# Problem Statement
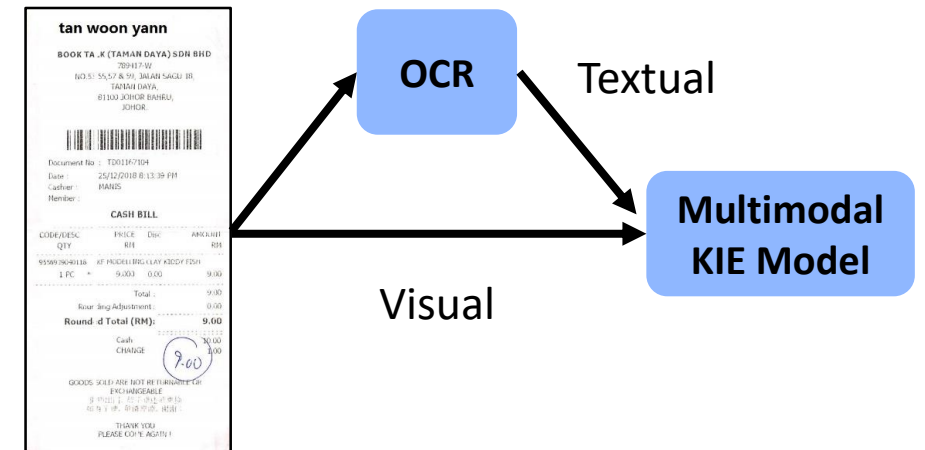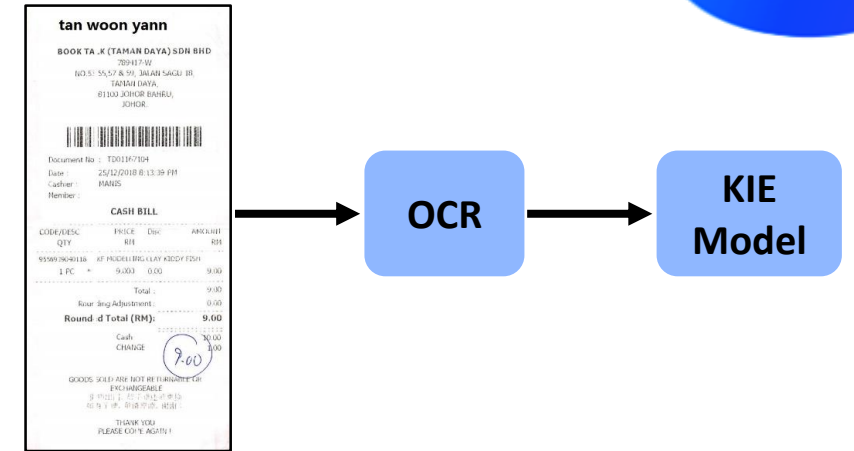
- **Key Information Extraction (KIE)** from document images is the automated process of retrieving relevant data presented visually in the document.

- **Optical Character Recognition (OCR)** tools are used to extract the text in the document.

# Problem Statement

- Using the OCR text, many studies model the KIE task as a **sequence tagging** problem and **solve using NER** (Yu et. al, 2021).

- **Multimodal** approaches combine the OCR output with visual information of the document image and position information of the tokens.

Yu, W., Lu, N., Qi, X., Gong, P., Xiao, R.: Pick: processing key information extraction from documents using improved graph learning-convolutional networks. In: 2020 25th International Conference on Pattern Recognition (ICPR). pp. 4363– 4370. IEEE (2021)

# Document Types

| Fixed-form / structured documents | Semi-structured documents | Unstructured documents |
|---|---|---|
| • Surveys<br>• Questionnaires<br>• Tests<br>• Claim forms | • Invoices<br>• Purchase orders<br>• Bills of lading<br>• EOBs | • Contracts<br>• Letters<br>• Articles<br>• Notes |

# Previous Methods and Limitations

## BiLSTM-CRF (Huang et. al, 2015)

- Pros:
  - Effective in sequence tagging and NER tasks.
  - Uses both past and future tokens.
  - Capture sentence level meaning.

- Cons:
  - Uses only textual information.
  - Unable to utilize layout or visual information directly.

**Works best on unstructured documents.**

Huang, Z., Xu, W., Yu, K.: Bidirectional lstm-crf models for sequence tagging. arXiv preprint arXiv:1508.01991 (2015)

# Previous Methods and Limitations

## Chargrid (Katti et. al, 2018)

- Pros:
    - Converts a page into 2D grid of chars.
    - Encodes spatial features.
    - Better in structured documents.

- Cons:
    - Using only char level information is not sufficient for understanding the semantics.
    - Thus, not as effective in unstructured documents as in structured ones.

**Raw data**

**Chargrid**

Katti, A.R., Reisswig, C., Guder, C., Brarda, S., Bickel, S., Höhne, J., Faddoul, J.B.: Chargrid: Towards understanding 2d documents. arXiv preprint arXiv:1809.08799 (2018)

# Previous Methods and Limitations

## BERTgrid (Denk et. al, 2019)

- Pros:
  - Converts a page into 2D grid of contextualized word embeddings obtained from BERT.
  - Possible to understand the semantics of document.

- Cons:
  - BERT language model is frozen during training.
  - Image of the document page is not directly utilized.



Raw image



BERTgrid representation

Denk, T.I., Reisswig, C.: Bertgrid: Contextualized embedding for 2d document representation and understanding. arXiv preprint arXiv:1909.04948 (2019)

# Previous Methods and Limitations

## ViBERTgrid (Lin et. al, 2021)

- Pros:
  - Incorporates BERTgrid with a CNN to process the raw image directly.
  - Jointly trains BERT and CNN.



- Cons from our observations:
  - Could not outperform a pure textual model (BiLSTM-CRF on BERT emb.) on unstructured money transfer order documents.

**Works best on structured documents.**

Lin, W., Gao, Q., Sun, L., Zhong, Z., Hu, K., Ren, Q., Huo, Q.: Vibertgrid: a jointly trained multi-modal 2d document representation for key information extraction from documents. In: International Conference on Document Analysis and Recognition. pp. 548–563. Springer (2021)

# Idea

What if we combine the best performing models?

# Proposed Approach: ViBERTgrid BiLSTM-CRF

# Proposed Approach: ViBERTgrid BiLSTM-CRF

# Datasets

## SROIE

- ICDAR SROIE dataset[1]: 973 receipts (626 training, 347 testing samples).
- Semi-structured documents.
- Four entity types: company, date, address, total.
- Dataset presents key information fields and OCR output separately.
- Previous studies[2,3] used text-based matching but results in poor matching.
- We manually annotated the entire dataset on the word-level.
- We publicly release the word-level annotations of SROIE
  dataset for use in multimodal transformers.
  (https://github.com/YKT-NLP/ICDAR-2019-SROIE-Token-Level-Annotations)
- Evaluation on test set is still problematic due to discrepancies
  between OCR output and key information fields, e.g.,
  mismatched punctuation, extra or missing white spaces, typos etc.
- Some of these errors have been documented and manually fixed.

[1]Huang, Z., Chen, K., He, J., Bai, X., Karatzas, D., Lu, S., Jawahar, C.V.: Icdar2019 competition on scanned receipt ocr and information extraction. In: 2019 International Conference on Document Analysis and Recognition (ICDAR). pp. 1516–1520 (2019). https://doi.org/10.1109/ICDAR.2019.00244
[2]Lin, Z.: Vibertgrid pytorch (2021), https://github.com/ZeningLin/ ViBERTgrid-PyTorch
[3]Delplace, A.: Chargrid model : Extraction of meaningful instances from document images (2020), https://github.com/antoinedelplace/Chargrid

# Datasets

**Transactional Documents**

- In-house dataset consists of unstructured Turkish money transfer order documents, introduced by Oral et. al (2022).

- Has two sets: **U**nstructured **T**ransactional **D**ocuments (**UTD**) and **U**nstructured **M**ulti-**T**ransaction **D**ocuments (**UMTD**).

- UTD has 3500 documents (2500 for training, 400 validation, 600 testing).

- UMTD has 1154 documents (954 for training, 200 testing).

- Within the UMTD test set, 54 out of 200 documents have tabular-like layouts (**TLL**), the rest has non-tabular-like (**noTLL**) documents.

- We used the same splits as in Oral et. al (2022) for consistency.

Oral, B., Eryiğit, G.: Fusion of visual representations for multimodal information extraction from unstructured transactional documents. International Journal on Document Analysis and Recognition (IJDAR) pp. 1–19 (2022)

# Experiments & Results

**SROIE**

**Table 1.** Performance comparison on SROIE.

| Model | Macro F1 Score (%) |
|---|---|
| ViBERTgrid | $93.56_{\pm 0.005}$ |
| ViBERTgrid BiLSTM-CRF | $\mathbf{93.85}_{\pm 0.003}$ |

# Experiments & Results

**Transactional Documents**

| Train Set | Test Set | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | UTD$_{test}$ | | UMTD$_{test}$ | | | | | | |
| | All | | All | | noTLL (73%) | | TLL (27%) | | |
| | w/o | w/ | w/o | w/ | w/o | w/ | w/o | w/ | |
| UTD$_{train}$ | $90.95_{\pm0.57}$ | $\mathbf{92.19_{\pm0.18}}$ | $91.04_{\pm0.33}$ | $\mathbf{92.42_{\pm0.14}}$ | $91.73_{\pm0.40}$ | $\mathbf{93.03_{\pm0.19}}$ | $89.61_{\pm0.38}$ | $\mathbf{91.12_{\pm0.30}}$ | |
| | $87.76_{\pm0.53}$ | $\mathbf{89.57_{\pm0.39}}$ | $86.71_{\pm0.82}$ | $\mathbf{89.29_{\pm1.03}}$ | $85.83_{\pm1.15}$ | $\mathbf{88.26_{\pm0.75}}$ | $85.53_{\pm0.72}$ | $\mathbf{87.97_{\pm0.70}}$ | |
| UTD$_{train}$ + UMTD$_{train}$ | $91.05_{\pm0.30}$ | $\mathbf{92.04_{\pm0.50}}$ | $93.28_{\pm0.34}$ | $\mathbf{93.98_{\pm0.29}}$ | $93.41_{\pm0.31}$ | $\mathbf{94.13_{\pm0.33}}$ | $93.11_{\pm0.52}$ | $\mathbf{93.70_{\pm0.70}}$ | |
| | $87.61_{\pm0.33}$ | $\mathbf{89.09_{\pm0.53}}$ | $90.15_{\pm0.71}$ | $\mathbf{91.78_{\pm0.26}}$ | $88.54_{\pm0.66}$ | $\mathbf{89.69_{\pm0.74}}$ | $90.66_{\pm0.65}$ | $\mathbf{91.77_{\pm0.57}}$ | |

vanilla ViBERTgrid

micro F1

macro F1

ViBERTgrid with BiLSTM-CRF (ours)

# Experiments & Results

**Transactional Documents**

| Train Set | UTD$_{test}$ All | | UMTD$_{test}$ All | | noTLL (73%) | | TLL (27%) | |
|---|---|---|---|---|---|---|---|---|
| | w/o | w/ | w/o | w/ | w/o | w/ | w/o | w/ |
| UTD$_{train}$ | 90.95$\pm$0.57 | **92.19$\pm$0.18** | 91.04$\pm$0.33 | **92.42$\pm$0.14** | 91.73$\pm$0.40 | **93.03$\pm$0.19** | 89.61$\pm$0.38 | **91.12$\pm$0.30** |
| | 87.76$\pm$0.53 | **89.57$\pm$0.39** | 86.71$\pm$0.82 | **89.29$\pm$1.03** | 85.83$\pm$1.15 | **88.26$\pm$0.75** | 85.53$\pm$0.72 | **87.97$\pm$0.70** |
| UTD$_{train}$ + UMTD$_{train}$ | 91.05$\pm$0.30 | **92.04$\pm$0.50** | 93.28$\pm$0.34 | **93.98$\pm$0.29** | 93.41$\pm$0.31 | **94.13$\pm$0.33** | 93.11$\pm$0.52 | **93.70$\pm$0.70** |
| | 87.61$\pm$0.33 | **89.09$\pm$0.53** | 90.15$\pm$0.71 | **91.78$\pm$0.26** | 88.54$\pm$0.66 | **89.69$\pm$0.74** | 90.66$\pm$0.65 | **91.77$\pm$0.57** |

Oral et. al (2022) obtained a micro and macro F1 score of **91.48** and **88.45** on UTD dataset using a textual only BiLSTM-CRF model pretrained on BERT embeddings.

# Experiments & Results

**Transactional Documents**

| Train Set | Test Set | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | UTD$_{test}$ | | UMTD$_{test}$ | | | | | |
| | All | | All | | noTLL (73%) | | TLL (27%) | |
| | w/o | w/ | w/o | w/ | w/o | w/ | w/o | w/ |
| UTD$_{train}$ | $90.95_{\pm0.57}$ | $\mathbf{92.19}_{\pm0.18}$ | $91.04_{\pm0.33}$ | $\mathbf{92.42}_{\pm0.14}$ | $91.73_{\pm0.40}$ | $\mathbf{93.03}_{\pm0.19}$ | $89.61_{\pm0.38}$ | $\mathbf{91.12}_{\pm0.30}$ |
| | $87.76_{\pm0.53}$ | $\mathbf{89.57}_{\pm0.39}$ | $86.71_{\pm0.82}$ | $\mathbf{89.29}_{\pm1.03}$ | $85.83_{\pm1.15}$ | $\mathbf{88.26}_{\pm0.75}$ | $85.53_{\pm0.72}$ | $\mathbf{87.97}_{\pm0.70}$ |
| UTD$_{train}$ + UMTD$_{train}$ | $91.05_{\pm0.30}$ | $\mathbf{92.04}_{\pm0.50}$ | $93.28_{\pm0.34}$ | $\mathbf{93.98}_{\pm0.29}$ | $93.41_{\pm0.31}$ | $\mathbf{94.13}_{\pm0.33}$ | $93.11_{\pm0.52}$ | $\mathbf{93.70}_{\pm0.70}$ |
| | $87.61_{\pm0.33}$ | $\mathbf{89.09}_{\pm0.53}$ | $90.15_{\pm0.71}$ | $\mathbf{91.78}_{\pm0.26}$ | $88.54_{\pm0.66}$ | $\mathbf{89.69}_{\pm0.74}$ | $90.66_{\pm0.65}$ | $\mathbf{91.77}_{\pm0.57}$ |

# Experiments & Results

**Transactional Documents**

| Train Set | Test Set | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| | $\text{UTD}_{test}$ | | $\text{UMTD}_{test}$ | | | | | |
| | All | | All | | noTLL (73%) | | TLL (27%) | |
| | w/o | w/ | w/o | w/ | w/o | w/ | w/o | w/ |
| $\text{UTD}_{train}$ | $90.95_{\pm 0.57}$ | $\mathbf{92.19}_{\pm 0.18}$ | $91.04_{\pm 0.33}$ | $\mathbf{92.42}_{\pm 0.14}$ | $91.73_{\pm 0.40}$ | $\mathbf{93.03}_{\pm 0.19}$ | $89.61_{\pm 0.38}$ | $\mathbf{91.12}_{\pm 0.30}$ |
| | $87.76_{\pm 0.53}$ | $\mathbf{89.57}_{\pm 0.39}$ | $86.71_{\pm 0.82}$ | $\mathbf{89.29}_{\pm 1.03}$ | $85.83_{\pm 1.15}$ | $\mathbf{88.26}_{\pm 0.75}$ | $85.53_{\pm 0.72}$ | $\mathbf{87.97}_{\pm 0.70}$ |
| $\text{UTD}_{train} + \text{UMTD}_{train}$ | $91.05_{\pm 0.30}$ | $\mathbf{92.04}_{\pm 0.50}$ | $93.28_{\pm 0.34}$ | $\mathbf{93.98}_{\pm 0.29}$ | $93.41_{\pm 0.31}$ | $\mathbf{94.13}_{\pm 0.33}$ | $93.11_{\pm 0.52}$ | $\mathbf{93.70}_{\pm 0.70}$ |
| | $87.61_{\pm 0.33}$ | $\mathbf{89.09}_{\pm 0.53}$ | $90.15_{\pm 0.71}$ | $\mathbf{91.78}_{\pm 0.26}$ | $88.54_{\pm 0.66}$ | $\mathbf{89.69}_{\pm 0.74}$ | $90.66_{\pm 0.65}$ | $\mathbf{91.77}_{\pm 0.57}$ |

# Conclusion

- We focused on the impact of using a multimodal transformer (i.e., ViBERTgrid previously explored on semistructured documents) on the NER task from **unstructured financial documents**.

- The initial results showed that the original ViBERTgrid has a **negative impact on unstructured** documents compared to a pure textual baseline.

- We presented an approach to enhance the performance of ViBERTgrid on unstructured documents by **extending it with a BiLSTM-CRF layer.**

- As a result, our proposed **ViBERTgrid BiLSTM-CRF** model demonstrated a significant improvement in performance (**up to 2 percentage points**) on unstructured documents, while maintaining its performance on semi-structured documents, in the **domain of financial and banking documents**.

- As an additional contribution, we publicly released token-level annotations for the SROIE dataset to pave the way for its use in multimodal sequence labelling models.

# Q & A ?