution 1: Text-based solutions?

Chemistry 1100 Online: Chemistry & Society
Spring 2017 class # 24935

Ever wonder what happens to your recycled plastics? Do you want to know the science behind hot-button topics like Climate Change, GMOs, Fracking, and how your favorite sleuth solves crimes using Chemical Forensics?

CHEM 1100 is a 3-credit distance course that explores significant social, political, economic and ethical issues involving chemistry.

The online format allows for an individualized learning experience, as well as opportunities for deeper understanding and inquiry!

For many students, CHEM 1100 counts as a GE Physical Science.

CHEM 1100 is intended for students not majoring in the physical sciences and does not include a laboratory.

For additional information, contact Dr. Barbara Pappas at pappas.42@osu.edu or the chemistry department undergraduate office at gc_office@chemistry.ohio-state.edu (chemistry.osu.edu)

THE OHIO STATE UNIVERSITY

- ❌ Identifying the semantic entities
- ❌ Determining context boundaries
 - Semantic role played by visual features not considered by text-based solutions
- ❌ Transcription errors



Possible Solution 1: Text-based solutions?


Chemistry 1100 Online: Chemistry & Society
Spring 2017 class # 24935

Ever wonder what happens to your recycled plastics? Do you want to know the science behind hot-button topics like Climate Change, GMOs, Fracking, and how your favorite sleuth solves crimes using Chemical Forensics?

CHEM 1100 is a 3-credit distance course that explores significant social, political, economic and ethical issues involving chemistry.

The online format allows for an individualized learning experience, as well as opportunities for deeper understanding and inquiry!

For additional information, contact Dr. Barbara Pappas at pappas.42@osu.edu or the chemistry department undergraduate office at gc_office@chemistry.ohio-state.edu (chemistry.osu.edu)

 THE OHIO STATE UNIVERSITY

Chemistry 1100 Online: Chemistry & Society
Spring 2017 class # 24935

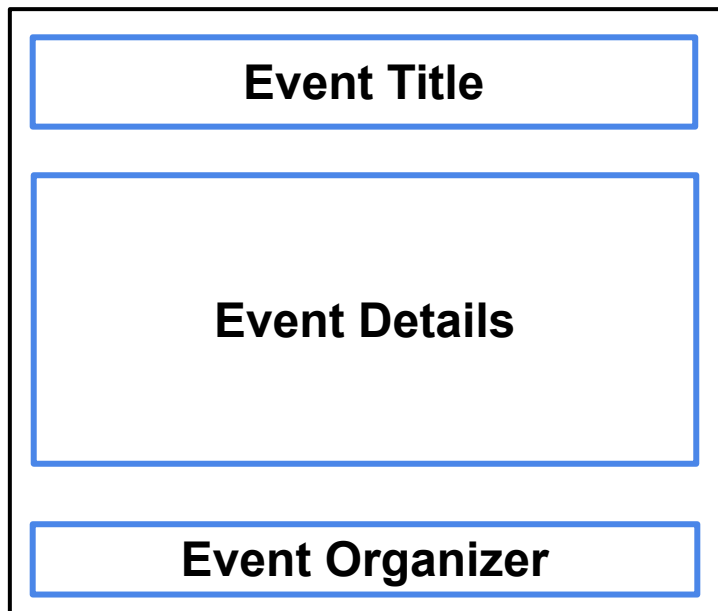
Ever wonder what happens to your recycled plastics? Do you want to know the science behind hot-button topics like Climate Change, GMOs, Fracking, and how your favorite sleuth solves crimes using Chemical Forensics?

CHEM 1100 is a 3-credit distance course that explores significant social, political, economic and ethical issues involving chemistry. The online format allows for an individualized learning experience, as well as opportunities for deeper understanding and inquiry!

For additional information, contact Dr. Barbara Pappas at pappas.42@osu.edu or the chemistry department undergraduate office at gc_office@chemistry.ohio-state.edu (chemistry.osu.edu)

THE OHIO STATE UNIVERSITY

Possible Solution 2: Visual rule-based solutions?



- **Template-based masks as rules**
 - Create and store masks to identify where different named entities appear
 - e.g. ReportMiner
- ⊘ **Cost of scalability**
 - Hard to maintain exact rules for all possible layouts
 - Expensive to deploy and maintain



Observations


Chemistry 1100 Online: Chemistry & Society
Spring 2017 class # 24935

Ever wonder what happens to your recycled plastics? Do you want to know the science behind hot-button topics like Climate Change, GMOs, Fracking, and how your favorite sleuth solves crimes using Chemical Forensics?

CHEM 1100 is a 3-credit distance course that explores significant social, political, economic and ethical issues involving chemistry. The online format allows for an individualized learning experience, as well as opportunities for deeper understanding and inquiry!

For many students, CHEM 1100 counts as a GE Physical Science. CHEM 1100 is intended for students not majoring in the physical sciences and does not include a laboratory.

For additional information, contact Dr. Barbara Pappas at pappas.42@osu.edu or the chemistry department undergraduate office at gc_office@chemistry.ohio-state.edu (chemistry.osu.edu)

 THE OHIO STATE UNIVERSITY

- Visually rich document = Bag of isolated coherent visual elements
 - Visual heterogeneity is significantly reduced in each of these coherent visual areas
- A few of these visual areas act as 'interest-points'
 - Areas with high semantic significance

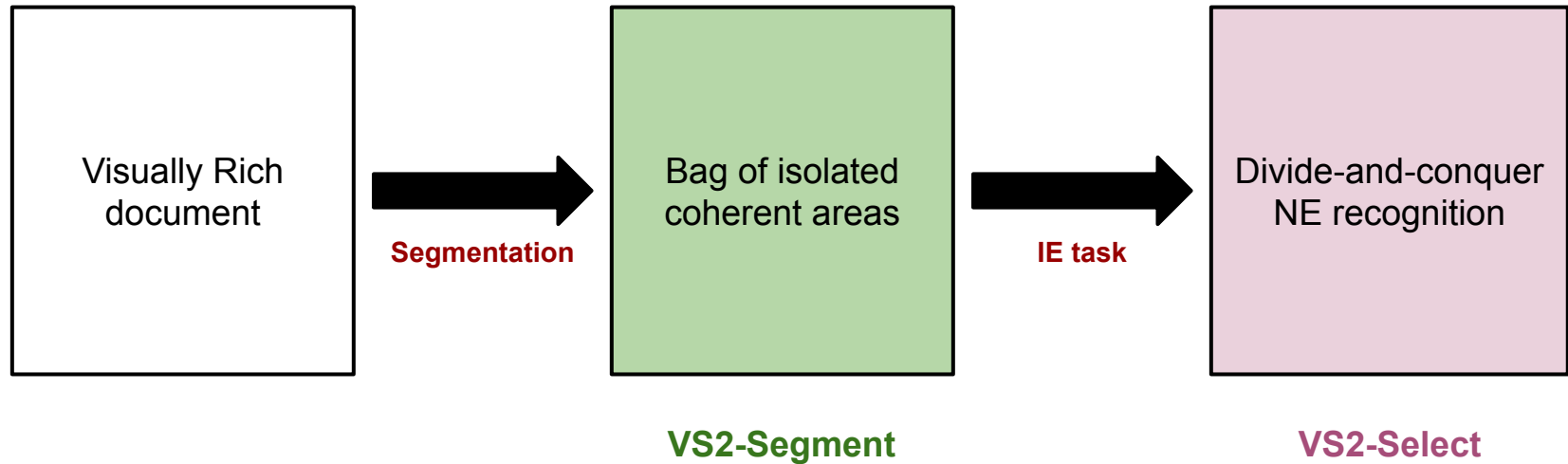


Our hypotheses

1. Decompose a document into a bag of visually isolated coherent areas
 - Incorporate both visual and semantic features
2. Divide-and-conquer by text-based IE methods within each segmented area



Proposed IE workflow for Visually Rich documents



Outline

- ❖ Visually Rich Documents : Lots of useful information
- ❖ Challenges and observations
- ❖ **Our contributions**
 - **VS2-Segment**
 - VS2-Select
- ❖ Experimental results
 - Event information extraction
 - Outperforms text-based baseline by an average F1-score of 5.07%



VS2-Segment: Visual Segmentation by layout analysis

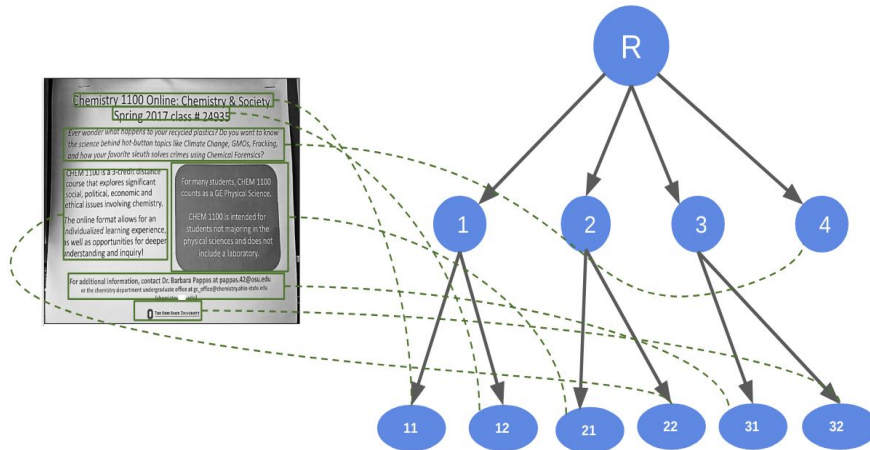
Input: Visually Rich document (D)

Output: Bag of isolated coherent visual areas

Constraint: Robust towards document formats



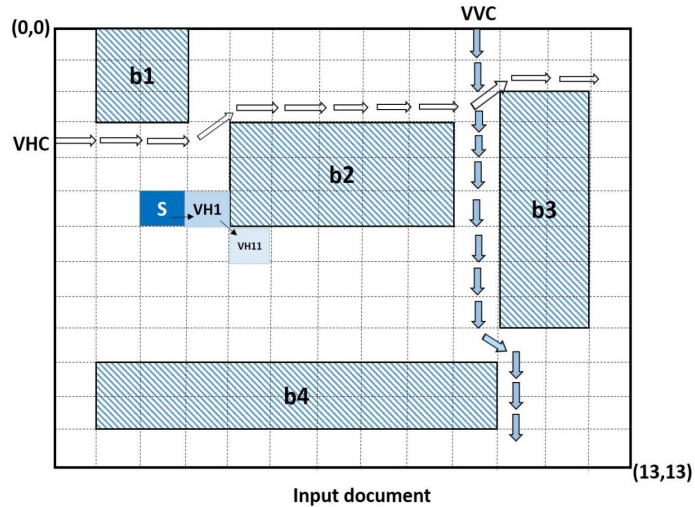
VS2-Segment: Segmentation by layout analysis



- A hierarchical layout-tree of D
 - Node → A visual area in D
 - Parent → Contains the visual area denoted by child
- Visual features: Explicit and Implicit
- Semantic features: Cosine similarity
- Return leaf nodes when terminates



VS2-Segment: Feature Descriptors



Explicit Visual Features

Index	Feature
1	Centroid coordinates
2	Height
3	LAB color
4	Angular distance from origin
5	Pairwise sum of angular distances

Implicit Visual Features



VS2-Segment: What did we learn so far ?

Chemistry 1100 Online: Chemistry & Society

Spring 2017 class # 24935

Ever wonder what happens to your recycled plastics? Do you want to know the science behind hot-button topics like Climate Change, GMOs, Fracking, and how your favorite sleuth solves crimes using Chemical Forensics?

CHEM 1100 is a 3-credit distance course that explores significant social, political, economic and ethical issues involving chemistry.

The online format allows for an individualized learning experience, as well as opportunities for deeper understanding and inquiry!

For many students, CHEM 1100 counts as a GE Physical Science.

CHEM 1100 is intended for students not majoring in the physical sciences and does not include a laboratory.

For additional information, contact Dr. Barbara Pappas at pappas.42@osu.edu or the chemistry department undergraduate office at gc_office@chemistry.ohio-state.edu (chemistry.osu.edu)

THE OHIO STATE UNIVERSITY

- A recursive algorithm to find semantically coherent areas
- Incorporates both visual and semantic features

Next up: IE by divide-and-conquer



Outline

- ❖ Visually Rich Documents : Lots of useful information
- ❖ Challenges and observations
- ❖ **Our contributions**
 - VS2-Segment
 - **VS2-Select**
- ❖ Experimental results
 - Event information extraction
 - Outperforms text-based baseline by an average F1-score of 5.07%



VS2-Select: Named entity recognition using distant supervision

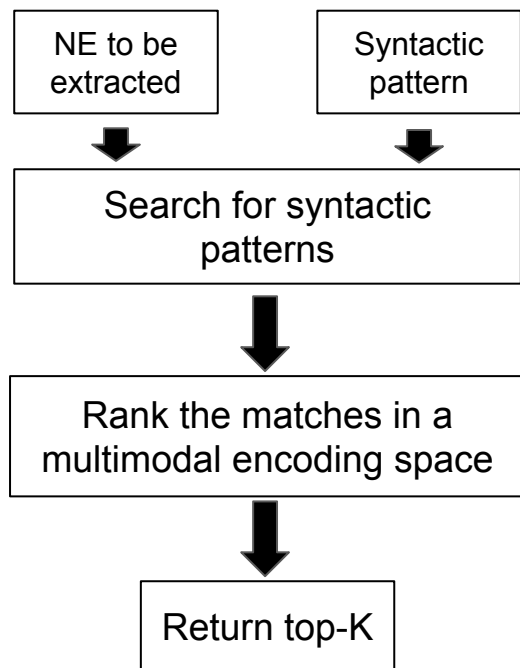
Input: Bag of coherent visual areas + NE to be extracted

Output: Text extracted corresponding to the NE

Constraint: Robust towards ad-hoc content



VS2-Select: Named entity recognition using distant supervision



Divide
&
Conquer !

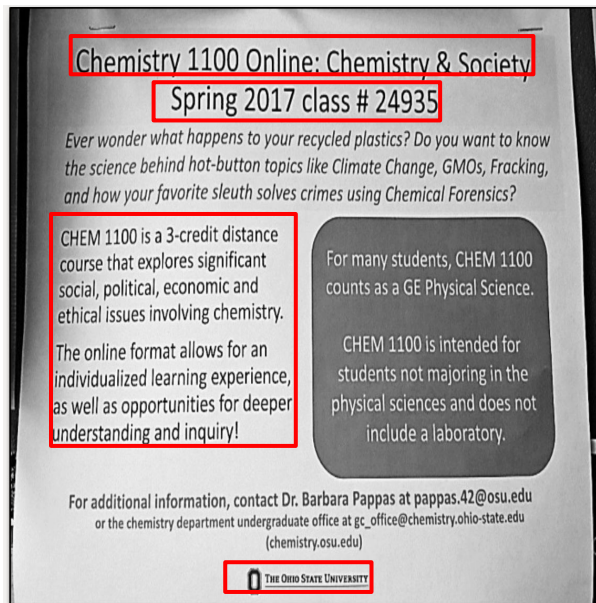


VS2-Select: Learning the syntactic patterns

- For each appearance **T** of an NE in a holdout corpus **H**:
 - Annotate and construct the dependency parse tree of **T**
 - Identify the patterns **P** represented by most frequent subtree(s)
 - **P** represents the NE
- Search for **P** within each segmented area



VS2-Select: Rank the matched patterns



- Rank the matches based on relative semantic significance
 - Metric: Weighted L1 distance from the closest ‘interest point’ in the document
 - Return top-K
- ‘Interest Point’ → Higher semantic significance
 - Subset of visual areas obtained from segmentation
 - Pareto-optimal front in a multi-objective optimization problem
 - Each objective represents a design principle

Some Interest points of the poster



VS2-Select: What did we learn so far ?

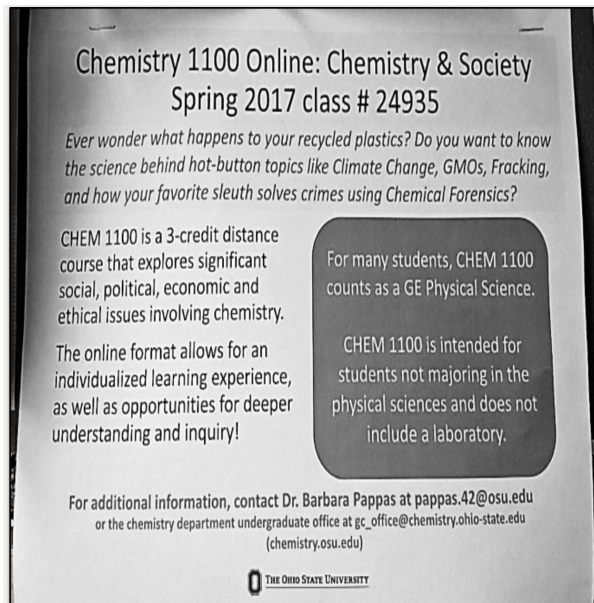
- Identifies NE's by leveraging the context boundaries determined by VS2-Segment
- Distantly supervised syntactic patterns represent each NE.

Outline

- ❖ Visually Rich Documents : Lots of useful information
- ❖ Challenges and observations
- ❖ Our contributions
 - VS2-Segment
 - VS2-Select
- ❖ **Experimental results**
 - **Event information extraction**
 - **Outperforms text-based baseline by an average F1-score of 5.07%**



Experiments



- 2190 event posters and flyers
 - Heterogeneous layouts
 - Digital images and PDF's
- 5 distinct named entities
 - Related to event information

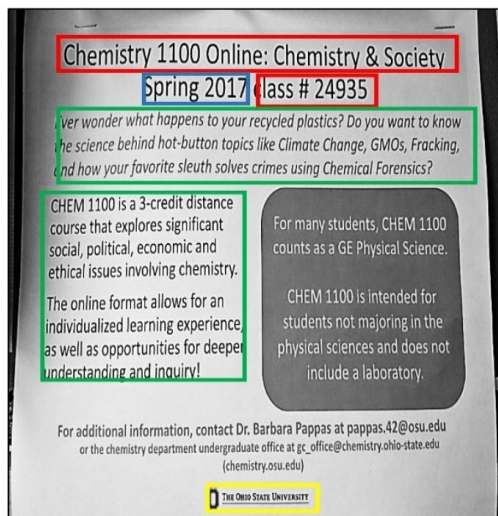


Experiment: Event information extraction

Named entity	Definition
Event Title	Short description of the event
Event Place	Full address of the event
Event Time	Time of the event
Event Organizer	Host of the event
Event Details	Other details of the event



Evaluation metrics: Two-phase evaluation



Color Legends

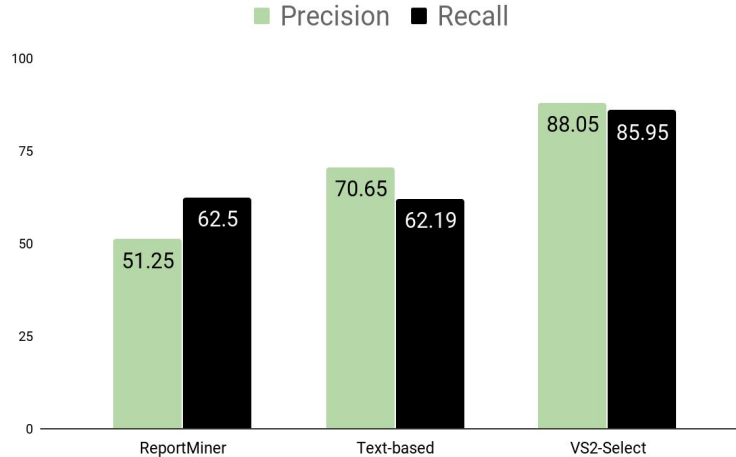
- Event title
- Event description
- Event organizer
- Event time
- Event place

- Text T identified as NE-type N is accurate:
 - a. Its position is accurate **and** the named entity type N is accurate too
 - b. Compared against human-annotated groundtruth



Quantitative evaluation: End-to-end

- Objective: End-to-end performance ?
- Accurate iff position coincides with ground-truth **and** inferred NE type is correct



Quantitative evaluation: End-to-end

- **Takeaways:**
 - a. VS2 demonstrates robust performance for all NE types
 - b. Both localization and classification capabilities are robust.
 - c. Improves over ReportMiner¹, a visual template driven rule-based method
 - d. Improvement over text-only baseline is significant

More experiments and analysis in paper!



Takeaways

- Prior segmentation helps IE from visually rich documents
 - Visual features need to be considered
- VS2: VS2-Segment and VS2-Select
 - Effectiveness of NE extraction explored
- Opens up a lot of exciting possibilities !
 - Explore the semantic relationship between the segmented visual areas
 - Utilize latent semantic information for more complex IE tasks -- QA, Summarization etc.

Thank you!

Questions?



Supplementary

Examples of learned lexical / syntactic patterns

Named entity	Learned patterns
Event Title	{ <i>VP</i> , <i>NP</i> with numeric (<i>CD</i>) and textual (<i>JJ</i>) modifiers, <i>SVO</i> }
Event Place	<i>NP</i> with valid geocode tags
Event Time	<i>NP</i> with <i>TIMEX3</i> tags
Event Organizer	{ <i>Captain</i> / <i>Create</i> / <i>Reflexive-appearance</i> verb-senses, <i>NP</i> with <i>Person</i> / <i>Organization</i> named entities }



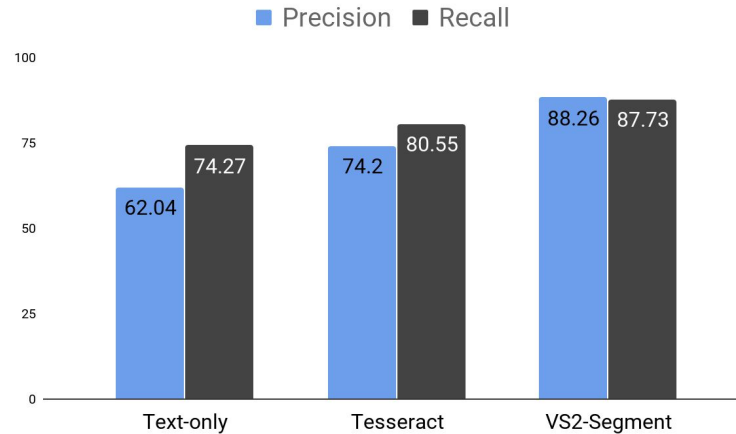
Learning the syntactic patterns

- Holdout corpus **H** created by scraping public-domain website
 - URL: allevents.in
 - Filters: Place and Date
- **H** contains 500 structured tuples
 - Tuple \rightarrow {Named entity type, Corresponding text}
- Derive lexico-syntactic patterns for every named entity type in **H**

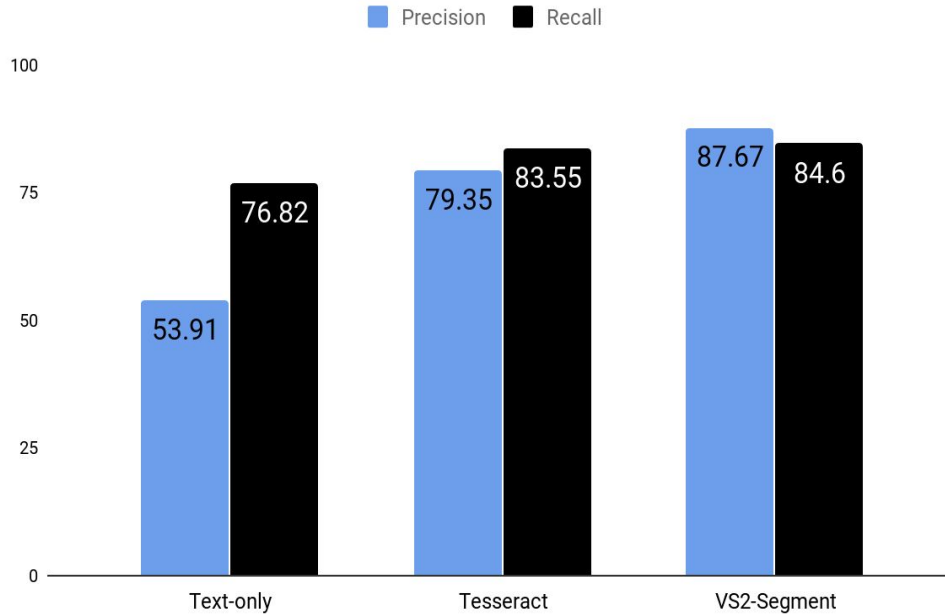


Quantitative evaluation: Localization

- Objective: Are the extracted named entities accurately localized ?
- Accurate if coincides with ground-truth



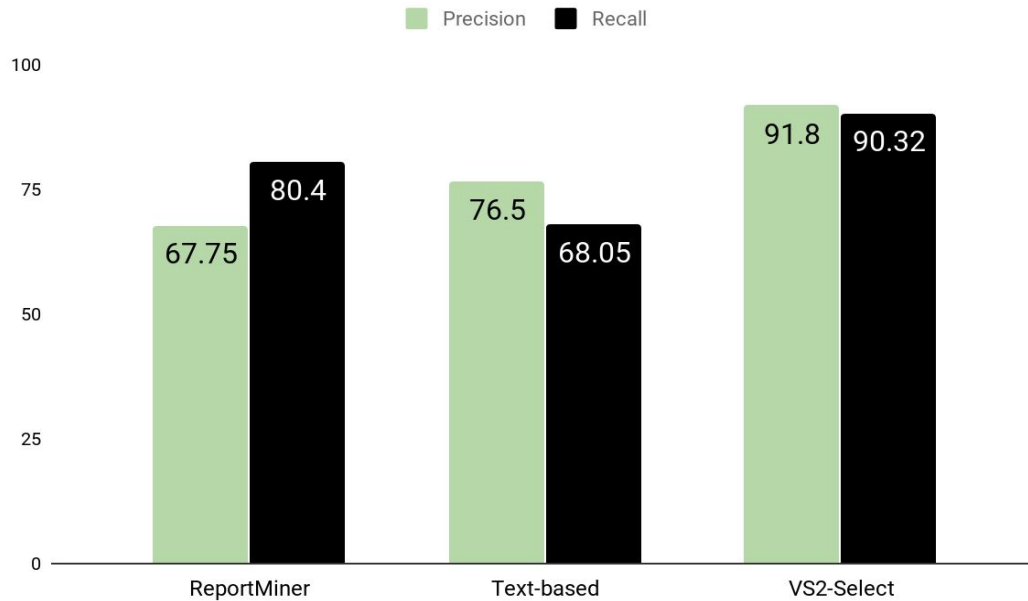
Property Information Extraction: Localization capability



- Accurate iff IoU score against ground-truth ≥ 0.65

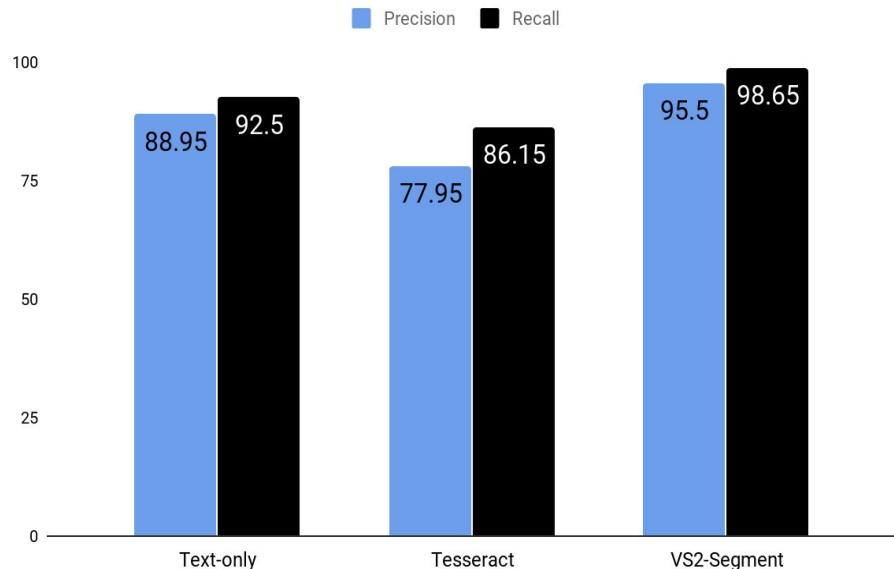


Property Information Extraction: Classification capability



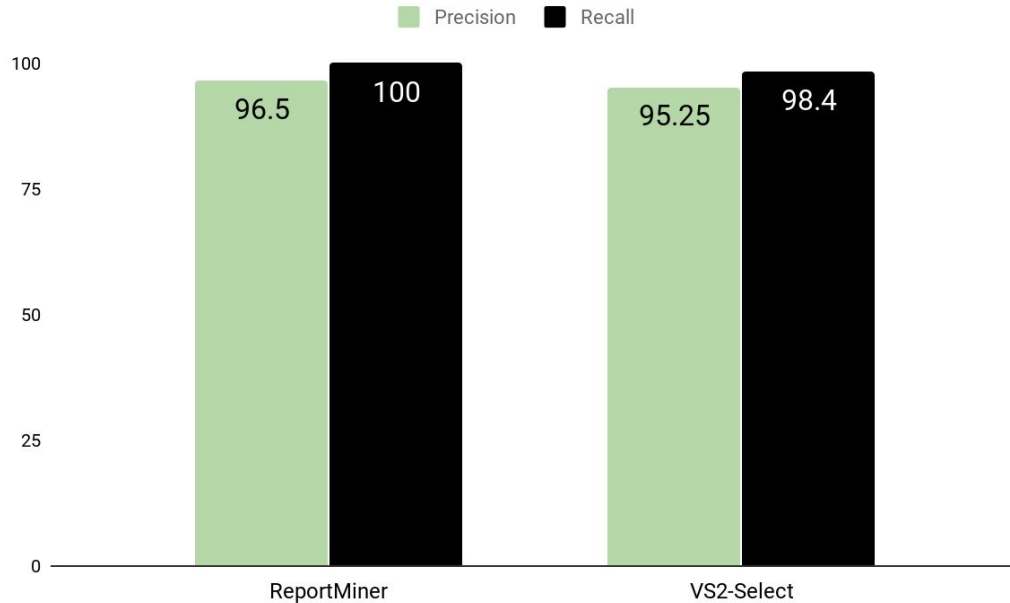
- Accurate iff IoU score ≥ 0.65 & inferred NE type is correct

Form field extraction: Localization capability



- Accurate iff IoU score against ground-truth ≥ 0.65

Form field extraction: Classification capability



- Accurate iff IoU score ≥ 0.65 & inferred NE type is correct



Ablation study

Index	VS2-Segment		VS2-Select	$\Delta F1$ (%)		
	<i>Visual feature</i>	<i>Semantic feature based merging</i>	<i>Entity disambiguation</i>	<i>D1</i>	<i>D2</i>	<i>D3</i>
A1	✓	×	✓	0.80	2.55	3.37
A2	×	✓	✓	1.07	4.22	3.84
A3	✓	✓	×	0.95	6.78	7.05
A4	✓	✓	Text-only	0.42	4.55	3.96



Syntactic Patterns learned for the Real-Estate flyers dataset

Named entity type	Description	Syntactic patterns to search
Broker Name	Full name of the listing broker	A bigram/trigram of NE's with <i>Person</i> / <i>Organization</i> tags
Broker Phone	Contact number of the listing broker	A regular expression containing digits, characters and separators such as '-', '(', ')', and '.'
Broker Email	Email address of the listing broker	An RFC-5322 compliant regular expression containing character and separators such as '@', and '.'
Property Address	Full address information of the listing	Noun phrase with valid geocode tags
Property Size	Size-attributes summarizing the size of a listing (e.g. 4 beds, 2,465 acres)	(1) Noun phrase with numeric (<i>CD</i>) or textual modifiers (<i>JJ</i>) and (2) Noun POS tags with senses <i>measure</i> / <i>structure</i> / <i>estate</i> in the Hypernym Tree [42]
Property Description	Mentions of the property type (e.g. building, floor, land/lot) and other essential details (e.g. parking, grocery)	Noun phrases with numeric (<i>CD</i>) or textual modifiers (<i>JJ</i>)



NE-specific breakdown for Event Information extraction

Index	Named Entity	Proposed method		
		<i>Pr. (%)</i>	<i>Rec. (%)</i>	$\Delta F1(\%)$
N1	<i>Event Title</i>	84.88	81.09	8.98
N2	<i>Event Place</i>	76.68	86.37	3.76
N3	<i>Event Time</i>	94.67	84.70	0.49
N4	<i>Event Organizer</i>	72.56	74.41	10.50
N5	<i>Event Description</i>	76.59	86.00	1.60
Overall		81.08	82.51	5.07

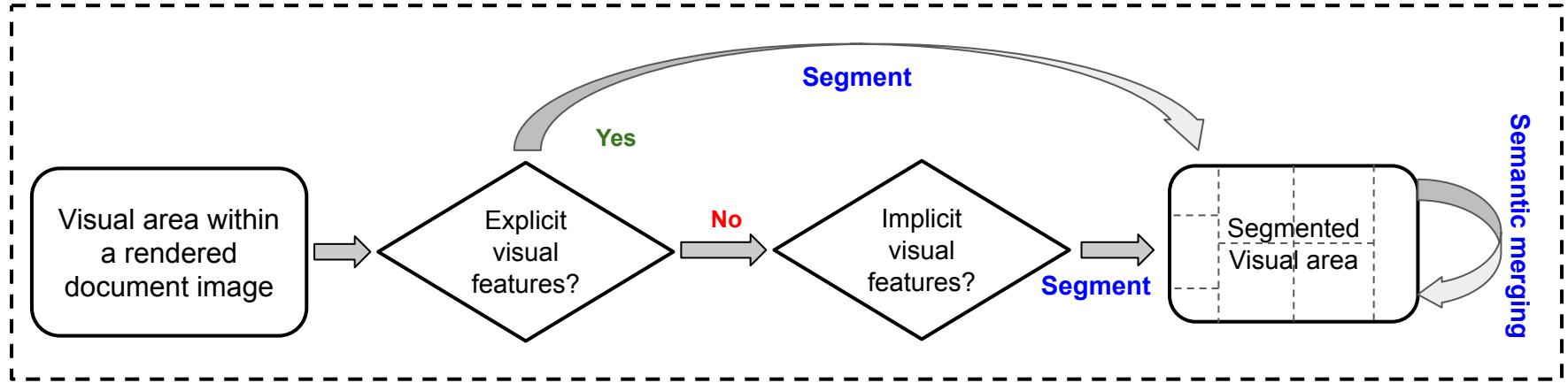


NE-specific breakdown for Property Information extraction

Index	Named Entity	Proposed method		
		<i>Pr.</i> (%)	<i>Rec.</i> (%)	$\Delta F1$ (%)
N1	<i>Broker Name</i>	94.72	90.85	10.18
N2	<i>Broker Phone</i>	96.15	82.25	1.63
N3	<i>Broker Email</i>	97.25	95.40	2.56
N4	<i>Property Address</i>	92.68	85.50	4.60
N5	<i>Property Size</i>	85.25	93.05	3.37
N6	<i>Property Desc.</i>	84.75	94.90	0.74
Overall		91.80	90.32	3.84



VS2-Segment: Visual Segmentation by layout analysis



- Recursive operation to decompose a document into smaller visual areas