# TRIE: End-to-End Text Reading and Information Extraction for Document Understanding

Peng Zhang[1]
Yunlu Xu[1]
Zhanzhan Cheng[21]
Shiliang Pu[1]
Jing Lu[1]
Liang Qiao[1]
Yi Niu[1]
Fei Wu[2]

1. Hikvision Research Institute, Hangzhou, China
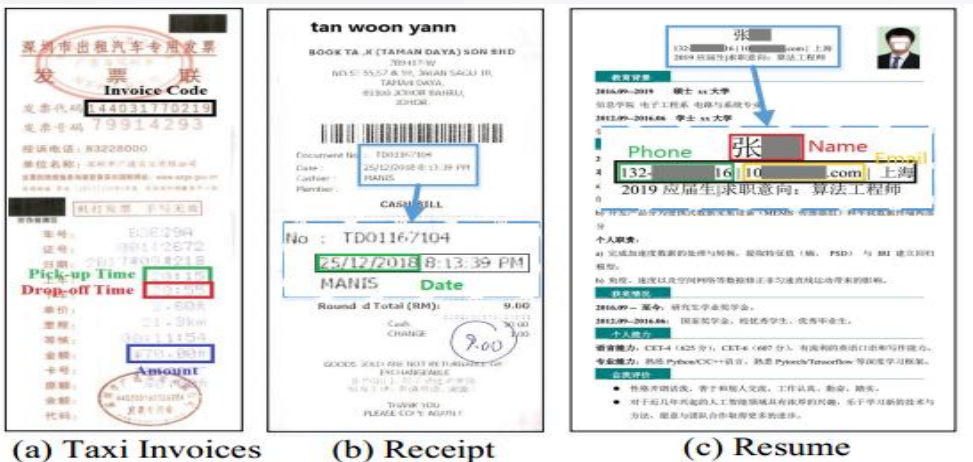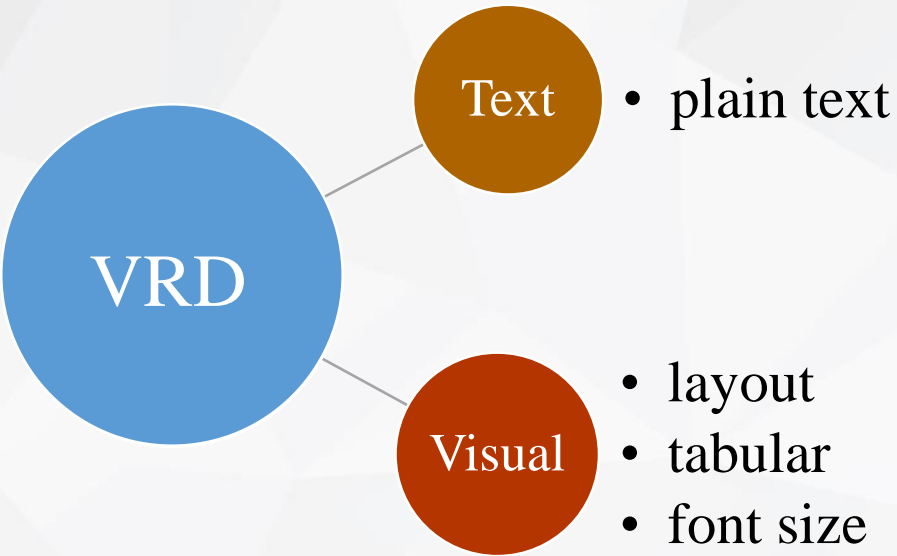2. Zhejiang University, Hangzhou, China

# Background

◆ **VRD (Visually Rich Document)**

VRD

Text
• plain text

Visual
• layout
• tabular
• font size

(a) Taxi Invoices  (b) Receipt  (c) Resume

◆ **VRD Understanding**

**VRD Understanding**

Receipt.

VRD Understanding Algorithm

| Name | OJC MARKETING SDN BHD |
| Date | 15/01/2019 |
| Total | 193.00 |

Entities to extract.

# Background

◆ **Problem of current framework**

◆ **Motivation**



**Limitation 1: Limited visual features in IE.**
Keep：
$x_0,y_0...,x_3,y_3$, '15/01/2019'
$x_0,y_0...,x_3,y_3$, '193.00'

**Lost：**
Font, Color, Layout etc.

**Limitation 2: Ignoring relations between OCR & IE.**

**Advantage 1: Multimodal fusion in IE.**
Keep：
$x_0,y_0...,x_3,y_3$, '15/01/2019'
$x_0,y_0...,x_3,y_3$, '193.00'

**Font, Color, Layout etc.**

**Advantage 2: Bridging OCR & IE,**
**Forward:** OCR boost IE
**Backward:** IE boost OCR

# Method

◆ **Overall Architecture.**

# Experiment

## ◆ Datasets

| Dataset | Training | Testing | Entities | Layout | Text Type |
|---|---|---|---|---|---|
| Taxi Invoices | 4000 | 1000 | 9 | Fixed | Struct |
| SROIE | 626 | 347 | 4 | Variable | Struct |
| Resumes | 1978 | 497 | 6 | Variable | Semi-struct |

Dataset Statics



(a) Taxi Invoices   (b) Receipt   (c) Resume

Dataset Examples

## ◆ Performance Summary

| Entities | Chargrid(TR) | NER(TR) | GCN(TR) | Our Model |
|---|---|---|---|---|
| Code | 89.4 | 94.5 | 97.0 | **98.2** |
| Number | 85.3 | 92.4 | 93.7 | **95.4** |
| Date | 89.8 | 82.5 | 93.0 | **94.9** |
| Pick-up time | 82.9 | 60.0 | **86.3** | 84.6 |
| Drop-off time | 87.4 | 81.1 | 91.0 | **93.6** |
| Price | 93.0 | 94.5 | 93.6 | **94.9** |
| Distance | 92.7 | 93.6 | 91.4 | **94.4** |
| Waiting | 89.2 | 85.4 | 91.0 | **92.4** |
| Amount | 80.2 | 86.3 | 88.7 | **90.9** |
| Avg | 87.77 | 85.59 | 91.74 | **93.26** |

Taxi Invoices Dataset

| Setting | Model | F1-Score |
|---|---|---|
| Setting 1: Prediction of bboxes and transcript of texts | Chargrid(TR) | 78.24 |
| | NER(TR) | 69.09 |
| | GCN(TR) | 76.51 |
| | Our model | **82.06** |
| Setting 2: Groundtruth of bboxes and transcript of texts | Character-Word LSTM [24] | 90.85 |
| | LayoutLM[54] | 95.24 |
| | PICK[58] | 96.12 |
| | Our model | **96.18** |

ICDAR2019 SROIE Dataset

# Experiment

◆ **Performance Summary**

| Entities | Chargrid(TR) | NER(TR) | GCN(TR) | Our Model |
|---|---|---|---|---|
| Name | 43.4 | 42.7 | 42.9 | **45.7** |
| Phone | 87.0 | 86.6 | 83.3 | **88.0** |
| E-mail | 70.9 | 69.6 | 68.0 | **74.9** |
| Edu-period | 77.1 | 68.7 | 62.2 | **81.4** |
| University | 74.7 | 86.0 | 82.3 | **87.4** |
| Major | 72.1 | 80.4 | 78.7 | **80.8** |
| Avg | 70.87 | 72.33 | 69.57 | **76.3** |
| Speed(fps) | 1.13 | 1.69 | 1.62 | **1.76** |

Resumes Dataset

Algorithm Architecture.

◆ **Discussion**

| | | | | |
|---|---|---|---|---|
| Text feat only | √ | √ | √ | √ |
| Textual Context feat | | √ | | √ |
| Visual Context feat | | | √ | √ |
| Accuracy | 74.33 | 92.30 | 92.70 | **93.26** |

Effects of multimodal features on IE.

| Text Reading Results | | IE Model | Accuracy |
|---|---|---|---|
| TR only | End-to-End (TRIE) | | |
| √ | | GCN [30] | 91.70 |
| | √ | GCN [30] | 92.60 |
| | √ | TRIE | 93.26 |

Effects of E2E framework on text reading.

| Datasets | Layers | Heads | | | |
|---|---|---|---|---|---|
| | | 2 | 4 | 8 | 16 |
| Taxi Invoices | 1 | 92.97 | 92.98 | 92.86 | 92.72 |
| | 2 | 93.00 | 92.98 | **93.26** | 92.71 |
| | 3 | 92.55 | 92.81 | 93.06 | 92.83 |
| Resumes | 1 | 75.20 | 75.21 | 75.47 | 74.53 |
| | 2 | 75.62 | 76.25 | 76.28 | 75.86 |
| | 3 | 75.55 | 75.74 | **76.35** | 76.35 |

Effects of layers and heads in textual context block.