

SCENE TEXT VIETNAMESE



"Nhận diện chữ tiếng Việt trong ảnh ngoại cảnh và sinh hoạt hằng ngày"

Nhóm trình bày: Nhóm 17

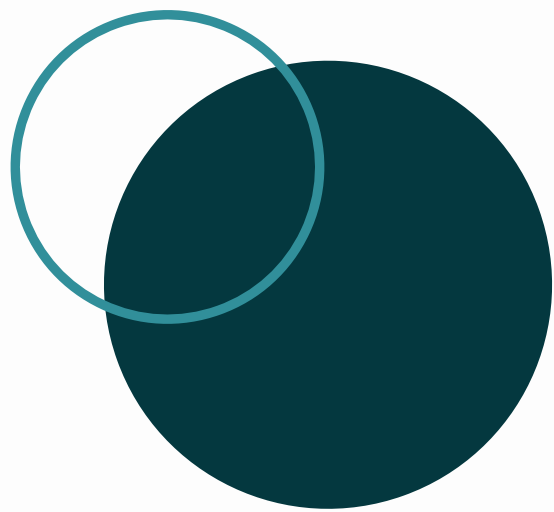
Sinh viên:

1. Nguyễn Lê Thanh - 19522238

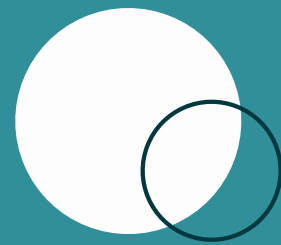
2. Phạm Vĩ - 19521101



CÁC CHỦ ĐỀ CHÍNH



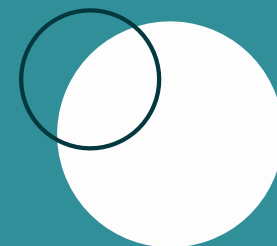
1. Tóm tắt lại bài toán Scene Text Vietnamses.
2. Các phương pháp đã thử nghiệm.
3. Các kết quả đạt được.
4. Tổng kết và hướng phát triển.



1. Tóm tắt bài toán *Scene text* Vietnamese.



Mục tiêu là phát hiện (*detect*) và nhận diện (*recognize*) chữ trong ảnh, cụ thể ở đây sẽ tập trung vào chữ trong khung cảnh (*scene text*) được thu lại từ nhiều nguồn camera khác nhau ở Việt Nam.



2. CÁC PHƯƠNG PHÁP ĐÃ THỬ NGHIỆM

THỬ NGHIỆM 2 HƯỚNG TIẾP CẬN TRÊN BỘ DỮ LIỆU VINTEXT.

Name	#imgs	#text instances	Examples
VinText	2000	About 56000	

- Hướng tiếp cận 1:
 - Detect: model SAST
 - Recognize: model SRN
- Hướng tiếp cận 2:
 - model Dict-guide

Vintext Dataset

- Folder Labels - chứa các file annotation của từng image.
- Folder train_images - chứa 1200 ảnh từ im0001 đến im1200.
- Folder test_image - chứa 300 ảnh từ im1201 đến 1500.
- Folder unseen_test_images - chứa 500 ảnh từ im1501 đến im2000
- File general_dict.txt
- File vn_dictionary.txt

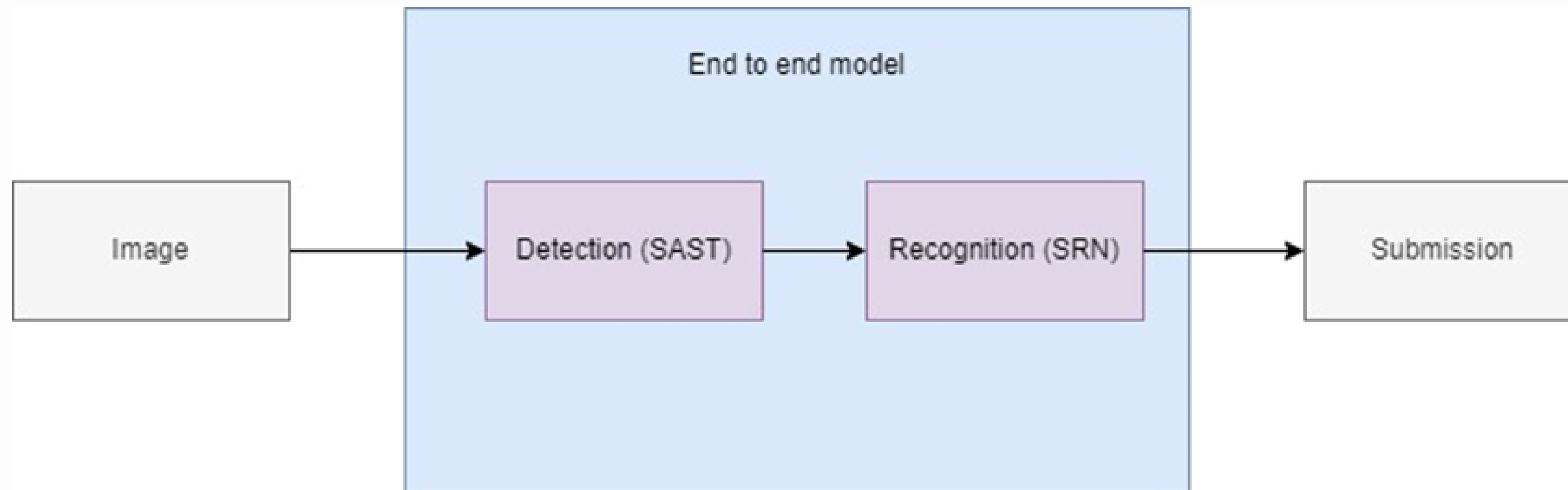


Ảnh trong folder train_image

```
690,402,712,402,716,420,692,420,###  
768,445,800,442,799,448,769,450,###  
782,506,799,506,799,514,784,515,###  
771,461,799,458,799,484,776,486,PHÒN  
356,113,363,114,363,130,351,129,Ở  
367,114,367,114,378,115,378,130,Y  
382,112,400,111,401,132,382,131,TẾ  
405,116,456,120,457,134,404,132,QUANG  
460,119,498,121,499,135,460,135,NINH
```

file .txt trong folder labels

HƯỚNG TIẾP CẬN 1



Quy trình xử lý hướng tiếp cận 1

Model Detect SAST

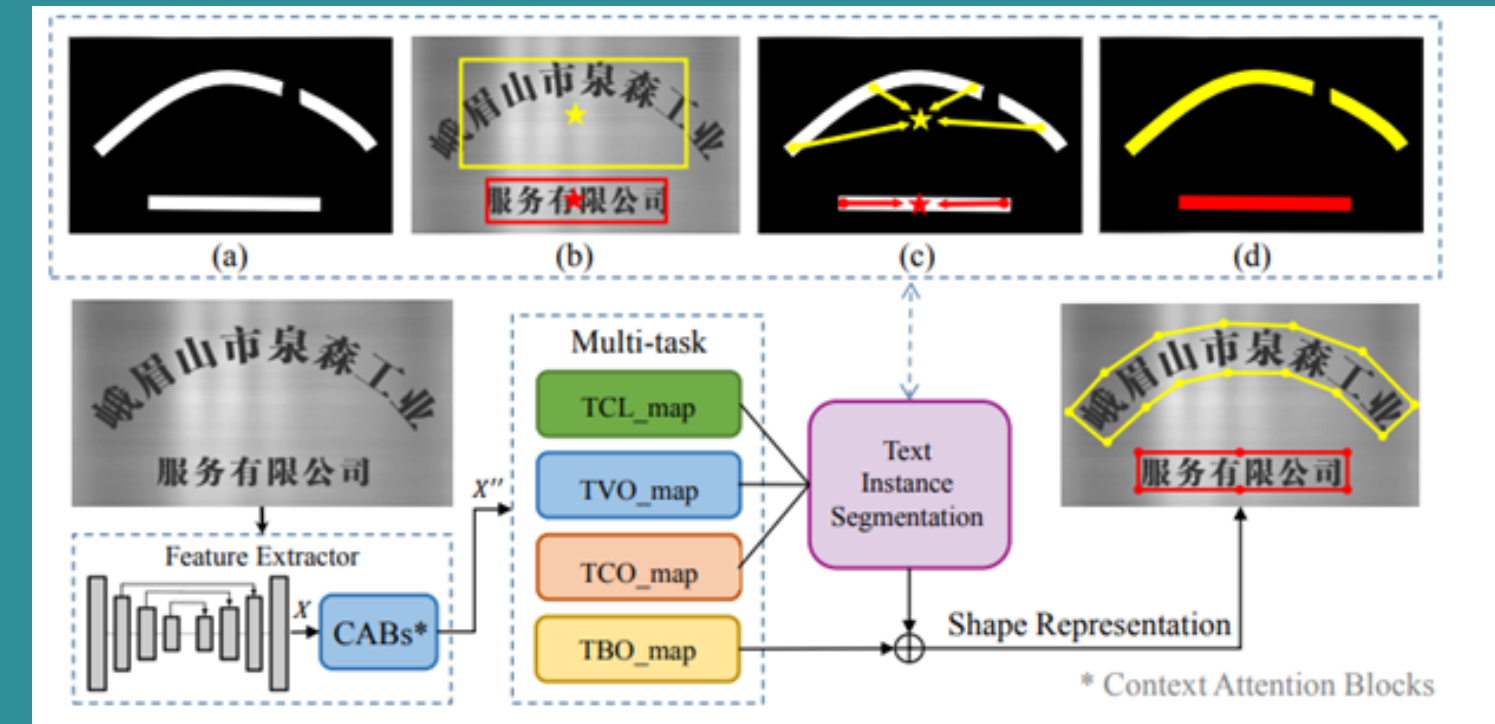
XỬ LÝ DỮ LIỆU

```
1 65,35,82,35,82,39,65,39,###
2 94,10,117,10,117,41,93,41,CHẤT
3 118,15,147,15,148,46,118,46,LƯỢNG
4 149,9,165,9,165,43,150,43,TỐT
5 167,9,180,9,179,43,167,42,ĐỂ
6 181,12,193,12,193,43,181,43,CỐ
7 195,13,215,14,215,46,196,46,VIỆC
8 217,13,237,14,239,47,217,46,LÀM,
9 240,14,265,14,266,45,240,44,NĂNG
10 267,11,291,11,291,46,267,45,SUẤT
11 293,15,310,16,311,46,292,45,CAO
12 312,11,327,10,324,44,312,44,ĐỂ
```

```
im1432.jpg [{"transcription": "cafe", "points": [[{"x": 277, "y": 32}, {"x": 367, "y": 22}, {"x": 364, "y": 60}, {"x": 278, "y": 70}]]}, {"transcription": "Nhỏ", "points": [[{"x": 252, "y": 39}, {"x": 294, "y": 60}, {"x": 275, "y": 155}, {"x": 239, "y": 127}]]}, {"transcription": "NGUYỄN", "point
im1297.jpg [{"transcription": "LÊ", "points": [[{"x": 487, "y": 107}, {"x": 512, "y": 105}, {"x": 518, "y": 120}, {"x": 488, "y": 122}]]}, {"transcription": "PHẦN", "poin
im1327.jpg [{"transcription": "ĐỊA", "points": [[{"x": 89, "y": 82}, {"x": 228, "y": 85}, {"x": 235, "y": 193}, {"x": 101, "y": 191}]]}, {"transcription": "XÁ", "points": [[{"x": 441, "y": 128}, {"x": 476, "y": 141}, {"x": 470, "y": 150}, {"x": 437, "y": 140}]]}, {"transcription": "###", "point
im1469.jpg [{"transcription": "BAN", "points": [[{"x": 96, "y": 88}, {"x": 114, "y": 88}, {"x": 114, "y": 97}, {"x": 96, "y": 97}]]}, {"transcription": "CHỈ", "points": [[{"x": 135, "y": 27}, {"x": 155, "y": 26}, {"x": 148, "y": 51}, {"x": 131, "y": 50}]]}, {"transcription": "X", "points": [[{"x": 265, "y": 116}, {"x": 358, "y": 117}, {"x": 358, "y": 142}, {"x": 263, "y": 144}]]}, {"transcription": "to", "po
im1408.jpg [{"transcription": "CHÙA", "points": [[{"x": 215, "y": 338}, {"x": 313, "y": 322}, {"x": 311, "y": 381}, {"x": 212, "y": 399}]]}, {"transcription": "LONG", "poi
im1377.jpg [{"transcription": "cafe", "points": [[{"x": 277, "y": 32}, {"x": 367, "y": 22}, {"x": 364, "y": 60}, {"x": 278, "y": 70}]]}, {"transcription": "Nhỏ", "points": [[{"x": 252, "y": 39}, {"x": 294, "y": 60}, {"x": 275, "y": 155}, {"x": 239, "y": 127}]]}, {"transcription": "NGUYỄN", "point
```

Ý tưởng:

- Dùng một vòng lặp for đọc từng file annotation của mỗi ảnh.
- Chạy một vòng lặp for đọc các dòng trong file annotation để lấy các text box.
- Dùng `split(',', 8)` để tách từng point và text ra sau đó lưu vào dictionary.



CHUẨN BỊ FILE CONFIG

- backbone ResNet50_vd

Precision	Recall	Hmean
0.8978295394388566	0.715008431703204	0.7960572635531565



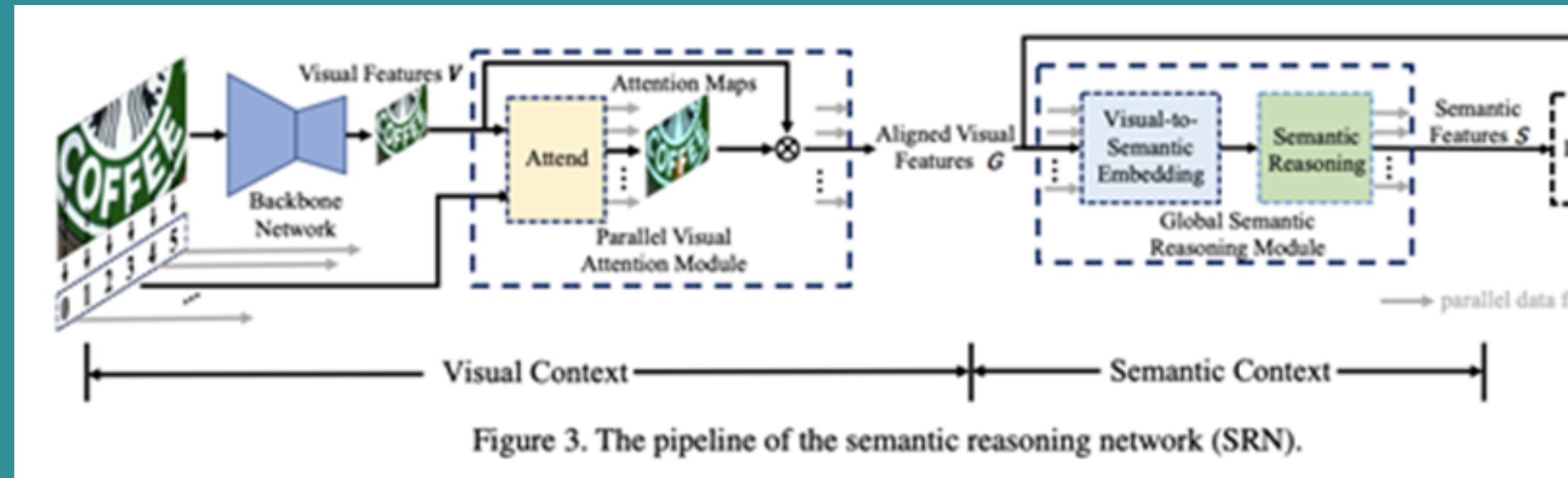
Ảnh sau khi được predict

```
Global:
    debug: false
    use_gpu: true
    epoch_num: 300
    log_smooth_window: 20
    print_batch_step: 2
    save_model_dir: ./output/SAST/
    save_epoch_step: 1
    eval_batch_step:
        - 40000
        - 50000
    cal_metric_during_train: false
    pretrained_model: ./pretrain_models/SAST/latest
    checkpoints:
    save_inference_dir: ./inference/SAST
    use_visualdl: True
    infer_img: null
    save_res_path: ./output/sast.txt
...
Train:
    dataset:
        name: SimpleDataSet
        data_dir: ./data #path_img_dir
        label_file_list: [ train_label.txt ]
        ratio_list: [1]
...
Eval:
    dataset:
        name: SimpleDataSet
        data_dir: #path_img_dir
        label_file_list:
            - unseen_label.txt
```


Model Recognize SRN

XỬ LÝ DỮ LIỆU

Data train model recognition sẽ là các ảnh nhỏ và một file txt lưu đường dẫn ảnh cùng với chữ trong ảnh đó.

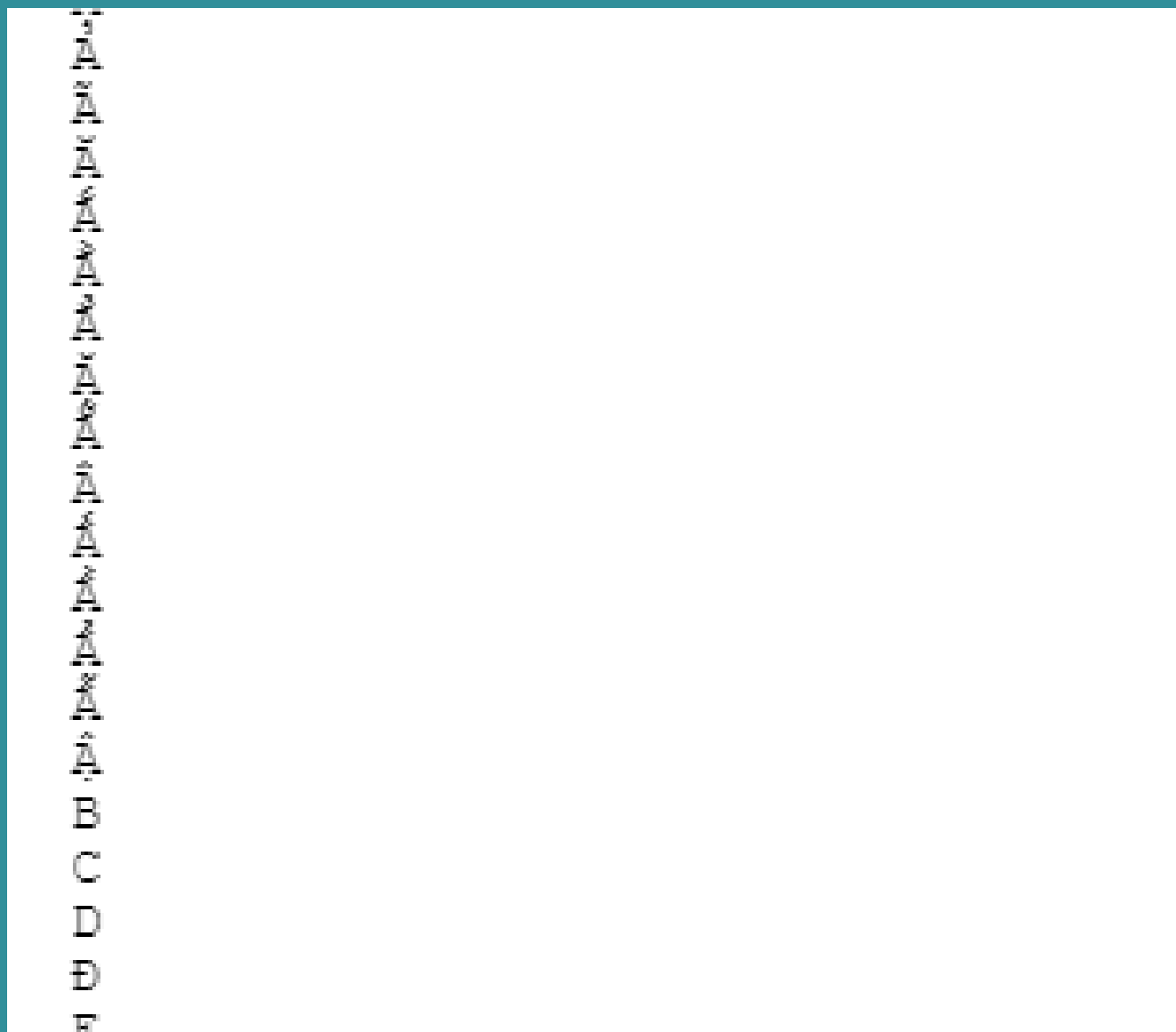


Ý tưởng:

- Dựa vào các point trong labels để cắt ảnh thành những ảnh nhỏ.

CHUẨN BỊ DICTIONARY

Để dự đoán được tiếng việt cần một file dictionary dành cho tiếng Việt chứa tất cả các kí tự.



CHUẨN BỊ FILE CONFIG

backbone là ResNetFPN.

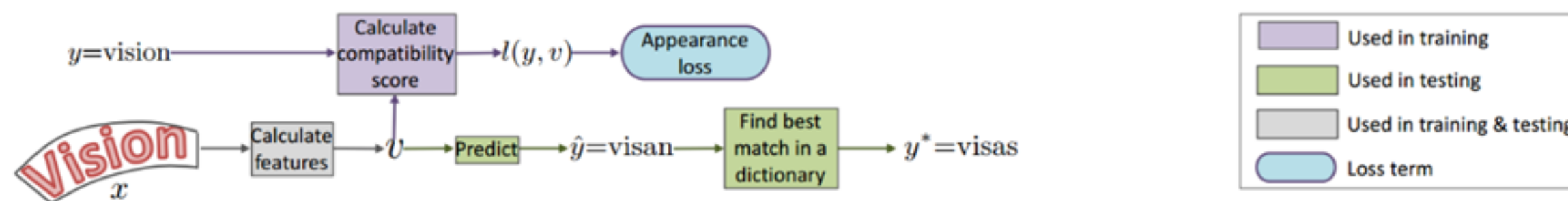
```
Global:
debug: false
use_gpu: true
epoch_num: 200
log_smooth_window: 20
print_batch_step: 5
save_model_dir: ./output/SRN
save_epoch_step: 1
eval_batch_step:
- 30000
- 40000
cal_metric_during_train: false
pretrained_model: ./pretrain_models/SRN/latest
checkpoints:
save_inference_dir: ./inference/SRN
use_visualdl: false
infer_img: doc/imgs_words/ch/word_1.jpg
character_dict_path: /content/drive/MyDrive/PaddleOCR/PaddleOCR/ppocr/utils/dict
character_type: ch
max_text_length: 25
num_heads: 8
infer_mode: false
use_space_char: true
...
```

KẾT HỢP 2 MODEL

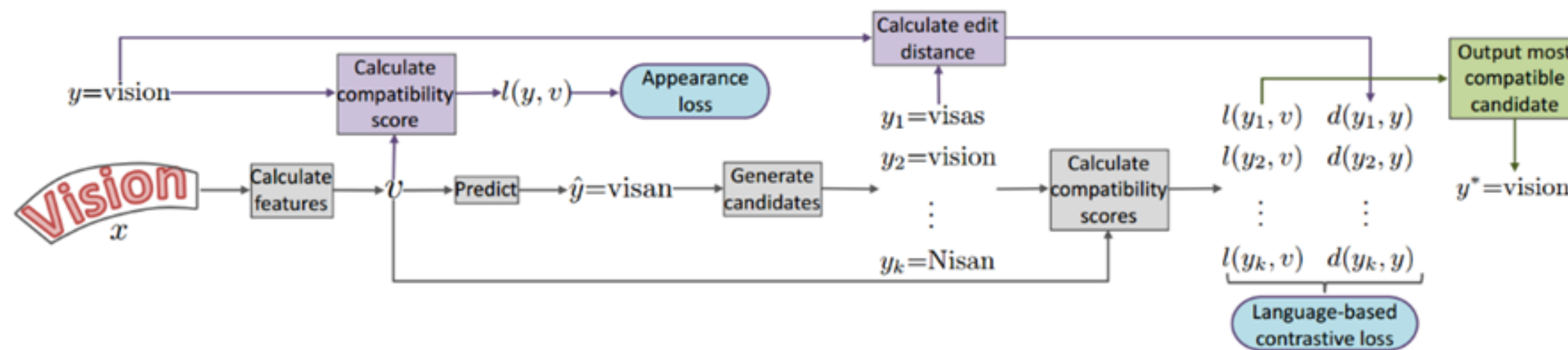


Hình ảnh predict sau khi kết hợp 2 model

HƯỚNG TIẾP CẬP 2



(a) The normal scene text recognition pipeline



(b) The proposed scene text recognition pipeline

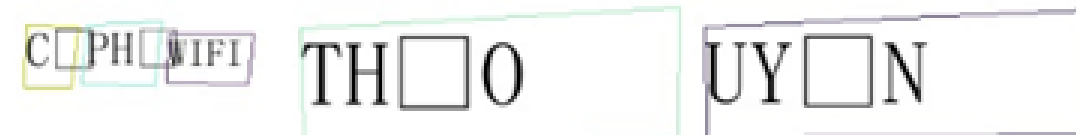
Quy trình hướng tiếp cận 2

3. KẾT QUẢ ĐẠT ĐƯỢC

	SAST+ SRN	Dict-Guide (VinAI)
Precision	0.51	0.9004
Recall	0.30	0.8014
Hmean	0.38	0.8480

4. Tổng kết và hướng phát triển

- Đã thử nghiệm được 2 hướng tiếp cận cho bài toán
- Hướng tiếp cận 1 (SAST + SRN) có kết quả đánh giá khá thấp, model SRN cho kết quả predict sai khá nhiều.
- Cải thiện bằng cách tăng cường dữ liệu với phương pháp cắt dán các bounding box có chứa chữ vào ảnh nền (không chứa văn bản) và thay đổi màu sắc, tính chất của bức ảnh.



DEMO

