SAMSUNG SDS

Foresee

# Techtonic 2021

Partner

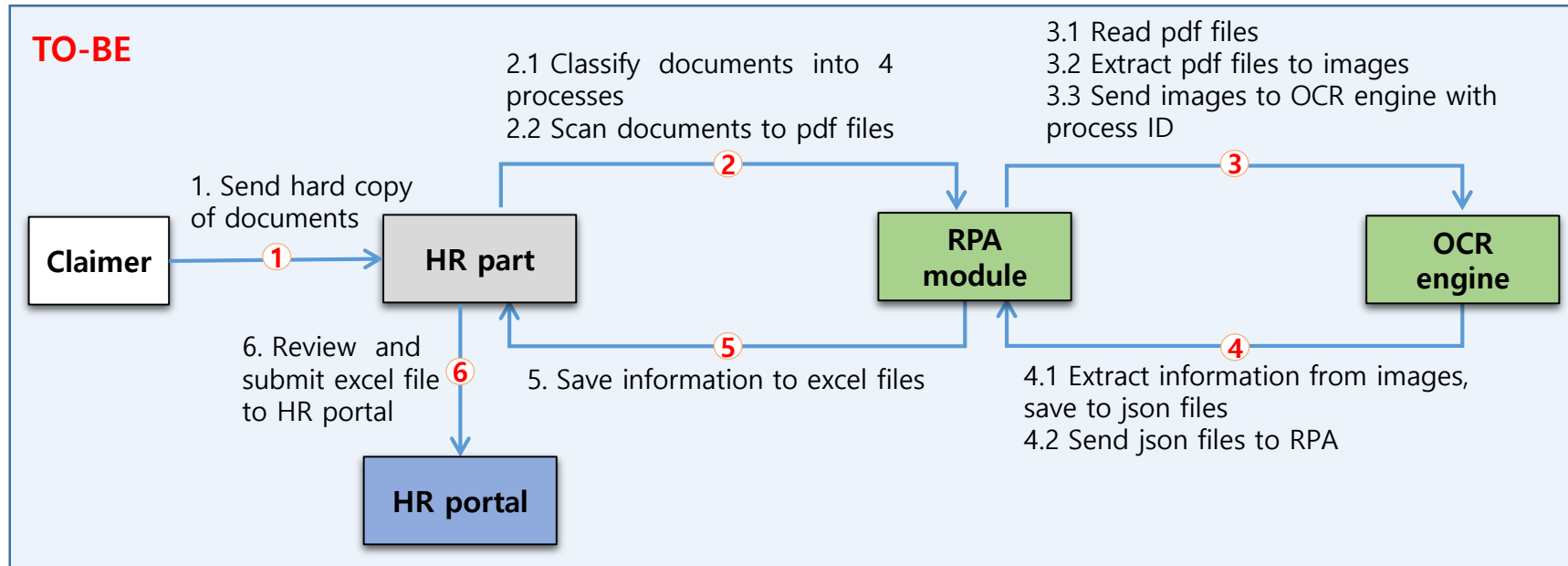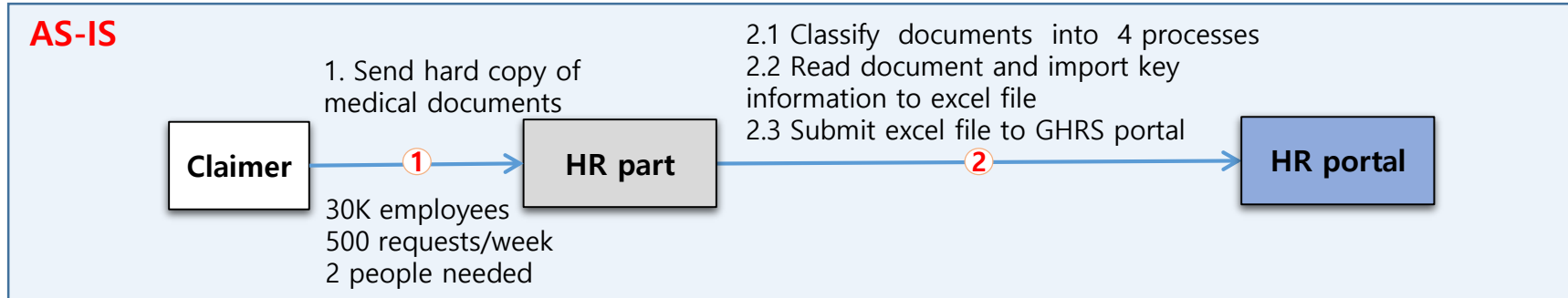Disrupt

Samsung SDS Research Vietnam

# An approach in Vietnamese handwritten OCR

Ngo T Dat

# Use case: data input automation



**AS-IS**

Claimer → **HR part** → **HR portal**

1. Send hard copy of medical documents
(①)

30K employees
500 requests/week
2 people needed

2.1 Classify documents into 4 processes
2.2 Read document and import key information to excel file
2.3 Submit excel file to GHRS portal
(②)

**TO-BE**

Claimer → **HR part** → **RPA module** → **OCR engine**

1. Send hard copy of documents
(①)

2.1 Classify documents into 4 processes
2.2 Scan documents to pdf files
(②)

3.1 Read pdf files
3.2 Extract pdf files to images
3.3 Send images to OCR engine with process ID
(③)

4.1 Extract information from images, save to json files
4.2 Send json files to RPA
(④)

5. Save information to excel files
(⑤)

6. Review and submit excel file to HR portal
(⑥)

**HR portal**

# Financial sector: Invoice, application form

# Limitations of the traditional OCR approach

- Bad image samples

The quick brown fox jumps over the lazy dog
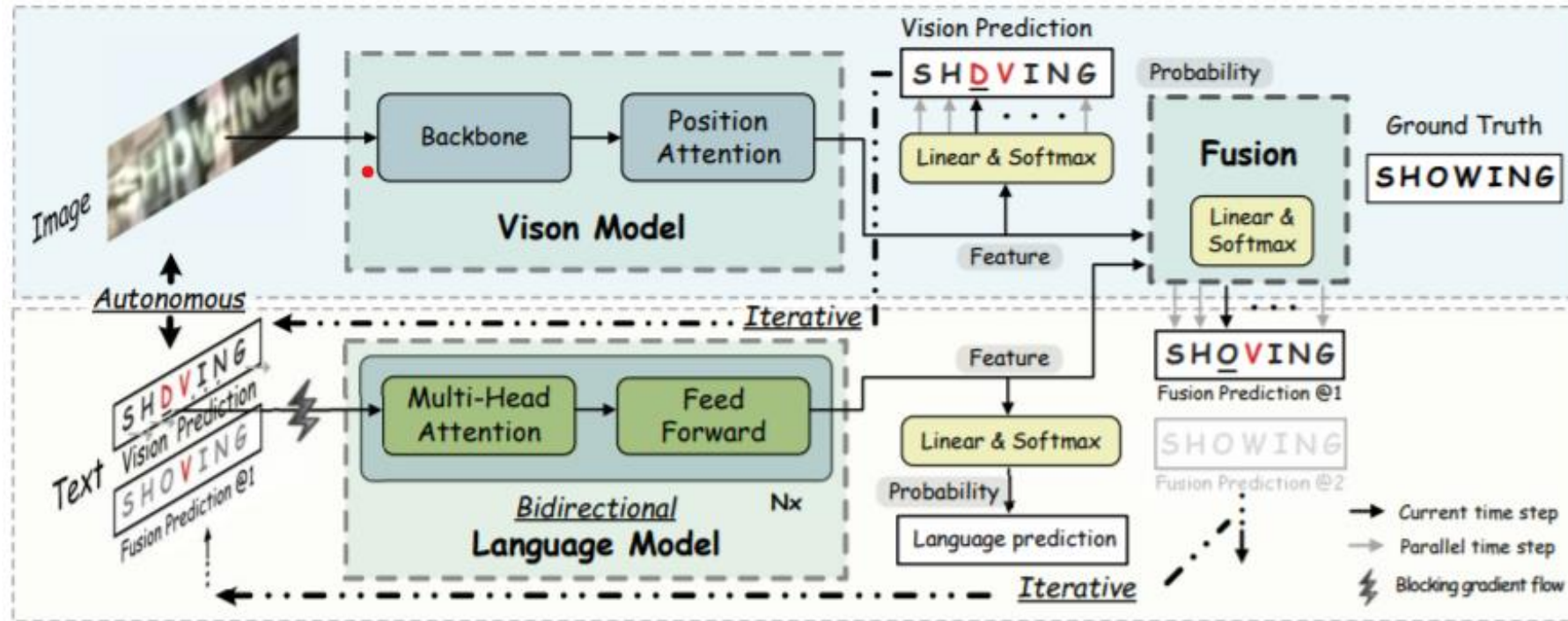
The quick brown fox jumps over the lazy dog

accountably cravenly say chagrining hino

- No knowledge about grammar rules or frequently used words.
- Do not understand semantic meaning of words.

⇒ Using Language Model to aid OCR Text Recognition Model can be a good solution.

# Overall architecture

The ABINet architecture has demonstrated a feasible way to ensemble Vision Model (VM) and Language Models (LM).
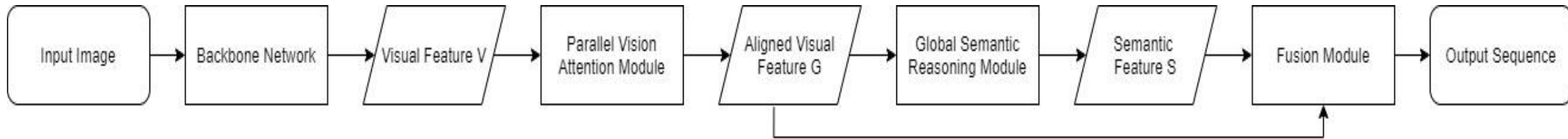


ABINet Architecture
Fang et. al., 2021. "Read Like Humans: Autonomous, Bidirectional and Iterative Language Modeling for Scene Text Recognition". In 2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)

# Phase 1 – Vision Model & Language Model
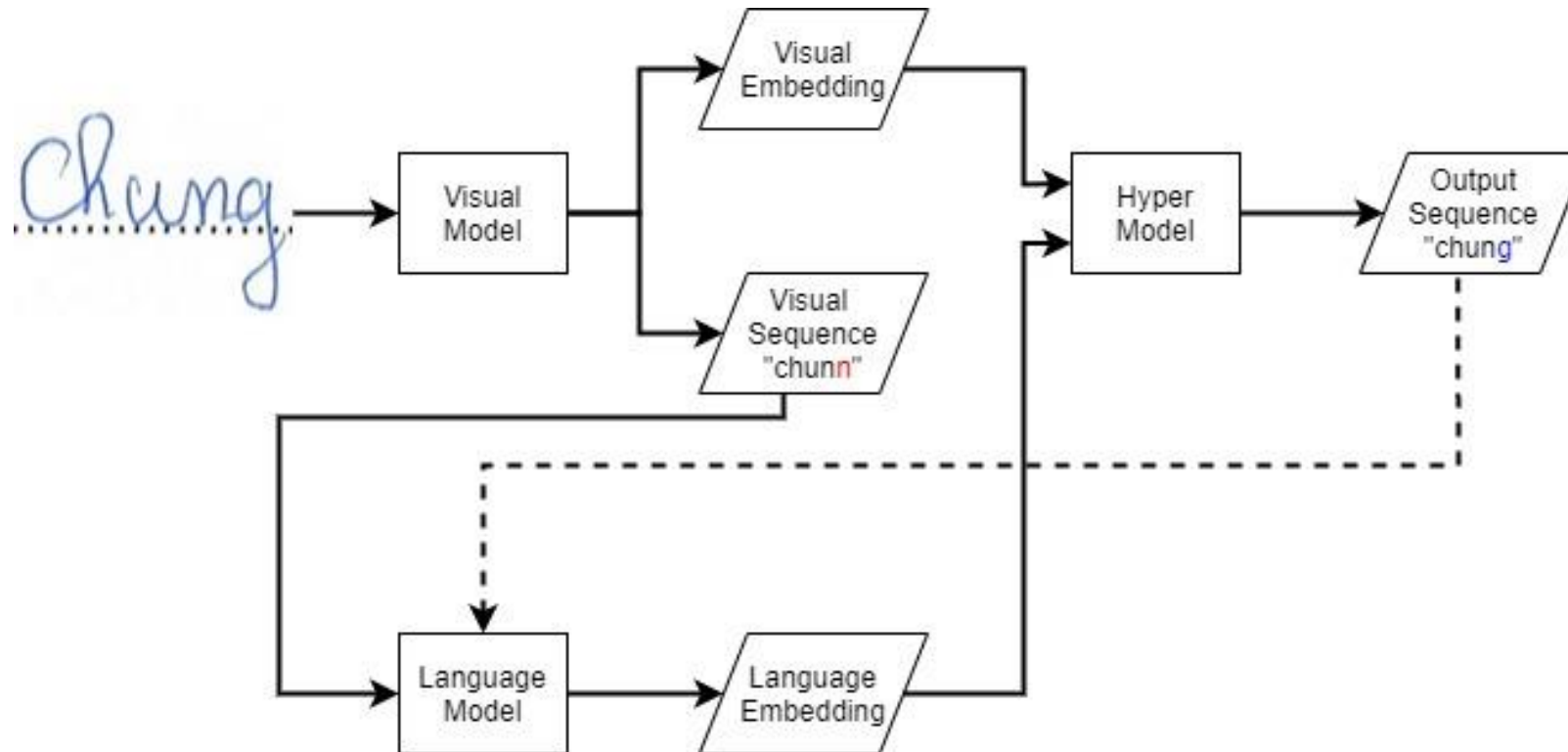
## Vision Model – SRNet



SRNet
Yu et. al., 2020. "Towards Accurate Scene Text Recognition with Semantic Reasoning Networks". In 2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)

## Language Model – Transformers

- Introduced in 2017
- A ground breaking achievement in NLP field and now expanding to computer vision.

# Phase 2 - Ensemble models



- Training procedure includes two phases.
- In phase 2, two models are jointly trained and a hyper-model (HM) was introduced to stack two phase 1 models' embedding vectors.
- Inference procedure can be single forwarded or iteratively forwarded.

# Single word benchmark

- Training data:
  - Two datasets: Vietnamese corpus (50K words, 26M samples), Vietnamese text images (20K images from our employees and 180K generated).
- Testing data:
  - 1300 handwritten words image dataset. One image contains a single cropped word (sized 32x128)
- The HM drastically improved recognition accuracy while not pushing to much computing burden.

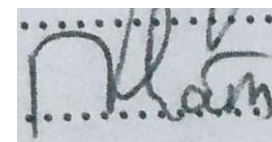| Model | Accuracy (percent) | Inference time | |
|-------|--------------------|----------------|---|
| | | (ms/word) | (ms/document) |
| VM (Vision Model) | 87.8 | 2.846 | 1343 |
| HM (Hyper Model) | 95.5 | 3.076 | 1452 |

# Sample demonstration

VM prediction: "cúy"
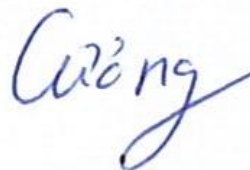HM prediction: "qúy"
Ground truth: "quý"

VM prediction: "phạn"
HM prediction: "phạm"
Ground truth: "phạm"

VM prediction: "thâđ"
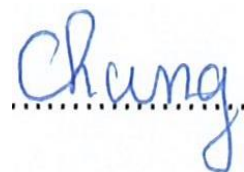HM prediction: "chai"
Ground truth: "nhâm"

VM prediction: "namc"
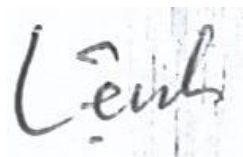HM prediction: "nam"
Ground truth: "nam"

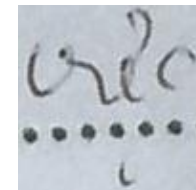VM prediction: "cuông"
HM prediction: "cương"
Ground truth: "cương"

VM prediction: "liitt"
HM prediction: "liiệt"
Ground truth: "việt"

VM prediction: "chunn"
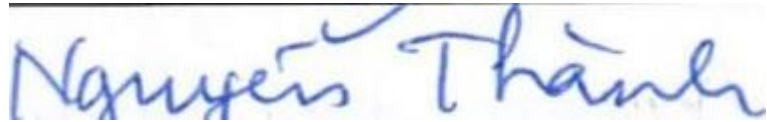HM prediction: "chung"
Ground truth: "chung"

VM prediction: "linh"
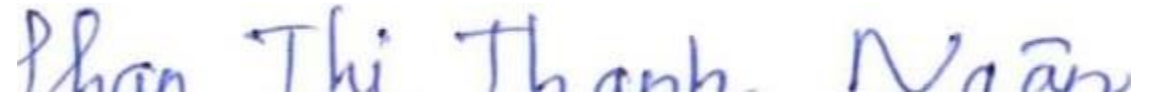HM prediction: "lệnh"
Ground truth: "lệnh"

VM prediction: "qiệc"
HM prediction: "đucc"
Ground truth: "vực"

# Full name benchmark

- Tested on 400 Vietnamese full names (2 to 4 words):
  - HM: 87% accuracy
  - VM: 77% accuracy



VM Prediction: "Nouyễn Thành"
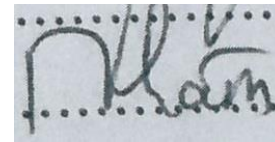HM Prediction: "Nguyễn Thành"
Ground truth: "Nguyễn Thành"



VM Prediction: "Phan Thi Thanh Noân"
HM Prediction: "Phan Thi Thanh Ngân"
Ground truth: "Phan Thị Thanh Ngân"

# Future works

- Improve the LM so it can performs well in cases where visual output has more than two errors.


VM prediction: "liitt"
HM prediction: "liiệt"
Ground truth: "việt"


VM prediction: "thâđ"
HM prediction: "chai"
Ground truth: "nhâm"

- Enhance the language model context knowledge even more by switching from character tokenizing mechanism to word tokenizing mechanism,.

# Thank you

# SAMSUNG SDS