

ПЮЮІ

2024 CUAU 중앙대학교 인공지능 학회 하계 컨퍼런스
Proceeding of 2024 Chung-Ang University Artificial Intelligence Summer Conference

Introduction

카메라로 촬영된 영상의 화질에 대한 정량 평가 점수를 예측하고, 그 평가 결과를 자연어로 상세하게 표현하는 알고리즘을 개발하는 것을 목표로 한다.

Methods

train.csv	test.csv	sample_submission.csv	train.csv	test.csv	sample_submission.csv
Views	Grid view	Hide fields	Filter	Group	Sort
img_name	img_path	mos	comments	img_name	img_path
1	41wy7upxzj	/train/41wy7upxzj.jpg	5.56923077	the pink and blue really c...	/test/j00zs3u6dr.jpg
2	yuyjiq6xxt	/train/yuyjiq6xxt.jpg	6.10317460	love rhubarb! great colors!	/test/ytv70so3zb.jpg
3	wk321130q0	/train/wk321130q0.jpg	5.54198473	i enjoy the textures and g...	/test/la9890oozp.jpg
4	w50dp2zjzp	/train/w50dp2zjzp.jpg	6.23484848	i like all the different colo...	/test/xs81t9px4a.jpg
5	l7rfqxeuh0	/train/l7rfqxeuh0.jpg	5.19047619	i love these critters, just ...	/test/f23994ghih.jpg
6	iapcid06sr	/train/iapcid06sr.jpg	5.93846154	excellent use of light, gre...	/test/oltx3kqzj.jpg
7	twvec6p4i1	/train/twvec6p4i1.jpg	7.38931298	double trouble! what grea...	/test/pzhjsj4uxo.jpg
8	h0wh5in52r	/train/h0wh5in52r.jpg	6.12307692	this is really meters from ...	/test/2mb38w9nb3.jpg
9	j70kuwif5y	/train/j70kuwif5y.jpg	6.86651473	i generally don't like overp...	/test/cz4tby7pfx.jpg
10	9qls7ros2m	/train/9qls7ros2m.jpg	5.75396825	awesome tones, great mo...	/test/8099bitbox.jpg
10 records			Sum: 60.71067621	10 records	

피처를 CNN에서 추출하는 기존의 캡셔닝 모델과 다르게, CPTR은 원본 이미지를 sequentialize 해 입력 데이터로 바로 사용한다. 인코더는 self-attention 매커니즘을 통해 패치들 간의 장기 의존성을 처음부터 활용할 수 있게 하고, 이미지를 고정된 크기 (16×16)로 작은 패치로 나눈 후, 각 패치를 flatten하고 이를 1D 패치 시퀀스로 변형하여 활용한다.

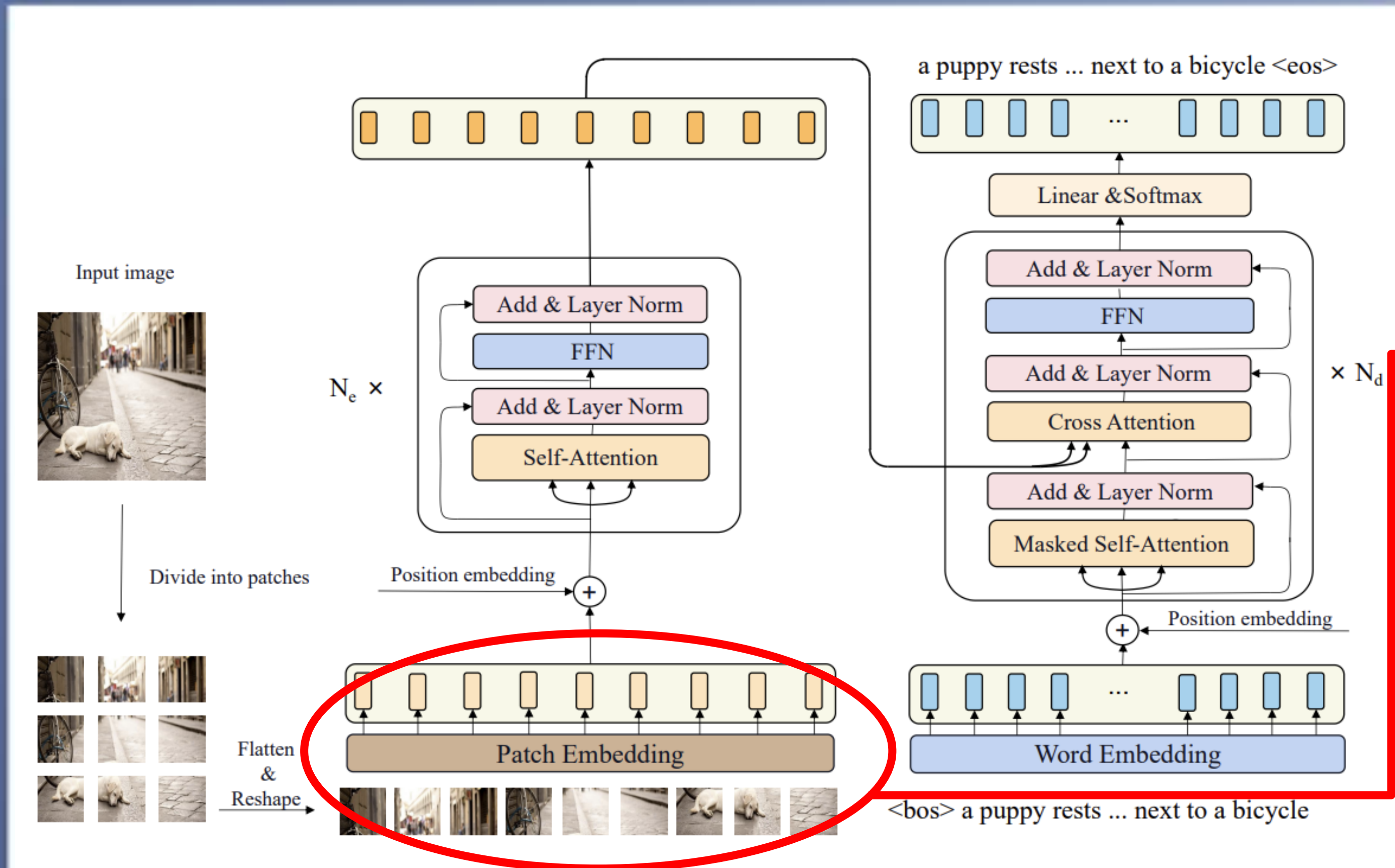


Fig. 1. The overall architecture of proposed CPTR model.

a. Processing

```
with open(file=path2captions, mode="r") as fp:
    img2captions = json.load(fp)

captions = list(img2captions.values())
captions = list(it.chain(*captions))

tokenizer = build_tokenizer(tok_name='spacy', lang='en_core_web_sm')
vocabulary = make_vocab(captions, tokenizer, SPECIALS2IDX)

logger.success('vocabulary has built')

serialize(path2vocabulario, vocabulary)
```

▽토큰화 처리

```

bos = torch.tensor([SPECIALS2ID['<bos>']])
eos = torch.tensor([SPECIALS2ID['<eos>']])

zip_imgtokenids = []
logger.debug('caption tokenization')
for key, val in track(itertools.items(*), 'build map_img2tokensids'):
    cap = cap_in[val]
    tok = tokenizer(cap.strip()).lower()
    idx = torch.tensor(vocabulary.get(tok))
    seq = torch.cat([bos, idx, eos]).numpy() # more effective for storage
    zip_imgtokenids.append((key, seq))

serialize(path2tokensids, zip_imgtokenids)

```

[illegible]

▽어휘집, 데이터 로드

```

logger.debug('load vocabulary')
vocabulary = deserialiize(path2vocabulary)
nb_tokens = len(vocabulary)

logger.debug('build dataset')
dataset = DatasetForTraining(path2trainfiles, path2features)
logger.debug('size of the dataset : %len(dataset) words')
dataset = OutDataset(dataset, batch_size=512, shuffle=False, collate_fn=custom_fn)
nb_data = len(dataset)
nb_data = len(dataset)

logger.debug('define network')
if path.isfile(path2checkpoint):
    net = th.load(path2checkpoint)
else:
    net = CaptionTransformer(
        in_dim=2048,
        hd_dim=256,
        ff_dim=512,
        nb_heads=8,
        num_encoders=5,
        num_decoders=1,
        pre_norm=False,
        seq_length=128,
        nb_tokens=nb_tokens,
        padding_idx=SPECIAL_2IDX['<pad>']
    )

net.to(device)
net.train()

```

< 트랜스포머 모델 생성

▽최적화와 학습, 손실 함수 설정

```
optimizer = th.optim.Adam(net.parameters(), lr=1e-5, betas=(0.9, 0.99), eps=1e-9)
criterion = nn.CrossEntropyLoss(ignore_index=SPECIALS2IDX['<pad>'])
logger.debug('training will begin ...!')
sleep(1)
```

```
memory = net.encode(src=src.to(device))
output = net.decode(
    tgt=tgt_input.to(device),
    memory=memory,
    tgt_mask=tgt_mask,
    tgt_key_padding_mask=tgt_key_padding_mask
)

logits = [net.generator(out) for out in output]
```

< encode/decode

[illegible]

```
vocab = deserialize(path2vocabulary)
logger.debug(f'vocab was loaded | len => {len(vocab)}')
```

< 어휘집 로드

```

logger.debug(f'load features extractor')

vectorizer = load_vectorizer(path2vectorizer)
vectorizer.eval()
vectorizer.set_device(device)

logger.debug('load ranker clip vit model')
ranker, processor = load_ranker(path2ranker, device)

logger.debug('features extraction by resnet152')

img_list = pull_images(os.path.dirname(path2image))
img_list = [os.path.basename(img) for img in img_list]
img_dir = os.path.dirname(path2image)

# 1. image = read_image(img_dir + '/' + img_name)
# 2. image = cv2.cvtColor(image, cv2.COLOR_BGR2RGB)
# 3. image = prepare_image(img, image)

# loading of dataset features(vectorizer, img_list,img,...) to device(1) : squeeze(0)
output_batch = fit_to_device(output_batch, 1) # 4, 0, 0, 0

response = beam_search(
    multi_infer,
    squeeze_output_batch(batch, ...),
    EOS_TOKEN_ID, num_beams=1,
    EOS_TOKEN_ID, num_beams=1,
    img_list=batch,
    beam_width=batch_width,
    beam_width=device,
    alpha=0.7)

# logger.debug(f'no generated : {len(response)}/')
scores = []
on_scores = []
caption = ''
caption = beam_search(response[1:-1]) # ignore those and ones
batch_caption = []
sentences.append(padded_caption)

< beam search

# 1. img = read_image(img_dir + '/' + img_name)
img_image = cv2.cvtColor(img, cv2.COLOR_BGR2RGB)
ranked_scores = rank_solutions(pil_image, sentences, ranker, processor, device)
ranked_response = list(zip(sentences, ranked_scores))
ranked_response = sorted(ranked_response, key=lambda t: t[1], reverse=True)

for caption, score in ranked_response:
    score = int(score * 100)
    # logger.debug(f'caption : {caption} | score : {score}/100')

best_caption, best_score = ranked_response[0]
logger.debug(f'filename : {img_name} | best caption : {best_caption} | score : {best_score}')

# csv 파일 저장
def save_file_path_writing(filename, results):
    # csv 파일 path = os.path.join('results/results.csv')
    with open(csv_file_path, mode='w', newline='') as file:
        writer = csv.writer(file)
        writer.writerow(['filename', 'caption', 'score'])
        writer.writerow([filename, best_caption,

```

Inspired by the above works, we consider solving the image captioning task from a new sequence-to-sequence perspective and propose CaPtion Transformer (CPTR), a full Transformer network to replace the CNN in the encoder part with Transformer encoder which is totally convolution-free. Compared to the conventional captioning models taking as input the feature extracted by CNN or object detector, we directly sequentialize raw images as input. Specifically, we divide an image into small patches of fixed size (e.g. 16×16), flatten each patch and reshape them into a 1D patch sequence. The patch sequence passes through a patch embedding layer and a learnable positional embedding layer before being fed into the Transformer encoder.

논문에서는 fully convolution free라는 점을 강조하나, 실제 코드에서는 CNN 기반 모델 (ResNet152) 을 사용해 이 이미지 피처를 추출하고 있다.

Results

#	팀	팀 멤버	화질 점수	캡처링	제출수	등록일
26	na_sanghyun	na	0.5594	0.66175	1	4일 전

데이콘에 베이스라인 모델의 결과와, 본문에서 소개한 모델의 결과를 제출하여 비교하려 하였으나 베이스라인 모델의 캡션 생성 오류로 진행하지 못하였다. 다만 본문 모델 캡션의 결과는 0.6175점으로, 캡셔닝 점수 기준 23등의 성과를 거두었다.

베이스라인 모델로 추론한 mos와 캡션을 데이콘에 제출하지 못하여 베이스라인 모델에 비한 본문 모델의 향상 정도를 수치적으로 나타내기 어려웠다. CNN으로 이미지 피처를 추출한 후 이를 트랜스포머 인코더에 입력으로 주는 방식을 적용하여 성능이 기대만큼 향상되지 않았으리라 생각한다. 향후 CLIP 등의 방대한 양의 데이터로 pre-trained ViT 인코더를 사용하여 성능 향상을 기대할 수 있을 것이다.

Conclusion

데이콘에서 주어진 이미지 데이터를 기반으로 MOS를 예측하고 캡션을 생성하였다. 이를 통해 트랜스포머 구조 기반 모델을 사용한 이미지 캡셔닝을 수행할 수 있었다. 자연어 처리에서 효과적인 성능을 보이는 트랜스포머 모델에 완전한 convolution-free 구조를 적용하기 위해 추후 pre-trained 인코더 등의 모델을 활용하는 연구를 지속하면 이미지 캡셔닝을 효과적으로 수행하는 모델을 개발할 수 있을 것으로 확인된다.

Reference

- 1) Wei Liu, et al. "CPTR: Full Transformer Network for Image Captioning", arXiv:2101.10804, 2021.
- 2) Tsung-Yi Lin, et al. "Microsoft COCO: Common Objects in Context", arXiv:1405.0312, 2015
- 3) Ashish Vaswani, et al. "Attention Is All You Need.", Google AI Language, 2017