Visual Segmentation for Information Extraction fr **Heterogeneous Visually Rich Documents**

Ritesh Sarkhel The Ohio State University sarkhel.5@osu.edu

Arnab Nandi The Ohio State University arnab@cse.osu.edu

ABSTRACT

Physical and digital documents often contain visually rich information. With such information, there is no strict ordering or positioning in the document where the data values must appear. Along with textual cues, these documents often also rely on salient visual features to define distinct semantic boundaries and augment the information they disseminate. When performing information extraction (IE)traditional techniques fall short, as they use a text-only representation and do not consider the visual cues inherent to the layout of these documents. We propose VS2, a generalized approach formation for information extraction from heterogeneous visually rich documents. There are two major contributions of this work. First, we propose a robust segmentation algorithm that de-Visually Rich document; Information Extraction; Named encomposes a visually rich document into a bag of visually iso-tity lated but semantically coherent areas, called logical blocks ACM Reference Format: Document type agnostic low-level visual and semantic features are used in this process. Our second contribution is a distantly supervised search-and-select method for identifying the named entities within these documents by utilizing perimental results on three heterogeneous datasets suggest that the proposed approach significantly outperforms its text-only counterparts on all datasets. Comparing it against Information extraction from documents has been widely the state-of-the-art methods also reveal that VS2 performs studied for a number of different applications in the past comparably or better on all datasets.

CCS CONCEPTS

 Information systems → Information extraction ument structure ontent analysis and feature selection; En-contemporary literature [13, 29] on information extraction tity resolution:

Permission to make digital or hard copies of all or part of this work for personalor classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org. or augment the semantics of various parts of the document. SIGMOD '19, June 30-July 5, 2019, Amsterdam, Netherlands

© 2019 Copyright held by the owner/author(s). Publication rights licensed to Association for Computing Machinery.

ACM ISBN 978-1-4503-5643-5/19/06. . . \$15.00 https://doi.org/10.1145/3299869.3319867







(a) Academic event poster with highlights on the class topic, class timing and scope

highlights on the property listing, and the broker name

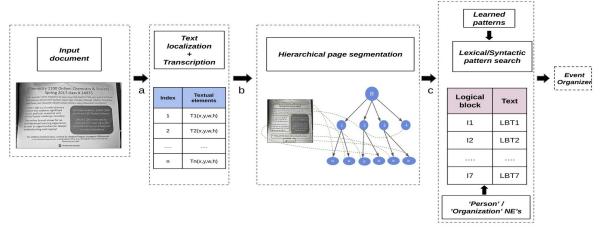
Figure 1:Samples of visually rich documents using salient visual features to highlight important event in-

KEYWORDS

Ritesh Sarkhel and Arnab Nand 1019. Visual Segmentation for Information Extraction from Heterogeneous Visually Rich Documents. In 2019 International Conference on Management of Data (SIG-MOD '19), June 30-July 5, 2019, Amsterdam, Netherlands. ACM, New the context boundaries defined by these logical blocks. Ex-York, NY, USA, 16 pages. https://doi.org/10.1145/3299869.3319867

INTRODUCTION

few years. This includes spam filtering event extraction, social-media-text mining, text classification, and web-based document retrievals. Most of these works, however, are developed and evaluated on a text-only corpus. Consequently, is heavily biased towards purely textual features, e.g. lexical features (e.g. stemming), contextual features (e.g., n-gram), and document features (e.g., TF-IDF). However, a significant amount of textual information that we encounter everyday is presented in the form of visually rich documents. Instead of relying solely on text-based cues, these documents use salient visual features such as text positioning, font-size distribution, typographical similarity, and color distribution to highlight A few examples are shown in Fig. 1. In many cases, these documents prove to be valuable resources of useful information, not readily available in an indexed database that can be searched for a quick lookup. For example, membership



A working example of VS2 extracting the 'Event Organizer' information from an academic event poster

Figure 2: An overview of VS2; Taking a visually rich document as input, it outputs the text correspondi named entity from the document. Upon input, the document is cleaned and its textual elements are loc a), we identify its logical blocks (Step b). This is achieved using the layout model of the document conusing a hierarchical segmentation algorithm. Finally, a set of predefined lexical or syntactic patterns (searched within each logical block to identify the text corresponding to the named entity to be extract

discount flyers from retailers comparing online commercial real-estate enlisting by different agenciescheduling local events from event posters; all of these scenarios requireclearly defined sentence structures and context boundaries extracting structured information from visually rich documents for downstream processing. We propose a generalized tactic properties. Due to atypical visual properties, defining method for automated information extraction from such documents in this work. To better explain our contributions. we demonstrate the limitations of a traditional information extraction (IE) method from visually rich documents in the following example.

Example 1.1; Alice, owner of a small event management She has collected a number of relevant event posters for this cially available document workflow management systems purpose and needs to extract a set of named entities { Event Title, Event Organizer } from these documentsor each named entit $\mathbf{Q}_i \in N, i = 1, 2$, she wants the corresponding text t_i extracted from these documents. In scenarios such as these, a traditional text-based IE system starts with clean-disambiguation methods30 are used to resolve conflicts ing (which includes perspective warping, skew correction, and binarization) the document first. Then the document is transcribed and its text is searched for some lexical and/or didates based on some contextual information of where they syntactic patterns, predefined for each named entity to be ex-appeared in the documents traditional disambiguation tracted. For example, when searching for Event Organizers instrategies fail to incorporate visual features of the document, the OCR'ed 36 transcription of an event poster, Alice may search for phrases that represent a 'Person' or 'Organization' Fig. 3) is also a challenging task itself. in the document. If there are multiple such candidates word sense disambiguation strateg 30 may also be employed at this stage. Although reasonable for unstructured text, there are two major challenges in following a similar approach for visually rich event posters.

Challenge 1 Most of the natural language libraries and semantic parsers used by traditional methods rely on in the input text for determining various lexical and syncontext boundaries in the transcribed text of a visually rich document may prove to be challenging (see Fig. 3). Errors introduced during optical character recognitio 37, 39 of the document may adversely affect the quality of the downstream extraction task too. Now, if the document template is known beforehandcustom rules can be generated and company in Columbus, wants to survey some local events.applied on a case-by-case basis. In fact, a majority of commerfollow this scheme. These approaches, however, require a significant amount of cost and effort to maintain, making it hard to scale for diverse document types.

> Challenge 2: In an information extraction workflow, entity if there are multiple matches for a named entity within the input document. Final selection is made by ranking the cantranslating their success to visually rich documents (refer to

The objective of this work is to propose a generalized approach for information extraction from visually rich documents. We propose VS2, a two-phase approach, to this end. Following the guidelines proposed by previous researchers [35] it should have the following properties.

P1.1: Ability to intake heterogeneous documents i.e., not relying on prior knowledge about document layout or format to perform information extraction.

P1.2: Robustness i.e., flexibility to be extended for different extraction tasks.

Based on the conceptual framework by Doan, Naughton, Ramakrishnan and Baid1[1], our goal is to extract a list of key-value pairs from the document. The keys originate from a predefined semantic vocabular VS2 retrieves the corresponding text entries from the input document. This list of key-value pairs can be loaded into a database after schema mapping. Along with traditional full-text queries, it also offers the capability to perform rich semantic queries on the document.

Upon input, VS2 starts with cleaning and localizing the textual elements of a document first. Next, it is decomposed into a number of semantically coherent visual areas, called rithm. Identifying the logical blocks helps define the context boundaries within the document prior to any semantic parsing. Once the logical blocks have been identified set of lexico-syntactic patternsdefined for the named entity to be extracted is searched within the text transcribed from each logical block. In the case of multiple matched patterns, matched patterns and their closest interest paint the docconflict resolution is performed using a multimodal entity named entity, using a text-only holdout corpus. This distant biguation strategy will be described in Section 5.2. supervision [1, 27] enables us to circumvent the necessity of learning extraction rules every time a document with unfamiliar layout or format is presented. In other words, it makes it easy to process large-scale heterogeneous datasats. overview of the proposed approach is presented in Fig.2.

Technical contribution 10ur first contribution for this geneous visually rich documents in a unified way. We pro-comparable for all three tasks. pose VS2-Segment, a robust, hierarchical page segmentation RELATED WORKS algorithm that decomposes a document into semantically coherent visual areas, called 'logical blocks' for this purpose. document layout model used to enable the segmentation pro- is wrapper induction [19, 20]. Layout specific custom masks the segmentation algorithm will be provided in Section 5.1.

Technical contribution 2Our second major contribution is a distantly supervised search-and-select method for ments. Most commercially available systems e.g., Shr@dtr [identifying the named entities to be extracted within the logical blocks of a visually rich document. We propose VS2 set of syntactic pattern(s) within each logical block and se-document, the most appropriate rule is selected manually lects the most optimal matched pattern. In the case of mul- for extracting relevant information on a case-by-case basis. tiple matches (see Fig.), conflict resolution is performed

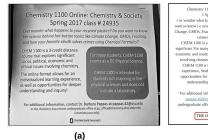


Figure 3:An illustration of the inherent challenges faced by a text-only approach for information extraction from an academieventpostersampled from dataseD2; (b) shows its transcription using Tesserac[41]; The red bounding boxes denote named entities belonging to the categories 'Person' or 'Organization' as recognized by the Stanford NER, logical blocks using a robust hierarchical segmentation algo-representing potential matches for the named entity 'EventOrganizer'Most of the false positives stem from errors during transcription and II-defined context boundaries in the transcribed text.

by a novel optimization-based multimodal entity disambiguation strategy. For each named entity, distances between the ument are minimized in a multimodal encoding space. The disambiguation strategy. The patterns are learned for each search-and-select method as well as the multimodal disam-

Summary of Results: We have evaluated VS2 on three heterogeneous datasets of visually rich documents. A comparative analysis against traditional text-based IE approaches (in Section 6.3 and 6.4 for separate IE tasks eveal that the proposed method performs significantly better than these approaches on all datasets. A comparison against the statepurpose is a robust encoding technique to represent heterof-the-art methods also reveals that VS2 performs better or

Most of the early works on IE from visually rich documents take advantage of the prior knowledge of the layout of a Each document is represented as a bag of logical blocks. The document. One of the most popular approaches among these cess will be introduced in Section 4. Further details regarding are defined to localize and extract information from the document. Researchers like Kushmerick et 1991 fand Mironczuk et al. [28] followed this approach for IE from HTML docu-ReportMiner 22 follow a similar approach. Relying on interactive interfaces, custom rules are designed for each layout Select, a distantly supervised approach that searches for aby experts and stored in a cloud-based server. For each test

¹An interest point is a visual area in the document that is visually and/or semantically significant

to generalize for diverse document types. For better gener-to extend their framework for different document types, iralization capabilities, some recent works, 21, 23, 43, 46 have also proposed heuristics-based extraction grammars Compared to the existing literature, one of the major conthat leverage visual information by analyzing the layout of the document. Contrary to these methods, VS2 does not as-diverse document layouts and format@ontrary to most sume any prior knowledge about the layout or format of the document, making it robust (refer to P1.1 in Section 1) to diverse document types.

Leveraging the homogeneity of the document format, a significant fraction of the existing work exploits high-level features defined by the document markup language2In [vert PDF documents into HTML format by assuming compli- information extraction tasks. ance with ISO 32000-1/32000 specifications. In a follow-upstudy, Gallo et al. 14 showed that this may be a strong assumption for many real-world documents that do not strictly conform with these specifications, as a slight misuse of the Suppose, we want to extract a set of named entitiesN= format operators in PDF stream may result in degraded vi- $\{n_1, n_2, ... n_p\}$ from a visually rich documen \mathcal{D} . Therefore, sual descriptors in the converted HTML document. Similar limitations can be observed in 4 too. Their extractors are trained on high-level features supported by the PostScript the text t_i and the named entit p_i , $\forall i \in [1, p]$ denotes the format, making it hard to generalize for heterogeneous doc- set of all textual elements in D. ument formats. Unlike these methods, VS2 does not make format specific assumptions in its feature design.

segmenting web pages into coherent visual blocks was pro-B. Our core insight here is that a visually rich document posed by Cai et al.4. Each web page is recursively decom- is a nested object comprised of distinct visualreas that posed into smaller blocks based on a set of carefully designe@re isolated from each other but semantically coherent by rules, defined using HTML tags and the vertical and/or horizontal whitespace separators that delineate them. Compared context boundaries of the document. We refer to eachs to [4], the segmentation algorithm proposed in VS2 (detail a logical block o₱. Once the logical blocks are identified, ment types. One of the major limitations of is its inability to be extended for various document formatshis is due to its reliance on HTML-specific tags to define various visual properties. Another major advantage of VS2 over [is its ability to segment overlapping blocks i.e., visual areas m, such that $n: N \to B$. Therefore, given a docume A, the that are not separated by a rectangular (vertical/horizontal) task $\mathfrak G$) of information extraction from D is decomposed whitespace separator. Gatterbauer et all also proposed a document-type agnostic approach for web-table construction by performing a layout analysis of rendered web pages, that represents D as a bag of logical blocks. VS2's scope of usability is much broader that [as it can be applied for non-HTML documents as well as a number ping $m: \mathbb{N} \to \mathbb{B}$, for each named entity $\epsilon \in \mathbb{N}$, that selects a set of non-trivial semantic tasks (refer to P1.2 in Section 1). In of textual elements within the boundaries defined by each some of their recent works, such as Fondutti and Deep-Dive [32], Re et al. have proposed machine-learning based solutions to this problem. A combination of visual and textual tion algorithm for the first sub-task. A hierarchical layout features are used to learn sequential patterns for extracting model is constructed by the segmentation algorithm to repren-ary relational tuples from each document. High-level features defined by document markup languages including XML and HTML were used to train their model for this purpose.

These approaches, however, are expensive to scale and hard S2 complements these methods by offering the flexibility respective of the layout or format of the input document. tributions of our work lies in the fact that it is robust to previous works, VS2 relies on a set of low-level visual and semantic features that can be extended to diverse document types (refer to P1.1 in Section 1) to localize and extract named entities from visually rich documents. This is the first work that proposes a generalized approach for IE from visually rich documents and reports promising results on three het-HTML-specific features are used by the researchers to con-erogeneous datasets (refer to P1.2 in Section 1) for separate

PROBLEM FORMULATION & DEFINITION

our objective is to return a set of textual eleme $f_{int}(\mathbf{x}_i \subset B)$ from D, such that there is a one-to-one mapping between

VS2 proposes a two-phase approach for information extraction from D. First, we represent as a bag of visual areas Although not for IE, a computer vision based approach for B_1 , B_2 , ... B_V , such that $\{B_1, B_2, ...B_V\}$ denotes a partition of themselves. Identifying these visual areas helps define the description in Section 5.1) is more robust to diverse docu- VS2 searches for some predefined lexical and syntactic patterns (q_i) for each named entity within the context boundaries defined by these blockblence, for a set of patterns defined for the named entity $\eta \in \mathbb{N}$, the task of extracting the named entity n_i can be defined as finding a mapping into two sub-tasks $\mathfrak{I} = \mathfrak{I}_1 \circ \mathfrak{I}_2$.

First sub-task (\S): Find a partition $P = \{B_1, B_2, ...B_N\}$

Second sub-task (2): Once is obtained, find a map- $B_i \in P$, conforming to the pattern g

We propose VS2-Segment, a hierarchical page segmentasent the diverse visual areas within a visually rich document.

 $^{{}^{2}\}forall i, j, B i \subseteq B \text{ and } B \cap B = \emptyset, i, j$

It is discussed in greater details in the following section. The second sub-task is undertaken by VS2-Select, a distance supvision approach that searches for predetermined syntactic patterns within the context boundaries defined by the logical blocks and selects the most optimal matched pattern.

THE DOCUMENT LAYOUT MODEL

We define the layout model of a visually rich document as a nested tuple $\mathcal{C}(T)$, where C denotes the set of visual contents in D and T denotes the visual organization of D.

4.1 Visual content and their properties

An atomic element denotes the smallest unit of the visual content appearing in It can be classified into two major categories: textual and image element.

- 4.1.1 **Textual element** e smallest element in a document that has textual attributes. A textual contentan be represented as a nested tuple,= (text-data, color, width, heidht) Heretext-data and color represent the text appearing within a_t and the average color distribution (in LAB colorspace) of the visual area contained in respectively. The attributeshei∂ht andwidth denote the height and width of the smallest bounding box that encloses espectively. We deem a 'word' as the textual element of a document.
- 4.1.2 **Image element:** is an atomic element that represents an image content in the documentAn image element a_i in documentD is represented as a nested tuple, $a_i = (ima\partial e - data, width, hei\partial ht)$. Here, $ima\partial e - data$ denotes the image bitmap in height and width denote the height and width of the smallest bounding box that encloses a

We have used Tesserattl, an open-source document processing software to obtain the textual elements of a doc-block consists of a number of textual and image elements ument for this work.

4.2 Organization of the visual content

The visual organization of a document is represented as a segmentation algorithm called VS2-Segment. nested structure. Each visual area appearing in the document is represented by the smallest bounding box (say B_{ν}) that encloses it. Following this approach, we obtain the set of textual (1/4) and image (1/4) elements appearing in and represent as a string of bounding boxes enclosing the atomic elements $A_T \cup A_I$) appearing in B_V . In this work, we represent the visual organization of a document as a tree, $T_D = \{V, E\}$, where E denotes the edges and denotes the nodes of the treeAn edge between a parent node and its child denotes that the visual area represented by the child segmentation algorithm and subsequent information extracnode is enclosed by the visual area represented by the partion steps are described in details in the following sections. ent node. Therefore, the *non-leaf node* irrepresent the non-atomic visual areas that contain multiple smaller, seman-5.1 VS2-Segment: Segmentation of visually tically diverse elements within themselves. In other words, they are nested. A leaf node, on the other hand, corresponds The objective of VS2-Segment is to decompose a visually to the smallest visual areas which are visually isolated but rich document into logical blocks i.e. visual areas that are sesemantically coherent.

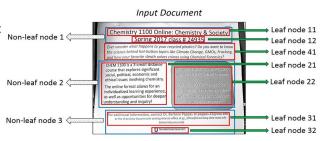


Figure 4: Each bounding box denotes a node in the lay out model of the academic event poster

We represent each node in T_D as a nested tuple = $(B, x, y, width, hei\partial ht)$ wherex, y, width, and hei ∂ht denote the (x,y) coordinates of the left-topmost point and width and height of the smallest bounding box that encloses B denotes the set of atomic elements that appear within For a visual are a_n , a_n can be easily obtained by performing a reverse lookup in the list of atomic elements (U A₁) appearing in D. The resulting layout tre \overline{e}_D , defined this way, not only encodes the visual and semantic properties of different visual areas appearing in the document, it also captures the hierarchical relationship between them. An illustration of the layout model for an academic event poster is shown in Fig.4. The layout tree is generated by a page segmentation algorithm, employed by VS2t will be described in more details in Section 5.1.

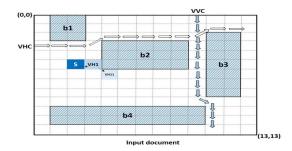
Takeaways. We propose a hierarchical layout model to represent the visual organization of visually rich documents in a unified way. The leaf-nodes of the tree-like structure represent the logical blocks of the document. Each logical of the document, provided that they are semantically coherent. We derive the layout model using a hierarchical page

OVERVIEW OF VS2

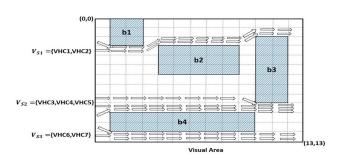
VS2 operates in two phases. In the first phase, a visually rich document is encoded as a bag of logical blocks, using VS2-Segment, a hierarchical page segmentation algorithm. In the next phase, VS2-Select searches for a set of lexical and syntactic patterns defined for each named entity, within the context boundaries of the logical blocks. We use distance supervision to learn a set of syntactic patterns for each named entity, using an isolated text-only corpus as the training dataset. The

rich documents

mantically coherent and isolated from each other. Following



ments in the document; S is a whitespace position at \$2,5)\/H 1 denotes a valid 1-hop horizontal movement from S;VH 1 → V H11 denotes a valid 2-hop horizontal movement from arrowtraces labeled &HCandVVC denote a valid horizontal and vertical cut from positions (0,3) and (10,0) respectively.



(a) £, £2, £3, and 4 are the bounding boxes of atomic content ele- (b) & 1 & 2 and V_S 3 denote sets of consecutive valid horizontal cuts; {b1, b2, b3}, {b3, b4}, and {b4} denote the set of neighboring bounding boxes and (0,2), (0,8), and (0,11) are the starting positions of V_{5 1} V_{5 2}and V_{5 3}respectively.

Figure 5: Illustrative example of the terminologies used in this work

the definitions introduced in Section 4.1, we represent each visual area by the smallest bounding box that encloses it. A bounding box $b \in B$ is defined as follows $b = (x_b, y_b, w_b, h_b)$, where X_b , Y_b denote the x-y coordinate of the top-left corner and W_b , h_b denote the width and height of the bounding box. In the following section, we will define a few key terms used in the segmentation algorithm before describing the algorithm itself in Section 5.1.2.

5.1.1 **Definitions:**

- If (x, y)denote the coordinates of a position on document D such that (x, y) < b, $\forall b \in \mathcal{B}$ where B is defined above, then (x, y) is called a **whitespace position**. $\forall b_1, b_2 \in V$, an edge exists betwee b_1 and b_2 , only if b_2 is
- If (x, y)and(x+1, y) are two whitespace positions A. then avalid horizontal movement (x, y) exists. If (x, y)'s a whitespace positio((x+1, y)) is not a whitespace position but either of (x+1)/(+1) and (x+1)/(-1)is a whitespace position in, then also availd hor**izontal movement**m (x, y)to that whitespace position exists. A valid vertical movement from, vito a whitespace position between, +1), (x+1, y+1) or (x-1, y+1)1, y+1) in D, can be defined in the same way **walld** horizontal vertical movement (x, y)s also referred to as a valid 1-hop movement.
- If a horizontal movement fron(x, y)to any one of the positions (x+1, y-1), (x+1, y) and (x+1, y+1) is valid and there also is a valid horizontal movement originating from that position, then valid 2-hop horizontal **movement** om the position(x, y)exists. Extending this definition, avalid k-hop horizontal vertical **movement** from (x, y) can be easily defined.
- For a document with height and width W, if a valid W-hop horizontal movement from (0, y), $y \in \emptyset$, H-1] exists, then **ahorizontal cut**iginating from (0, y)

exists. Similarly, if a valid H-hop vertical movement from (x, 0), $x \in [0, W-1]$ exists, then **vertical cub**m (x, 0) exists.

Illustrative examples of the terms introduced above are presented in Fig. 5.a and Fig. 5.b. Grid lines represent the rectangular coordinate system with the origin at left-top corner.

5.1.2 The Segmentation Algorithm

Let, T = (V, E) denotes the layout tree of a visually rich document D, where V denotes the set of atomic and non-atomic elements appearing i and E denotes the set of edges representing pairwise relationships between these elements. completely contained within 1. We hypothesize that a visually rich document is a nested object comprised of smaller semantically coherent visual areas, called logical blocks. The objective of our segmentation algorithm is to identify the logical blocks of diverse visually rich documents in a generalizable way. This is achieved by recursively decomposing a document into smaller visual areas by identifying the explicit and implicit visual modifiers used to augment/highlight an area within a visually rich document. A set of empirically selected low-level visual and semantic features are used to encode each area for this purposte layout tree acts as a unified data structure during the segmentation process. If a visual area () in D, represented by the node () in T, is segmented into a set of smaller areas, $v_2, ...v_t$, then nodes n_i , $\forall i \in [1, t]$ are added as child nodes n_i in T. The same steps are again repeated for these newly added nodes n_i , $\forall i \in [1, t]$ as more nodes representing visual elements in *Vi* are added as child nodes? ofto *T*. At each iteration of the segmentation algorithm, the leaf nodes of frepresent a set of isolated visual areas. Each notion T is represented as a nested tuple(v_c , v_t), where $v_c \subset C$ denotes the set of atomic

Visual Attribute	Description			
centroid-position	Position of the bbox centroid			
height	Height of the bounding box			
color	Average color in LAB col-			
	orspace			
angular distance	Angular distance of the bbox			
	centroid from origin			
sum of angular distance	SThe sum of angular dis-			
	tances between two bbox			
	centroids			

elements within the visual are and vt represents the complete sub-tree of with n_V as root. After convergence, the visual areas represented by the leaf nodes of onstructed this way, represent the logical blocks of the document. Each iteration of the segmentation algorithm involves identifying the explicit and implicit visual modifiers within a visual area in the document, followed by a semantic merging step.

At every iteration, the algorithm begins by searching for explicit visual delimiters within a visual area. Each visual areaV is scanned from top to bottom and left to right to identify sets of consecutive valid horizontat (s) and/or vertical cuts V_s) (refer to Fig. 5.b) that may act as potential visual separators for semantically diverse visual elements appearing in V. If such separators exist, the visual area is divided into smaller areas delimited by those separators. For example, if $V_{\rm S,1}$ and $V_{\rm S,2}$ are visual separators in Fig. 5.b, the visual area is divided into three smaller are \$4\$, \$\nu_2\$, and \$\nu_3\$, containing the bounding boxe\$1, \$\psi_2\$, \$\mathbb{B}\$} and \$\psi_4\$ respectively. Whether a set of consecutive, valid horizontal or vertical cuts should be con-8: sidered as a visual separator is decided using Algorithm 1. As-9: suming that, (a) distribution of the inter-area distance between 10: textual elements is different from the distribution of intra-area separation, and (b) font size is uniform within a semantically 11: coherent area, this algorithm scans for irregularities in the dis-12: tribution of correlation between the width (cardinality of the set of consecutive valid cuts) of a set of consecutive valid cuts₁₄. and the maximum height of its neighboring bounding boxes in a topologically sorted order. A neighboring bounding box for a set of consecutive valid cuts is the bounding box which is at minimum Euclidean distance from the set of consecutive valid cuts (refer to Fig. 5.b). The correlation distribution between width and maximum neighboring bounding box heights for all consecutive valid cuts is scanned in a topolog- often leads to over-segmentation. Its effects are worse for ical order (left to right and top to bottom) as the set closest to the first inflection point of the distribution is identified to be a visual delimiter. Although the cut-based segmentation

Table 1: Visual features used for clustering described above identifies explicit visual delimiters such as whitespace separators, it fails to recognize the implicit modifiers such as proximity, negative space, alignment, balance and symmetry that are often used to augment or highlight the semantics of areas within a visually rich document. To address this, a clustering of visual elements withins introduced at this stage. If no visual delimiters are found at the end of the previous step, each atomic element is encoded using a set of low-level features and grouped into clusters based on pairwise similarity. The features used for this purpose are empirically selected and shown in Table 1. To initialize the clustering process, a 22 equal-partition grid is assumed on V and one atomic element from each cell of the grid is selected as the cluster center. We choose the atomic elements as cluster-center which are at the minimum average distance from the rest of the atomic elements in each grid ceAt each iteration of the clustering step, pairwise distances are computed for each cell and the atomic elements and b_2 are assigned to the same clustefbif b) is the closest neighborpair in the encoding space that is not visually separated by another atomic element. The clustering step terminates when no new element can be assigned to a different cluster.

Algorithm 1 Identification of visual delimiters in D

```
1: procedure segment(S, B)
        S = \{s_1, s_2, ...s_n\}
                                   ▷ S:Consecutive valid cuts
        B = \{b_1, b_2, ...b_1\}
                                   ▶ B:Textual elements in D
 3:
        width = \Phi
 4:
        for i = 1 to m do width_i = |s| \times \frac{ar\partial max \ _k \ (hei\partial ht \ (nei\partial hbor -bbox \ _k \ (s)))}{ar\partial max \ _j \ (hei\partial ht \ (b_j))}
 5:
 6:
 7:
         Topologically sort S on (x,y) starting positions
        for i = 2 to m do
             W = \{width j, j \in [1, i-1]\}
             H = \{ar\partial max \ k(height(neidhbok(g))), j \in [1, i-
    1]}
             correlation_i = \rho(W, H)
         Sort s \in S on width in decreasing order
13:
        for i = 1 to m do
             C = C \cup correlation
15:
        t = inflection-point(i, correlation_i), t \in [1,m-1]
16:
        VD = \{s_1, s_{11}...s_m\}
                                       ▶ VD: Visual delimiters
17:
         return V D
```

We observed that this recursive segmentation process based on identifying visual delimiters, as described above, heterogeneous datasets. To address this issue we introduce a semantic merging operation in our workflow. The semantic contribution of textual elements within a visual area is computed for this purpose. If the semantic contributions of two visual areas are similar, they are merged together. The

³We derive the inflection points by solving $fo^{2}(f) = 0$, where f is the distribution of separator width vs. maximum neighboring-bbox-height

semantic contributior (a) of a visual area, represented by node n_i in the document layout tre \overline{e} is defined as follows:

$$SG = \Sigma_i \cos - similarity(n, \eta_i) - \Sigma_k \cos - similarity(n, \eta_k),$$
 (1)

In Eq. 1, $\forall j$, $n_j = siblin\partial(n_i)$, $\forall k$, n_k , $siblin\partial(n_i)$ and n_i , n_k denote nodes on the same level the layout tree (T = (V, E)). We have used a pre-trained Word2V266 Embedding to compute the cosine similarities in this work. If the semantic contribution of a node() is greater than a threshold⁴, it is merged with its sibling nod $\theta_{\mathcal{A}}$), with which it has the highest semantic similarity among all of its sibling nodes, provided that n_i and n_p are not visually separated. In other words, $\forall a \in V_{i,i} \Rightarrow a \in n_p$ or $a \in \eta$. Following this operation, nodes n_i and n_p are replaced by the merged node in the updated layout tree. The insight behind defining the semantic contribution of an area using Eqn. 1 is to ensure that each node in the layout tree represents a semantically distinct area within the document. An illustration of the logical blocks with respect to both the local and global context of where it appears in the document. The merging step terminates when no new nodes in the layout tree can be updated.

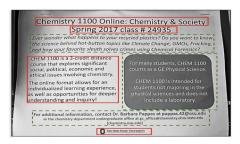


Figure 6: Each bounding box represents a logicalablood multiple matched patterns in the previous step, we in the academic event posterual areas enclosed propose an optimization-based entity disambiguation stratby solid red bounding boxes represent interest polyterior to these semantic operations, the transcribed text within the document

tion between distinct areas in the document by grouping elements that are visually similar and not separated by any visual delimiters. The explicit visual delimiters such as whiteof the document are identified during the beginning of every iteration. Implicit visual modifiers, on the other hand, are taken into consideration in the next phaseduring a bottom-up clustering stepWe observed that considering only the visual features during this process, often leads to over-segmentation. Therefore, a merging operation is undertaken that maximizes the semantic similarity between these notated text entries \mathcal{N}_i) for every named entity \mathcal{N}_i) related groups by merging the visual areas that are semantically similar. At the end of this recursive process, we obtain the separate datasets in this paper. layout-model T. Leaf-nodes of represent the logical blocks

Table 2: Constructing the holdout corpus

Datase	t Website	Query	Filter
D1	irs.gov	1988	1040
D2	allevents.in	NY	04/01-05/31
D2	dl.acm.org	Talks	Sorted by views
D3	fsbo.com	NY	None
D3	homesbyowner.com	NY	None

of the document. VS2-Segment does not assume any prior knowledge about the document template or format in any of its iterations, making it easier to be generalized for diverse document types. It is robust to rotation (upto)45 ind page artifacts that are common in many real-world scenarios. The logical blocks obtained from VS2-Segment helps define contextual boundaries, enabling effective semantic parsing obtained for an academic poster is shown in Fig.6.

5.2 VS2-Select: Information extraction from the logical blocks

Once the layout tree () has been constructed, the extraction task resolves to a search-and-select operation. For every named entity to be extracted, a set of lexico-syntactic patterns is searched within the text transcribed within the contextual boundaries defined by the logical blocks. In the following section, we will describe how these patterns are learned. Furthermore, as text-based disambiguation strategies do not work well for visually rich documents, to resolve conflicts

is normalized, its stopwords are removed, dependency trees Takeaways:VS2-Segment maximizes the visual separa- are constructed, and named entities are recognized. We have used publicly available natural language processing (NLP) tools for this purpose.

5.2.1 Learning the syntactic patterns

space separators across horizontal and vertical directions VS2-Select performs information extraction by searching for some predefined syntactic patterns within the context boundaries defined by the logical blocks. These are lexical and/or syntactic patterns learned from a holdout corput) (H is a readily annotated, structured, text-only corpus, constructed for the extraction task by scraping relevant public domain websites as a preprocessing step. Evidehtly onsists of anto that task i.e. $H = \sum_{i} (N_i, T_{N_i})$. We evaluate VS2 on three

Constructing the holdout corpusopulating the holdout corpusH with text entries for a named entitW₁ consists of four simple steps: (a) first, an expert identifies a public domain website(s) (using a web search engine) that maintain(s)

⁴For a layout tree of heigh?, the threshold paramete θ_b) is defined as follows, $\theta_h = \theta_{min} + \frac{\theta_{max} - \theta_{min}}{10} \times h$, where $\theta_{min} = 0$ and $\theta_{max} = 1$.

an indexed list of web pages where appear within a fixedformat HTML environment in diverse semantic contexts similar to the IE task, (b) second, select from the available filters may result in multiple matches. This is a known phenometo query the list such that the set of results returned is maximized; store the results to an HTML file, (c) extract the text T_{N_i} corresponding to N_i from all appearances N_i in the fixed-format HTML file using a custom web-wrappet9 and (d) finally, insert the tuple($\mathbf{s} \mathbf{h}_i$, T_{N_i}), $\forall i$ to H. For each N_i , tuples returned by querying the list were inserted Ho until the distribution of distinct syntactic patterns defined by the tuples in was approximately normal or there were no more tuples to be inserted. Holdout corpus for the first IE task contained 20 tables, each with two columns, an identifier of the named entity to be extracted and its corresponding field descriptor. The holdout corpus for the second IE task was constructed from the first 500 results obtained of an 'interesting' visual area in the document. More details from the search queries mentioned in Table 2. The corpus on identifying the interest points and the disambiguation consisted of a single table with two columns, an identifier for the named entity and its corresponding text. The holdout corpus for our third task was constructed in a similar fashion by collecting the top 100 results for each search query men-that is either visually prominent or semantically significant their corresponding IE tasks will be presented in Section 6.

Frequent sub-tree mining for learning the patterns: identify the syntactic patterns relevant to a named entity N_i , its corresponding entry in the holdout corpus T_{N_i} is annotated with a number of handcrafted lexical and syntactic in our work. These are selected empirically from a number features using publicly available NLP toolsirst, the text was chunked and dependency parse trees were obtained. ment or highlight the semantics of an area in a visually rich Named entities in every chunk were identified. The named entities with category 'Location' were further augmented with a geocode tag24. The noun POS tags were annotated with their respective Hypernym 42 senses. Verbnet β8 senses were extracted for every Verb POS tags as well. Once these features were extracted, the maximal frequent subtrees across the chunks were obtained. We used TreeManagera popular frequent subtree mining algorithm for this purpose. The syntactic patterns P_i) obtained this way represent the syntactic patterns for the named entit N_i . The patterns obtained this way for D2 and D3 are listed in Table 3 and 6. In case of D1, exact string match against the field descriptors in the holdout corpus was carried out.

Takeaways: A set of syntactic patterns are learned from a holdout corpus for each named entity to be extracted. This distance supervision approach circumvents the necessity of of that document. Interest points of an academic poster, obprior knowledge about the template or format of the document, a necessity in directly supervised approaches. This also Takeaways Interest points denote an optimal subset of helps avoid the curse of heterogeneity, making the proposed logical blocks obtained from the segmentation algorithm. approach easier to generalize for diverse document types.

5.3 Entity disambiguation by optimization

Searching for a syntactic pattern within the transcribed text non [16,34] in IE workflows. In these scenarios, traditional IE approaches employ word-sense disambiguati@filstrategies to rank all the matches using contextual information of where they appear in the document. Due to atypical visual properties, the traditional text-based techniques, however, do not work well for visually rich documents. Hence, we propose an optimization-based disambiguation strategy in this work. Every matched pattern is encoded using a set of visual and semantic descriptors. Disambiguation is performed by minimizing the distance between a match and its closest interest point in a multimodal encoding space. The key insight here is to prioritize those matches, that are in close proximity strategy will be discussed in the following sections.

5.3.1 Determining the interest points

An interest point 44 represents a visual area in the document tioned in Table 2. A detailed description of these datasets and both. We formulate this problem as an optimal subset selection B problem in this paper. Our objective is to select the most optimal subset from the set of all logical blocks) (obtained from the document. For a logical blo§k€(♦), we define 'optimality' using three visual and semantic objectives of commonly used visual or semantic modifiers used to augdocument. They are as follows:

- (1) maximizing the height of the bounding box enclosing larger font size is typically used to highlight significant areas in a visually rich document
- (2) maximizing semantic coherence i.e., the sum of pairwise cosine similarities between all text elements s', ∀s, s∈ \$, s, s'
- (3) minimizing the average word density; sparsely worded visual blocks covering a significant area of the document highlight semantically significant areas in a visually rich document

We solve the subset selection problem by non-dominated sorting [25] of the universal set of logical blocks obtained by VS2-Segment. The subset of logical blocks that constituted the first-order pareto-front, is selected as the interest points tained this way, are shown by red bounding boxes in Fig. 6.

⁵In multi-objective optimization paradigm, the pareto-front represents a state where the optimal value of one objective cannot be improved without worsening other objectives

We identify them by optimizing some visual and semantic properties, used to augment or highlight the semantics of a visual area in the document.

5.3.2 Distance based optimization

The semantics of a visually rich document is part of at least two modalities: textual and visual. Hence, to disambiguate among multiple matches found from the previous step, we encode every matched pattern using a set of visual and tex-Figure 7Sample documents from our experimental measure is computed using E2. The candidate which is closest to an interest point in the document, is selected as the racted from each of them will be presented in the following optimal match for that named entity. The features used for this purpose are empirically selected, similar to the features used to determine the interest points in a document. In the multimodal encoding space, the distarted between two visual areas s and c is defined as follows:

$$F_{s, c} = \alpha \Delta D(s, c) + \beta \Delta H(s, c) + \gamma \Delta Sim(s, c) + \nu \Delta Wd(s, 2)$$

where, $\alpha + \beta + \gamma + \nu = 1$ and α , β , γ , $\nu \in [-1]$. In Eqn. (3), $\Delta D(s, c)$ denotes the L1 distance between two centroids and est bounding-boxes enclosing the text-element§ amd C. $\Delta Sim(s, c)$ denotes the cosine similarity between text elements appearing within and and $\Delta W d(s, c)$ enotes the difference between distance-normalized word-densities of the smallest bounding-boxes enclosing and c respectively. The model parameters, β , γ and γ reflect the relative importance of visual saliency vs. textual verbosity in a document. For example, if the documents are not verbose but visually ornate (e.g. our second dataset), the $p_{\mathcal{B}}$, $\nu \geq y$. Similarly, if the corpus is not visually rich but verbose, then α , β , ν . For a balanced corpus (e.g. first and third datasets), it is safe^{collected} and prepared for this work. to assume $\alpha \approx \beta \approx \nu \approx \gamma$. A sample from each dataset is shown in Fig. 7.

of lexico-syntactic patterns within these blocks, for every named entity to be extracted. These patterns are learned for sponds to a form field in the document. A complete list of only corpus. To disambiguate between multiple matches, the https://s3.amazonaws.com/nist-srd/SD6/SD06_users_guide.pdf. distance between every match and its closest interest point is minimized in a multimodal encoding space. Eq. 2 is used to compute a weighted L1 distance between two visual areas in the document for this purpose.

EXPERIMENTS

We evaluate VS2 for three IE tasks on three separate datasets: NIST Tax dataset (D1), Event posters dataset (D2), and Realestate flyers dataset (D3). These datasets are heterogeneous i.e., the documents in these datasets originate from differ-



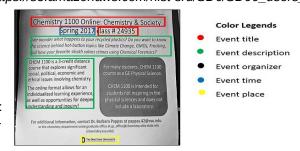
tual features and rank them based on their distances from thedatasets; Figures in (a), (b) and (c) represent documen closest interest point in that encoding space. The distance sampled from experimental datasets D1, D2 and D3

detailed description of the datasets and named entities exsection. We seek to answer three key questions in this study: (a) how does VS2 perform against a traditional text-only counterpart? (b) how does it compare against other state-ofthe-art methods?, and (c) what are the individual effects of various components used in VS2 on its end-to-end extraction quality? We answer the first two questions by following a two-phase evaluation strategy. First, we evaluate the performance of VS2-Segment (refer to Section 6.3) in accurately locating the positions of named entities within the document. ΔH (s, c) enotes the difference between heights of the small- End-to-end performance is evaluated by measuring the accuracy to classify (refer to Section 6.4) the named entities accurately. In both cases, we have compared against a textonly baseline and respective state-of-the-art methods. answer the third question, we perform an ablation study (refer to Section 6.5) to evaluate the individual effects of various components on the extraction quality for all three tasks.

6.1 Experimental datasets

We evaluate VS2 on three heterogeneous datasets (D1, D2 and D3) of visually rich documents. Datasets D2 and D3 were

NIST Tax datasetur first dataset (D1) is the NIST Tax dataset 33. It contains 5595 images of structured tax TakeawaysOnce the logical blocks have been obtained forms, representing 20 different form faces, all of which befor a visually rich document, VS2-Select searches for a setlong to the IRS 1040 package of 1988. The IE task defined for this dataset was to extract every named entity that correeach task, using distance supervision from an isolated text- the 1369 form fields defined for this dataset is available at:



ent sources, and/or belong to different types or formats. A Figure 8Ground-truth annotations of the academic event poster shown in Fig. 3

Table 3: Named entities extracted from D2

Named entity typ	pe Description	Syntactic patterns to search
Event Title	Short description or the event	(1) Verb phrase, (2) Noun phrase with numeric (<i>CD</i>) or textual modifiers (<i>JJ</i>), and (3) SVO
Event Place	Full address of the event	Noun phrases with valid geocode tags
Event Time	Time of the event	Noun phrases with valid TIMEX3 [5] tags
Event Organizer	Person/organization responsible for the e	even(t1) Verb phrase with captain/create/reflexive_appearance verb-senses3[8], (2) Noun phrase with Person/Organization as named entities
Event Description	Essential details of the event (what to end from the event if planning to attend, who be present)	expe \$ VO or Verb phrase or Noun phrase with modifiers (<i>CD/J</i> .) will

Table 4: Named entities extracted from D3

Named entity typ	e Description	Syntactic patterns to search
Broker Name	Full name of the listing broker	A bigram/trigram of NE's with Person / Organization 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.g4 beds,2,465 acres)	(1) Noun phrase with numeric (CD) or textual modifiers (JJ) and (2) Noun POS tags with senses measure / structure / estate in the Hypernym Tr42 [
Property Description	Mentions of the property type (e.g. building,floor,land/lot) and other essential details (e.g. parking, grocery	

Event posters data\$batesecond dataset (D2) is a col- this corpus are in HTML format. This IE task defined on D3 lection of event posters and flyers, advertising various local was to extract various attributes of the listed property. The uments, collected randomly from various sources, including used commonly accepted NLP lexicons/[to denote the local magazines, bulletin boards, and event hosting websites syntactic patterns representing each named entity. It contains both mobile captures of event flyers (1375 out of 6.2 Evaluation metrics 2190) as well as digital flyers in PDF format (815 out of 2190) Ground-truth construction (S2 is evaluated in two phases. ties that convey important event information. A complete list of the named entities is presented in Table 3. We have localization capabilities i.e., locating the position of a named used commonly accepted lexicon\$7 by NLP researchers to represent the syntactic patterns for each named entity.

Real-estate flyers dataset inal dataset (D3) comprises of online flyers containing commercial property listings in counties surrounding a major U.S. city. It was con-

and US national events. It contains a total of 2190 event doc-list of the named entities is presented in Table 3. We have

The performance of VS2-Segment is measured based on its entity in a document. Whereas, the end-to-end performance of VS2 is evaluated based on the accuracy of VS2-Select to correctly identify the named entity type, post localization. We evaluate VS2 against manually annotated ground-truth data. Every document in our experimental datasets was annotated structed by collecting 1200 commercial real-estate flyers fromby three experts. Annotation guidelines were developed and 20 different real-estate broker websites. The documents in the experts were asked to provide(a) coordinates of the smallest bounding boxes that contained a named entity in

Table 5: Evaluation of VS2-Segment on experimental datasets

Index	Algorithm	D1			2		D3
		Precision (%)Recall (%)	Precision (9	%)Recall (%)	Precision	(%)Recall (%)
A1 A2 A3 A4 A5 A6	Text-only XY-Cut Voronoi-tessellation VIPS Tesseract VS2-Segment	88.95 90.88 92.55 - 77.95 95.50	92.50 97.72 98.25 - 86.15 98.65	62.04 67.25 80.45 70.28 74.20 88.26	74.27 72.85 87.30 72.15 80.55 87.73	53.91 52.12 81.62 86.62 79.35 87.67	76.82 65.55 81.33 84.75 83.55 84.60

the document, and (b) a mapping between each bounding box and the named entity it contained. Annotations were performed using a specially designed web-based tool. The positional information from three experts was then averaged to derive the final coordinates. The final mapping between a bounding box and the named entity it contains was performed by majority voting among the tags assigned to that bounding box. An academic event poster from our second dataset (D2), annotated this way, is shown in Fig.8.

Two-phase evaluation We evaluate VS2-Segment by computing intersection-over-union (IoU) between the bounding box proposals and the corresponding ground-truth annotations. Following the benchmark proposed by Everingham et cursively segments an input document into smaller Voronoiat. [12] for evaluating visual object segmentation algorithms, areas. Summary statistics such as the distribution of font a bounding box proposal by VS2-Segment was deemed to besize, area ratio, angular distance are taken into consideration accurate if its IoU score against a labeled bounding box in for this purpose. VIPS by Cai et a4][(A4) exploits HTMLthe ground-truth data was greater than 0.65. The labels are specific features to identify visual delimiters that separate not considered at this stage o measure the end-to-end ex- visual blocks within an HTML document. All non-HTML doctraction performance, the predicted label for all localized and uments were converted to HTML format. Evidently, A4 could semantically classified named entities are compared against not be applied on dataset D1. Our final baseline method (A5) their corresponding ground-truth labels. If matched, the pro- is Tesseract4[1], an opensource document processing softposal is considered to be accurate. We report precision and ware that performs hierarchical layout analysis of an input recall values for both phases on all experimental datasets.

6.3 Evaluation of VS2-Segment

An evaluation of VS2-Segment's performance in accurately localizing the named entities for all experimental datasets has been presented in Table 5. Results show that it achieves, method, stemmed from under-segmentation of the logical satisfactory performance for all three IE tasks. We observe relatively better performance for D1, compared to D2 and D3. This can be attributed to higher structural variability in documents belonging to datasets D2 and D3. An exhaustiveWe evaluate the end-to-end performance of VS2 by measurerror-analysis of the final results also revealed that about 80% of the errors stemmed from over-segmentation of the VS2-Select within an input document post localization. For logical blocks due to low-quality transcription inhibiting semantic merging at later iterations of the algorithm.

Comparison against state-of-the-art methods: compare VS2-Segment against five contemporary page segthe text transcribed from each segmented area. Entity dismentation algorithms (refer to Table 5). Our first competitor (A1) is a text-based baseline method that groups words text-only entity disambiguation method. with similar word-embeddings into the same clusters. The second baseline (A2) is a visual segmentation algorithm

Table 6: End-to-end evaluation of VS2 on D2

Index	Named Entity	Pro _l Pr. (%)	posed n Rec. (%)	nethod)∆F1(%)
N1 N2 N3 N4 N5	Event Title Event Place Event Time Event Organizer Event Description	84.88 76.68 94.67 72.56 776.59	81.09 86.37 84.70 74.41 86.00	8.98 3.76 0.49 10.50 1.60
	Overall	81.08	82.51	5.07

that divides a document into smaller visual areas by finding vertical and/or horizontal cuts. Our third competitor (A3) redocument to segment it into blocks. Results show that we were able to outperform A1, A2, A3 and A5 on all datasets. We significantly outperformed A4 on dataset D12lost of the errors in the final segmentation result, for this baseline blocks that were not delineated by a rectangular whitespace separator. We observed competitive results on D3.

6.4 Evaluation of end-to-end performance

ing the accuracy of accurately classified named entities by each dataset, we compare the performance of our method against a text-only baseline. Using Tesser4dt to segment the input document, it searches for syntactic patterns within ambiguation is performed using Les® a state-of-the-art

Evaluation on D1: The objective of this IE task was to extract 1369 named entities corresponding to every form

Table 7: Comparison of end-to-end performance against existing methods on all datasets

Index	Algorithm		01		D2		D3
		Precision (%)Recall (%)	Precision ((%)Recall (%)	Precision	(%)Recall (%)
A1 A2 A3 A4 A5 A6	ClausIE FSM ML-based Apostolova et al ReportMiner VS2	85.0 - 92.20 96.50 95.25	90.75 - 96.25 100.0 98.4	70.65 77.25 83.92 85.25 51.25 88.05	62.19 79.05 81.0 85.66 62.50 85.95	76.50 84.50 92.65 87.28 67.75 91.80	68.05 82.95 86.40 90.42 80.40 90.32

field in D1. Results show that VS2 achieved an overall average precision o95.25%nd recall of98.4%or this dataset. Compared to the text-only baseline, we observed an overall improvement of 2.84% in average F1-score.

Evaluation on D2The objective of this task was to extract five different named entities from a corpus of visually rich event posters. The named entities and their corresponding lexico-syntactic patterns have been presented in Table 3. Table 6 presents the end-to-end evaluation of VS2 for this IE task. The final column in the table represents the average improvement in F1-score against the text-only baseline. Compared to its text-only counterpart, significant improvement in average F1-scores were observed for named entities 'Event errors introduced in our workflow for the document images Title' (8.98% in F1-score) and 'Event Organizer' (10.5% in F1-score). Whereas, marginal improvements were observed images also inhibits the semantic merging step at later iterfor 'Event Time' (0.41% Further inspection revealed that the text-only approach performed wellif: (a) the syntactic pattern defined for a named entity contained a regular expression with partial string matching capabilities; (b) there was only a single matched pattern for that entity before five contemporary information extraction methods. Results disambiguation. Both of these were true in this case...

Evaluation on D3The objective of this task was to extract six different named entities from a corpus of online realestate flyers. The named entities and their corresponding lexico-syntactic patterns are presented in TablePerformance of VS2 for this task is shown in Table 8. Compared traction task defined for dataset D1. Our second baseline is to the text-only baseline, significant improvements in average F1-scores were observed for the named entities 'Broker'ery named entity to be extracted, it finds the most frequent Name' (10.18%) and 'Property Address' (4.60%). Both of thesubtrees within the dependency trees for entries against named entities were among the most visually rich entities in D3. Smaller improvements were observed for the entities terns defined by these subtrees are then searched within the 'Broker Phone' and 'Broker Email'. We observed that for most transcribed text of a test document to identify the named documents in D3, these named entities appeared only once in entities. VS2 outperforms this method on all datasets. Our the document. Marginal improvement was observed for the entity 'Property Description' also. This was due to the low inter-annotator agreement on what constituted "essential information" of a listed property.

ment in performance using VS2 was statistically significant (t-that are in PDF formatfor a fair comparison against this test reveals p < 0.05) for all datasets. Results in Table 6 and method. Following this approach, an SVM based classifier Table 8 also reveal that end-to-end performance was better

Table 8: End-to-end evaluation of VS2 on D3

Index	Named Entity	Pro Pr.(%)	posed Rec. (%	method %) ∆ <i>F1 (%)</i>
N1 N2 N3 N4 N5 N6	Broker Name Broker Phone Broker Email Property Addres Property Size Property Desc.	94.72 96.15 97.25 892.68 85.25 84.75	90.85 82.25 95.40 85.50 93.05 94.90	10.18 1.63 2.56 4.60 3.37 0.74
	Overall	91.80	90.32	3.84

on D3 compared to D2. This is due to the over-segmentation in D2. Low-quality transcription of some of the document ations of our segmentation algorithm leading to incorrect semantic parsing, affecting the downstream extraction task.

Comparison against existing metWedsampare the end-to-end performance of VS2 on all datasets against of this study are shown in Table 7.

Our first competitor is ClausIE10, a text-only approach that constructs a set of clause-based rules for every named entity to be extracted. VS2 significantly outperforms ClausIE on both D2 and D3. It does not apply for the form field exa Frequent Subtree Mining (FSM)1[48] approach. For evthat named entity in the holdout corpus. The syntactic patthird competitor (Zhou et al.[49]) proposes a supervised machine-learning approach. Every non-HTML document needs to be converted to HTML format for this approach. Hence it could not be applied for the first dataset D1. Due Compared to the text-only baseline, the average improve- to similar reasons, we only consider those documents in D2

Table 9: Evaluating individual components in VS2 by ablation study

Index	1	VS2-Select	Δ	\F1 (%	<u>6)</u>	
	Visual feature	Semantic feature based mergi	n ∉ ntity disambiguation	D1	D2	D3
A1 A2 A3 A4	√ × √	× √ √	√ √ × Text-only	1.07 0.95	2.55 4.22 6.78 4.55	3.84 7.05

method on D2 and provides comparable performance on D3. textual boundaries obtained from the prior segmentation of Our fourth competitor is a multimodal IE approach proposed the document. The most useful insight gathered from this by Apostolova et al[2]. They proposed a combination of textual and visual features to train an SVM classifier. Results measures the effects of the proposed entity disambiguation show that the proposed approach outperforms this method for all tasks. This is attributed to better semantic parsing capabilities exhibited by VS2as it leveraged the context boundaries obtained from the prior segmentation step. Finally, we compared our method against ReportMin22[a commercially available, human-in-the-loop document workflow management tool. It allows its users to define custom masks for each named entity in the document. Information extraction is performed by manually selecting the most appropriate rule for a document. We randomly selected 60% of tion extraction from visually rich documents in this work. Usthe dataset to generate the rules and evaluated our perfor-ing a set of empirically selected low-level features to encode mance on the rest. Results show that this approach did not variability in document layouts increased. VS2 performed competitively or better on all datasets, with lesser human effort required in its end-to-end workflow.

6.5 Ablation study

the end-to-end extraction quality, we have performed an ab- end extraction quality. To the best of our knowledge his effect of a critical component in VS2 on overall F1-score of from heterogeneous visually rich documents and reports its the downstream extraction task. The final column in Table 9 quantifies the end-to-end effect of these changes on the over tend our work to incorporate more complex documents. For all average F1-score. A1 investigates the effects of semantic example, the assumption of font size similarity within each merging in VS2-Segment algorithm on the overall extraction quality. Results show that although this affects the overall F1-score of all datasetiss effects are most prominent for datasets D2 and D3. This is attributed to over-segmentation of our workflow would improve the robustness of our method of the documents, adversely affecting the localization of named entities within a document. Scenario A2 investigates the effects of incorporating visual features for segmenting a visually rich document. Similar to A1, this leads to imprecise localization of the logical blocks, contributing to poor overall F1-scores for all of our datasets improved results by incorporating visual features also establish our design

was trained on the dataset (60%-40% split) using some visualchoice of a two-phase IE method for visually rich documents. and textual features of the document. VS2 outperforms this Better overall F1-scores are achieved by leveraging the conablation study, however, stems from scenarios A3 and A4. A3 strategy on the end-to-end extraction quality. Significant effects of this simulation are observed for all datasets. We have also compared our disambiguation strategy against Lesk [a popular text-based entity disambiguation method. Experimental results revealed significant improvements over the text-based disambiguation method for datasets D2 and D3.

CONCLUSION

We have proposed VS2, a generalized approach for informaeach visual area, a hierarchical segmentation algorithm is perform well for D2 and D3. Performance worsened as the proposed to divide each document into logical blocks. Named entities are extracted by following a distantly supervised search-and-select method within the contextual boundaries defined by these logical block \$.52 is evaluated on three heterogeneous datasets for separate IE tasks. Results suggest that careful consideration of visual and semantic features To investigate the effects of individual components in VS2 on can outperform current state-of-the-art methods in end-tolation study in this section. Each row in Table 9 measures theis the first work that proposes a generalized approach for IE performance on three IE tasks. In future, we would like to exblock in our current implementation can be addressed by introducing a generalizable feature to identify font-type. Addressing the issue of transcription errors during both phases towards processing real-world documentsxtending our feature library to include sophisticated contextual semantic features (e.g. n-gram features), learning to weight each feature based on observed data, language-agnostic multimodal embedding to encode each document, would further increase the robustness of our method. We also have plans to extend this work on multilingual and nested documents in future.

REFERENCES

- [1] Gabor Angeli, Julie Tibshirani, Jean Wu, and Christopher D Manning. 2014. Combining distant and partial supervision for relation extraction. In Proceedings of the 2014 conference on empirical methods in natura [21] Wei Liu, Xiaofeng Meng, and Weiyi Meng. 2010 ide: A vision-based language processing (EMNLP). 1556-1567.
- [2] Emilia Apostolova and Noriko Tomuro. 201@ombining Visual and Textual Features for Information Extraction from Online Flyers.. In EMNLP. 1924-1929.
- [3] Satanjeev Banerjee and Ted Pedersen. 20102 adapted Lesk algorithm for word sense disambiguation using WordNet. In International [23] Tomohiro Manabe and Keishi Tajima. 20 Estracting logical hierarconference on intelligent text processing and computational linguistics. Springer, 136-145.
- [4] Deng Cai, Shipeng Yu, Ji-Rong Wen, and Wei-Ying Ma. 20063: a vision-based page segmentation algorith(2003).
- [5] Angel X Chang and Christopher D Manning. 20 Sutime: A library for recognizing and normalizing time expressions.. In LREC, Vol. 2012. 3735-3740.
- and Tapan S Parikh. 2012hreddr: pipelined paper digitization for low-resource organizations. In Proceedings of the 2nd ACM Symposium on Computing for Development. ACM, 3.
- [7] Laura Chiticariu, Yunyao Li, Sriram Raghavan, and Frederick R Reiss. [27] Mike Mintz, Steven Bills, Rion Snow, and Dan Jurafsky. 2018 Pant 2010. Enterprise information extraction:recent developments and open challenges. In Proceedings of the 2010 ACM SIGMOD International Conference on Management of data. ACM, 1257-1258.
- [8] Thomas H Cormen, Charles E Leiserson, Ronald L Rivest, and Clifford Stein. 2009Introduction to algorithmsMIT press.
- [9] Valter Crescenzi Giansalvatore Mecca Paolo Merialdo, et al. 2001. Roadrunner: Towards automatic data extraction from large web sites. In VLDB, Vol. 1. 109-118.
- [10] Luciano Del Corro and Rainer Gemulla. 20 CauslE: Clause-based Open Information Extraction. In Proceedings of the 22Nd International [29] David Nadeau and Satoshi Sekin2007. A survey of named entity Conference on World Wide Web (WWW '13). ACM, New York, NY, USA, 355-366.https://doi.org/10.1145/2488388.2488420
- [11] AnHai Doan, Jeffrey F Naughton Raghu Ramakrishnan Akanksha Baid, Xiaoyong Chai, Fei Chen, Ting Chen, Eric Chu, Pedro DeRose, Byron Gao, et a2009. Information extraction challenges in managing unstructured data ACM SIGMOD Record 37, 4 (2009), 14-20.
- [12] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. 2010The pascal visual object classes (voc)
- [13] Emilio Ferrara, Pasquale De Meo, Giacomo Fiumara, and Robert Baumgartner. 2014. Web data extractionapplications and techniques: A survey. Knowledge-based systems 70 (2014), 301-323.
- [14] Ignazio Gallo, Alessandro Zamberlettiand Lucia Noce2015. Content extraction from marketing flyers. In International Conference on Computer Analysis of Images and Patterns. Springer, 325–336.
- [15] Wolfgang Gatterbauer Paul Bohunsky, Marcus Herzog, Bernhard Krüpl, and Bernhard Pollak.2007. Towards domain-independent information extraction from web tables. In Proceedings of the 16th international conference on World Wide Web. ACM, 71-80.
- [16] Jing Jiang. 2012nformation extraction from text.In Mining text data. Springer, 11-41.
- [17] Dan Jurafsky. 200 Speech & language process Pregarson Education India.
- [18] Mukkai Krishnamoorthy, George Nagy, Sharad Sethand Mahesh Viswanathan, 1993Syntactic segmentation and labeling of digitized pages from technical journald EEE Transactions on Pattern Analysis and Machine Intelligence 15, 7 (1993), 737-747.
- [19] Nicholas Kushmerick. 2000/rapper induction: Efficiency and expressivenessArtificial Intelligence 118, 1-2 (2000), 15-68.

- [20] Alberto HF Laender, Berthier A Ribeiro-Neto, Altigran S Da Silva, and Juliana S Teixeira. 2002. brief survey of web data extraction tools. ACM Sigmod Record 31, 2 (2002), 84-93.
- approach for deep web data extraction. IEEE Transactions on Knowledge and Data Engineering 22, 3 (2010), 447-460.
- [22] Astera LLC. 2018ReportMiner: A Data Extraction Solutionhttps: //www.astera.com/products/report-miner. (2018)ccessed: 2018-09-
- chical structure of HTML documents based on headingsoceedings of the VLDB Endowment 8, 12 (2015), 1606-1617.
- [24] Google Maps. 201&Google Maps Apihttps://developers.google.com/ maps. (2018)Accessed: 2018-09-30.
- [25] R Timothy Marler and Jasbir S Arora. 2008 urvey of multi-objective optimization methods for engineering. Structural and multidisciplinary optimization 26, 6 (2004), 369-395.
- [6] Kuang Chen, Akshay Kannan, Yoriyasu Yano, Joseph M Hellerstein, [26] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corradoand Jeff Dean. 2013. Distributed representations of words and phrases and their compositionality. In Advances in neural information processing systems, 3111-3119.
 - supervision for relation extraction without labeled data. In Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP: Volume 2-Volume 2. Association for Computational Linguistics, 1003-1011
 - [28] Marcin Michał Mirończuk. 2018The BigGrams: the semi-supervised information extraction system from HTML: an improvement in the wrapper induction. Knowledge and Information Systems 54, 3 (2018), 711-776.
 - recognition and classification Linguisticae Investigationes 30, 1 (2007). 3-26.
 - [30] Roberto Navigli. 2009Word sense disambiguation: A surveyCM computing surveys (CSUR) 41, 2 (2009), 10.
 - [31] Dat PT Nguyen, Yutaka Matsuo, and Mitsuru Ishizuka. 2008 ation extraction from wikipedia using subtree mining. In Proceedings of the National Conference on Artificial Intelligence, Vol. 22. Menlo Park, CA; Cambridge, MA; London; AAAI Press; MIT Press; 1999, 1414.
 - challenge. International journal of computer vision 88, 2 (2010), 303- [32] Feng Niu, Ce Zhang, Christopher Ré, and Jude W Shavlik. 2042p-Dive: Web-scale Knowledge-base Construction using Statistical Learning and Inference VLDS 12 (2012), 25-28.
 - [33] National Institute of Standards and Technology. 2018ST Special Database 6https://www.nist.gov/srd/nist-special-database-6. (2018). Accessed: 2018-09-30.
 - [34] Benjamin Roth, Tassilo Barth, Michael Wiegand, Mittul Singh, and Dietrich Klakow. 2014Effective slot filling based on shallow distant supervision methodsarXiv preprint arXiv:1401.1158 (2014).
 - [35] Sunita Sarawagi et a 2008. Information extraction. Foundations and Trends® in Databases 1, 3 (2008), 261-377.
 - [36] Ritesh Sarkhel, Nibaran Das, Aritra Das, Mahantapas Kundu, and Mita Nasipuri. 2017. A multi-scale deep quad tree based feature extraction method for the recognition of isolated handwritten characters of popular Indic scripts Pattern Recognition 71 (2017), 78-93.
 - [37] Ritesh Sarkhel, Nibaran Das, Amit K Saha, and Mita Nasipuri. 2016. multi-objective approach towards cost effective isolated handwritten Bangla character and digit recognition Rattern Recognition 58 (2016),
 - [38] Karin Kipper Schuler. 2005. VerbNet: A broad-coverage, comprehensive verb lexicon.(2005).

- [39] Asif Shahab, Faisal Shafait, and Andreas Dengel. 2000DAR 2011 robust reading competition challenge 2: Reading text in scene images. [45] Sen Wu, Luke Hsiao, Xiao Cheng, Braden Hancock, Theodoros Rekatsi-In Document Analysis and Recognition (ICDAR)11 International Conference on. IEEE, 1491-1496.
- [40] Samuel Sanford Shapiro and Martin B Wilk1965. An analysis of variance test for normality (complete samples)Biometrika 52,3/4 (1965), 591-611.
- [41] Ray Smith. 2007An overview of the Tesseract OCR engine. In Document Analysis and Recognition, 2007. ICDAR 2007. Ninth International Conference on, Vol. 2. IEEE, 629-633.
- [42] Rion Snow, Daniel Jurafsky, and Andrew Y Ng. 2005. Learning syntactic patterns for automatic hypernym discoveryln Advances in neural information processing systems. 1297-1304.
- [43] Fei Sun, Dandan Songand Lejian Liao. 2011. Dom based content extraction via text density. In Proceedings of the 34th international ACM SIGIR conference on Research and development in Information Retrieval. ACM, 245-254.
- [44] Tinne Tuytelaars, Krystian Mikolajczyk, et. 2008. Local invariant feature detectors: a surveyFoundations and tren@sin computer graphics

- and vision 3, 3 (2008), 177-280.
- nas, Philip Levis, and Christopher Ré. 20 Randuer: Knowledge base construction from richly formatted data.In Proceedings of the 2018 International Conference on Management of Data (SIGMOD)M,
- [46] Yudong Yang, Yu Chen, and Hong, Jiang Zhang 2003. HTML page analysis based on visual cuds. Web Document Analysis: Challenges and Opportunities. World Scientific, 113-131.
- [47] Mohammed J Zaki. 200 Efficiently mining frequent trees in a forest. In Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining. ACM, 71-80.
- [48] Xiaowen Zhang and Bingfeng Cher2017. A construction scheme of web page comment information extraction system based on frequent subtree mining. In AIP Conference Proceedings, Vol. 1864. AIP Publishing, 020059.
- [49] Ziyan Zhou and Muntasir Mashuq2014. Web Content Extraction Through Machine Learning(2014).