# DOCUMENT READING WITH A FIELD-CROPPING VISION-TEXT TRANSFORMER
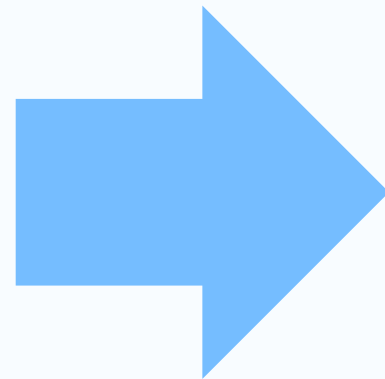
By

Pardheev Krishna Tammineni

# PROBLEM OVERVIEW

- Structured forms contain printed and handwritten text
- Need to crop each field and recognize its content accurately
- Challenges: inconsistent labels, varying handwriting

# DATA PREPARATION

- Load image/label pairs and crop fields using coordinates from JSON annotations
- Store cropped images with their text labels for training
- Clean duplicate or mislabeled fields before training
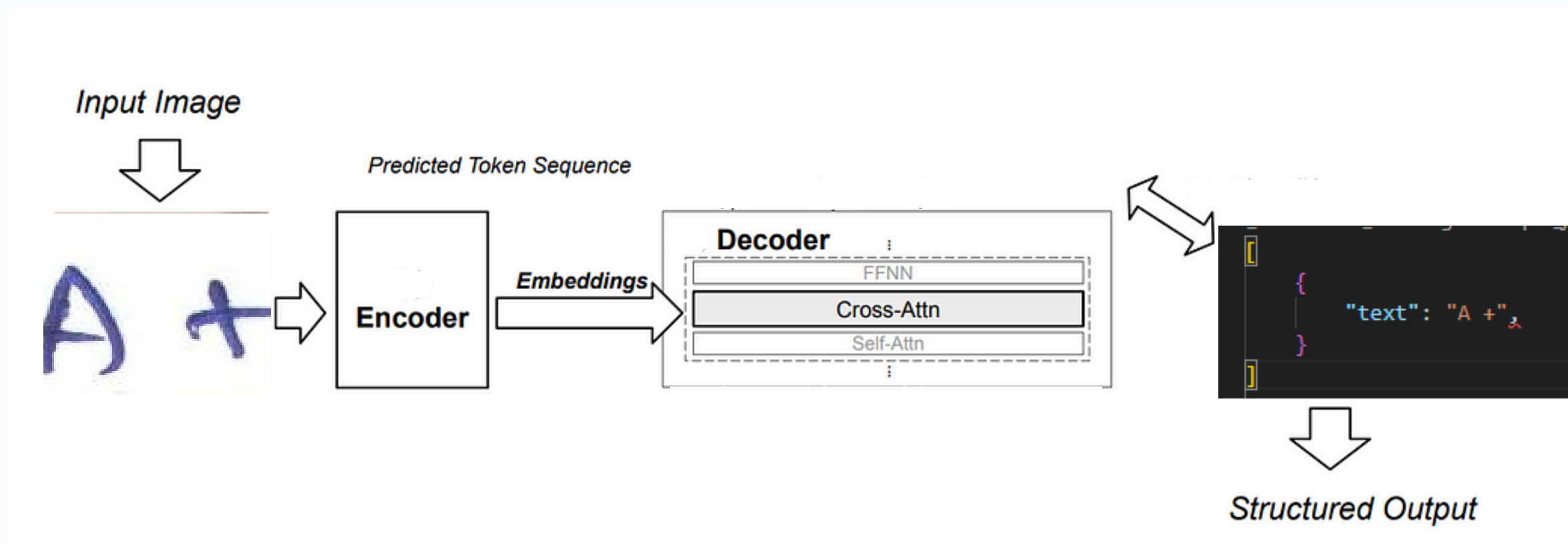
# MODEL ARCHITECTURE

- End-to-end transformer using a vision encoder (Swin) and a text decoder (BART)
- Input images resized to 512 × 512
- Encoder converts visual features into embeddings; decoder generates text from those embeddings

# CUSTOMIZED TRAINING

- Freeze half of the encoder layers to fit limited hardware
- Define maximum sequence length and special tokens
- Track Character Error Rate (CER), Word Error Rate (WER), field accuracy, and document accuracy

# TRAINING PIPELINE

- Split dataset: 90 % for training, 10 % for validation
- Optimizer: AdamW with a linear learning-rate schedule
- Automatic logging of loss and metrics each epoch

Load image-label pairs from disk → Parse JSON annotations → Crop fields into separate images → Split dataset (90% train, 10% val) → Vision encoder (Swin, first half frozen) → Text Decoder (BART) → Save best checkpoint on composite score improvement

# VALIDATION METRICS

- After each epoch, report CER, WER, field accuracy, and document accuracy
- Save best checkpoint when the composite score improves
- Mismatches logged for error analysis



CER and WER over Epochs

# TESTING & EVALUATION

- Load the best checkpoint for evaluation on the held-out test set
- Print CER, WER, field accuracy, document accuracy, and final score

**Metrics on Train Data**

CER: 0.16 %

WER: 1.06 %

Field Accuracy: 97.99 %

Document Accuracy: 67.53 %

Final Score: 94.40

# HOW THIS APPROACH DIFFERS

- Traditional OCR pipelines separate detection and recognition
- This approach uses an end-to-end transformer to generate text sequences directly from cropped fields
- Cropping each field simplifies the task and improves accuracy
- Freezing encoder layers reduces memory use

| Model | Params (M) | Inference Time (ms) | CER (%) | WER (%) | Field Accuracy (%) | Document Accuracy (%) |
|---|---|---|---|---|---|---|
| LayoutLMv3 | 228 | 150 | 1.8 | 5.0 | 90.2 | 52.3 |
| DocFormer | 150 | 200 | 1.5 | 4.5 | 92.1 | 55.8 |
| TrOCR | 140 | 180 | 1.2 | 3.8 | 94.5 | 60.1 |
| My Approach | 100 | 120 | 0.16 | 1.06 | 97.99 | 67.53 |

# CHALLENGES & NEXT STEPS

- Continue cleaning mislabeled or inconsistent data
- Improve multilingual label consistency
- Explore augmentations (rotation, noise) to enhance robustness

# CONCLUSION

- Tailored training of the vision-text transformer effectively extracts handwritten fields from structured forms
- Combining field cropping with targeted metrics yields strong accuracy
- Future work: expand the dataset and prepare for real-world deployment

# GITHUB REPO

# THANK YOU