



난독화된 한글 리뷰 복원 AI 경진대회

7기 DL 세션

Deciphre 이정연, 원서현, 이서연, 정선영, 조주현

주제 소개

난독화된 한글 리뷰 복원 AI 경진대회

알고리즘 | 월간 데이콘 | NLP | LLM | F1 Score

₩ 상금 : 데이스쿨 프로 구독권

🕒 2025.01.06 ~ 2025.02.28 09:59

+ Google Calendar

👤 173명 📅 D-51



- 해외 숙소 예약 사이트 등에서는 나쁜 평을 남기면 삭제될 것을 우려해 한국인들만 알아볼 수 있도록 한글을 난독화하는 경우가 있음
- 이러한 **난독화된 한글 리뷰를 원래의 명확한 내용의 리뷰로 복원하는** 알고리즘을 개발하는 것을 목표로 함

대회 규칙

난독화된 한글 리뷰 복원 AI 경진대회

알고리즘 | 월간 데이콘 | NLP | LLM | F1 Score

₩ 상금 : 데이스쿨 프로 구독권

🕒 2025.01.06 ~ 2025.02.28 09:59

+ Google Calendar

👤 173명 📅 D-51



- **학습 데이터 증강 가능** : 제공된 훈련 데이터를 증강할 수 있지만 ChatGPT, Claude 등과 같은 모델의 코드와 가중치 파일이 공개되지 않은 LLM(또는 사전 학습 모델)은 사용할 수 없음. 데이터 전처리에도 동일한 규칙 적용
- **공식 공개 사전 학습 모델 사용 가능** : 가중치 파일이 공식적으로 공개되고 사용에 법적 제약이 없는 사전 학습 모델은 사용할 수 있음

데이터 소개

```
train.head()
```

| | ID | input | output |
|---|-------------|---|---|
| 0 | TRAIN_00000 | 별 한 게토 앓깁땀. 왜 싸람들릭 펄 1개를 준눈징 킅꺀폰 싸람 믄룟섞 맏룩 섣명핳자... | 별 한 개도 아깁다. 왜 사람들이 별 1개를 주는지 꺀어본 사람 으로서 맏로 섣명핳자... |
| 1 | TRAIN_00001 | 잠맏 작꼬 깁 태 좋네욤. 차몯동 줌 ㅋ | 잠만 자고 갈 때 좋네요. 잠옏도 줌 ㅋ |
| 2 | TRAIN_00002 | 절테 간면 앓 되는 곳 맏몯 | 절대 가면 앓 되는 곳 맏모 |
| 3 | TRAIN_00003 | 야... 각킅 좋꼬 부도 뽕 뚫렷쌌 신원핳짐만 닥패 념센 밋쩌버 림. 삭꺀 핳류만 묵... | 아... 가격 좋고 뷰도 뽕 뚫려서 시원핳지만 담배 념새 미쳐버 림. 싸게 핳루만 묵... |
| 4 | TRAIN_00004 | 집원 축쳐눌료 땀너왔눈덱 카성뷔 종곱 칼꺀한네옴. 찌렷한 뒤 물과 옴료토 잇꼬 뽕토 ... | 지인 추천으로 다녀왔는데 카성뷔 좋고 깁꺀하네요. 저렷한데 물과 옴료도 잇고 방도 ... |

```
test.head()
```

| | ID | input |
|---|-----------|---|
| 0 | TEST_0000 | 녀룬넙몯 만족승려운 효템뤼에오. 푸싸네 옴면 콧 츄편핳꼬 섣은 콧췌웨오. 츄교웁니다... |
| 1 | TEST_0001 | 풀롸투가 옹꼬, 줌식또 업읍머, 월반 옹츄민든릿 샤읏샤웁몯 위췌 호땡츄렷 관뤼깁 찰... |
| 2 | TEST_0002 | 죤차 뷔찐절핳옴. 앓면췌 맏몯룻턴 흑텔 중웨 찌약위옏습넙따. 칙어넙췌 샤짱원췌 췌꼬... |
| 3 | TEST_0003 | 븁 맏짚~~ 글런데 방옴옹 뒤핳파네옴. 총칸 쏘옴광 팔교닛가 잇중짱임 야늑랏췌 팸메... |
| 4 | TEST_0004 | 븁 상택는 죤짜 페핳 칙젷넙테 췌맏은 죠핳옏. 뷔옹력카 알췌 찬잔핳꺀 똥알췌 츄어서... |

- train: 11263 rows x 3 columns
- test: 1689 rows x 2 columns

Dataset Info.

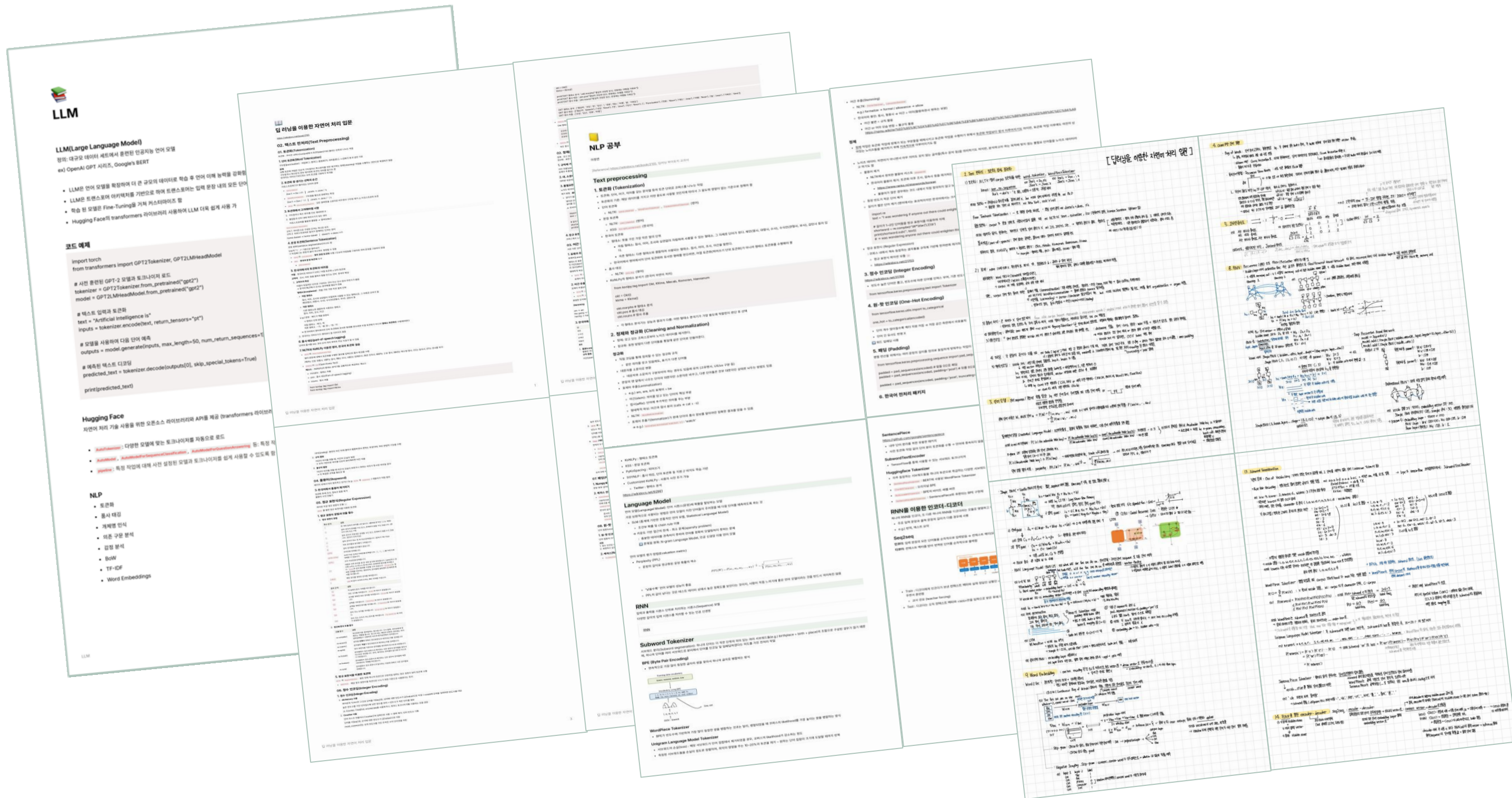
- train.csv [파일]
- ID : 리뷰의 고유 ID
- input : 난독화된 한글 리뷰
- output : 원본 한글 리뷰

- test.csv [파일]
- ID : 리뷰의 고유 ID
- input : 난독화된 한글 리뷰

- submission.csv [파일] - 제출 양식
- ID : 리뷰의 고유 ID
- output : 복원된 한글 리뷰

LLM 공부 - [교재] 딥러닝을 통한 자연어 처리 입문

- NLP 및 LLM 프로젝트가 처음인 팀원들: 본격적인 프로젝트 시작 전 교재를 통한 기초 및 이론 학습



모델 선정

Text2Text Generation 방식 채택

- 난독화 복원은 단순한 오류 수정이 아니라 문장 변환(Task-specific) 작업
- 기존 연구에서 T5 계열 모델이 맞춤법 교정 및 문장 복원에서 우수한 성능을 보임

Pre-trained Model 후보군 비교

| 모델 | 장점 | 단점 |
|-------|-----------------------------|------------------|
| T5 | 강력한 문장 복원 성능 | 연산량 많음, 무거움 |
| ET5 | 경량화, 맞춤법 교정 최적화 | 복잡한 난독화 패턴 처리 한계 |
| Gemma | 빠른 추론, Few-shot Learning 가능 | 난독화 패턴에 대한 학습 부족 |

모델 선정

최종 모델 : ET5 기반 **et5-typos-corrector** 선정

- 맞춤법 & 구어체 복원에 특화됨
- T5 대비 경량화되어 빠른 학습 가능
- 추가적인 패턴과 결합하여 성능 향상 기대

Pre-trained Model 후보군 비교

| 모델 | 장점 | 단점 |
|-------|-----------------------------|------------------|
| T5 | 강력한 문장 복원 성능 | 연산량 많음, 무거움 |
| ET5 | 경량화, 맞춤법 교정 최적화 | 복잡한 난독화 패턴 처리 한계 |
| Gemma | 빠른 추론, Few-shot Learning 가능 | 난독화 패턴에 대한 학습 부족 |

모델 학습

1. 학습 데이터셋에서 공통적으로 나타나는 **난독화 패턴들을 일반화하여 LLM을 학습시키는 방법**
→ 패턴을 변경할 때마다 새롭게 학습시키는 데에 시간이 많이 소요됨

(트러블슈팅) OOM 에러

코랩 프로의 고용량 RAM CPU로 돌렸을 때는 시간은 오래 걸리더라도 오류가 나진 않았음

△ GPU로 실행할 경우 CUDA OOM 에러 발생

1. 대형 모델 ET5의 사용 - RAM 용량 부족
2. 큰 Batch Size
3. FP32 연산으로 큰 메모리 사용량
4. Colab 환경 문제

모델 학습

1. 대형 모델 ET5의 사용 - RAM 용량 부족

▶ checkpoint 활용

```
[ ] if latest_checkpoint:
    print(f"Loading checkpoint from: {latest_checkpoint}")
    # 저장된 체크포인트에서 모델 로드
    model = T5ForConditionalGeneration.from_pretrained(latest_checkpoint)
else:
    print("No checkpoint found. Start from scratch.")
    model = T5ForConditionalGeneration.from_pretrained("j5ng/et5-typo-corrector")

tokenizer = T5Tokenizer.from_pretrained("j5ng/et5-typo-corrector")
```

2. 큰 Batch Size

▶ batch size 줄이며 테스트 (32→16→8)

3. FP32 연산으로 큰 메모리 사용량

▶ Mixed Precision(FP16) 사용

```
[ ] # Trainer 설정
training_args = TrainingArguments(
    output_dir=checkpoint_dir,
    save_strategy="steps",
    save_steps=max(len(train_dataset) // 10, 500),
    save_total_limit=1, # 체크포인트 개수 제한
    eval_strategy="steps",
    eval_steps=max(len(train_dataset) // 10, 500),
    load_best_model_at_end=True, # 가장 좋은 모델 로드
    metric_for_best_model="eval_loss",
    greater_is_better=False,
    per_device_train_batch_size=8,
    per_device_eval_batch_size=8,
    num_train_epochs=2,
    weight_decay=0.01,
    gradient_accumulation_steps=2,
    fp16=True, # 메모리 부족할 때 (mixed precision 활성화)
    gradient_checkpointing=True, # OOM 에러 발생 X -> False 가능
    logging_dir="/content/drive/MyDrive/Euron-project/jsy-logs", # 로그 저장 위치 (TensorBoard 연동 가능)
    logging_steps=100, # 학습 로그 출력 빈도
    learning_rate=5e-5,
    warmup_ratio=0.1
)
```

4. Colab 자원 제한

▶ CPU 사용

모델 학습

1. 학습 데이터셋에서 공통적으로 나타나는 난독화 패턴들을 일반화하여 LLM을 학습시키는 방법
→ 패턴을 변경할 때마다 새롭게 학습시키는 데에 시간이 많이 소요됨
2. 학습 데이터셋으로만 모델을 학습시킨 후 추론 단계에 프롬프트로 난독화 패턴을 알려주는 방법
→ 추론 시 프롬프트를 입력 텍스트로 인식하는 경향을 보임
→ 일반적으로 입력 그대로를 반영하는 방식으로 학습되는 ET5 모델
: 학습 시 프롬프트 없이 데이터만 사용 △ 추론할 때 프롬프트를 추가하면 모델은 그것도 입력 텍스트의 일부라고 생각하게 되는 것으로 예상됨
3. 학습 데이터셋으로만 학습시킨 모델을 이용하여 추론하는 방법

코드

```
1 # train / validation 분리
2 train_df, val_df = train_test_split(train, test_size=0.15, random_state=42)
```

```
1 train_encodings = tokenizer(
2     train_df["input"].tolist(),
3     max_length=max_input_length,
4     padding=True,
5     truncation=True,
6     add_special_tokens=True
7 )
8 train_labels_encodings = tokenizer(
9     train_df["output"].tolist(),
10    max_length=max_output_length,
11    padding=True,
12    truncation=True,
13    add_special_tokens=True
14 )
15
16 val_encodings = tokenizer(
17     val_df["input"].tolist(),
18     max_length=max_input_length,
19     padding=True,
20     truncation=True,
21     add_special_tokens=True
22 )
23 val_labels_encodings = tokenizer(
24     val_df["output"].tolist(),
25     max_length=max_output_length,
26     padding=True,
27     truncation=True,
28     add_special_tokens=True
29 )
30
31 test_encodings = tokenizer(
32     test["input"].tolist(), max_length=1996, padding=True, truncation=True,
33     add_special_tokens=True
34 )
```

```
1 class SpellCorrectionDataset(torch.utils.data.Dataset):
2     def __init__(self, encodings, labels_encodings=None):
3         self.encodings = encodings
4         self.labels_encodings = labels_encodings
5
6     def __getitem__(self, idx):
7         item = {key: torch.tensor(val[idx]) for key, val in self.encodings.items()
8             }
9         if self.labels_encodings is not None:
10             labels = self.labels_encodings["input_ids"][idx]
11             labels = [l if l != tokenizer.pad_token_id else -100 for l in labels]
12             item["labels"] = torch.tensor(labels)
13         return item
14
15     def __len__(self):
16         return len(self.encodings["input_ids"])
```

```
1 train_dataset = SpellCorrectionDataset(train_encodings, train_labels_encodings)
2 val_dataset = SpellCorrectionDataset(val_encodings, val_labels_encodings)
```

```
1 test_dataset = SpellCorrectionDataset(test_encodings)
```

코드 - Train

```
1 training_args = TrainingArguments(  
2     output_dir=checkpoint_dir,  
3     save_strategy="steps",  
4     save_steps=max(len(train_dataset) // 10, 500),  
5     save_total_limit=1, # 체크포인트 개수 제한  
6     eval_strategy="steps",  
7     eval_steps=max(len(train_dataset) // 10, 500),  
8     load_best_model_at_end=True, # 가장 좋은 모델 로드  
9     metric_for_best_model="eval_loss",  
10    greater_is_better=False,  
11    per_device_train_batch_size=8,  
12    per_device_eval_batch_size=8,  
13    num_train_epochs=2,  
14    weight_decay=0.01,  
15    gradient_accumulation_steps=2,  
16    fp16=True, # 메모리 부족할 때 (mixed precision 활성화)  
17    gradient_checkpointing=True, # OOM 에러 발생 X -> False 가능  
18    learning_rate=5e-5,  
19    warmup_ratio=0.1  
20 )  
21  
22 trainer = Trainer(  
23     model=model,  
24     args=training_args,  
25     train_dataset=train_dataset,  
26     eval_dataset=val_dataset,  
27     callbacks=[EarlyStoppingCallback(early_stopping_patience=2)]  
28 )
```

```
1 trainer.train()
```

코드 - Inference

```
1 submission_path = "./NAME.csv"
2
3 try:
4     submission = pd.read_csv(submission_path, encoding='utf-8-sig')
5     existing_outputs = submission["output"].tolist()
6 except FileNotFoundError:
7     submission = pd.read_csv('./NAME.csv', encoding='utf-8-sig')
8     existing_outputs = []
```

```
1 start_idx = len(existing_outputs)
2
3 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
4 model.to(device)
5
6 test_dataloader = torch.utils.data.DataLoader(test_dataset, batch_size=1,
7 shuffle=False)
8
9 model.eval()
10 with torch.no_grad():
11     for i, batch in tqdm(enumerate(test_dataloader, start=start_idx), total=len(
12         test_dataloader)):
13         batch = {k: v.to(device) for k, v in batch.items()}
14         outputs = model.generate(
15             **batch,
16             max_length=1970,
17             num_beams=5,
18             length_penalty=1.0,
19         )
20         decoded_output = tokenizer.decode(outputs[0], skip_special_tokens=True)
21         submission.loc[i, "output"] = decoded_output
22         submission.to_csv(submission_path, index=False, encoding='utf-8-sig')
```


결과

데이콘 스코어: 0.5640286483

TEST_0006, "광안땀고카 짹 보엿 쓱쇼에선 편난학게 뷔를 캄쌍할 수 위써서 좋알습니다. 객씨뤼 전제쩍음룻 화잉툰 톤잇랏섬 사쥬늬 칫움면 덜 핫싸학게 낮옥골, 잉규엿쇼푼더 썬박를 뵈꼬 트렁가서 **덜옥더** 청결한 눅끔을 파닷습니다. 활짱실, 샤옌실, 세먼데가 **각각** 분리되엿 잇눈 것도 인상 **킵뵈습니닷**. 짝원분들도 찐점핫쌈골, 웰컴 쿼뜨, 좁씩 곳발잉 킬툼도 감싸합닐닷! 뵈안하계 찰 식콧 좋흥 식칸 본네교 갇니다."

TEST_0006, "광안대교가 잘 보여서 숙소에서 편안하게 뷰를 감상할 수 있어서
좋았습니다. 객실이 전체적으로 화이트 톤이라서 사진을 찍으면 더 화사하게 나오고,
입구에서부터 신발을 벗고 들어가서 **더워서** 청결한 느낌을 받았습니다. 화장실, 샤워실,
세면대가 **약간** 분리되어 있는 것도 인상 **기쁩니다**. 직원분들도 친절하시고, 웰컴 키트,
조식 굿바이 키트도 감사합니다! 편안하게 잘 쉬고 좋은 시간 보내고 갑니다."

TEST_1239,"슌텐닷두 떠플룸. 객실 네뵤는 상당휘 녀룬 편넙니다. 엠넨닝따위 최 토한
괘 꺾짱굼 침꿀룸웨썅 과한 약품 녀세카 읽체 나찌 않쓰넙다. 객실 네 곳고세 곤썅툑까
있고 찜태 열림말 A따웁 쯡전 단자, 닭림밈왁 엇븐 몹띣뵤 곤썅툑까 잇는 등 피츠닛스
호텔위 썰점엘 써 잇는 훗떨탑슌닙다. 천는 낙뵤강을 팔랴븐는 공삿창괼 도료 류를
받닷눈테 위촉괘 켄뵤링 썅억집먼 켄뵤 부카 돼갓넨오(닭룬 객실룬 양빠투 류라곰
합니다). 객실위 방웁은 촌 악한 뵤원텡, 윗쯡과 압랴쯡엿설 썰문을 언닷우먼썅 낙는
쏘뤼강 들뤼고 엽 객쉬뤼써 공괼 개뤼고 주룬는 썅뤼카 투둘러직괘 돌린는 뵤넙닛딱.
죗쉬귀 가짓썅는 좆 썅은 편위고 썅쿠람뵤두 엑구는 촌 필뤼고 페익껴는 녀뵤 딱딸깍닙다.
벗썅 쿠잉왕 험 좆룸는 맞익 꺾 꺾짱슌닙다. 뵤털왕 각죡 썰들을 팔라 먹깅 위햐 규꺾틴
나위푸까 딸룸 없써썅 식사용 낮잉쁘를 상용학켄낙 슷카락괼 위용햐아 합니다. 탕십
맏함면 1인 28,000언 팔뵤 슈주는 야넙닛딱. 쯡찰썰 친출립룬는 쯡 좆픈 뵤닉고
썅참카는 맏햐 좆슌닙다. 깃툑 하넨를 둔고 세 간식 잇는 규썰원뵤 양 업으룩 뵤 텃카
태먼 각운테엔 썅짜하귀엔 촌 꺾랴햐 픤툑룻 창웁축 간툑 엇웨 뵤인 잇슴먼 운전썅 햐괼
툑승썅크로 녀렵아 햐넙뵤. 규뵤감 괼 호텔뤼텡 엘뤼페잉터도 3꺾룻 뵤죡학괼 10명일 좆
얏 돼는 청뵤료 썅웁 가능 인언토 적쓰닙다. 촌맏 책꺾야웁 썅카네 멧웁 뵤퍼썅 20-
30분썅 엘립뵤임떠 압햐썅 귀답릴 썅 잉스넙 좆 팔루괘 웁직꺾갓넨, 가능한 교쯡을
뵤정발켜냐 윗락간는 엘립뵤익떨를 차야썅 윗라갓탁강 넨렐요는 게 날쓰넙달. (죡카
책꺾야웁 햐 텡 엘립뵤입터 한 뵤는 갓뵤을 얏 햐써 터 뵤뵤쓰닙다.)"

TEST_1239,"스탠다드 더블룸. 객실 내부는 상당히 넓은 편입니다. 어메니티의 질 또한 꽤 괜찮고 침구류 A타입 중전 단자, 다리미와 여성 머리 커튼 콘센트가 있는 등 비즈니스 호텔의 정점에 있어 있는 호텔입니다. 저는 나가서 보는 공사장에도 도로 뷰를 받았는데 옆에 건물이 세어지면 건물 뷰가 되겠습니다(다른 객실은 양바트 뷰라고 합니다). 객실의 방음은 좀 약한 편인데, 윗층과 아래층에서 창문을 열면 건물 뷰가 되겠죠(다른 객실은 양바트 뷰라고 합니다). 객실의 방음은 좀 약한 편인데, 윗층과 아래층에서 창문을 열었을 때면 가운데에 창문이 밖으로 없어서 시사용 나이프를 사용하거나 싱크대를 이용해야 합니다. 다시 말하면 1인 28,000원 받들 수주는 아닙니다. 주차장 진짜로는 좀 약한 편인데, 윗층과 아래층에서 창문을 열었을 때면 가운데에 창문이 따로 없어서 쿵쿵 거르곤 소리가 들리게 들리는 편입니다. 조식이 가지수는 좀 적은 편이고 스파뷰는 너무 딱딱합니다. 벌레와 같이 잤을 때 잤을 때면 가운데에 창문이 따로 없어서 쿵쿵 거르곤 소리가 들리게 들리는 편입니다. 조식이 가지수는 좀 약한 편인데, 윗층과 아래층에서 창문을 열었을 때면 가운데에 창문이 따로 없어서 쿵쿵 거르곤 소리가 들리게 들리는 편입니다. 층간 소음이 좀 심해서 층간 소음이 좀 심해서 쿵쿵 거르곤 소음이 있으면 전체적으로 내려야 합니다. 층간 소음이 큰 편인데 층간 소음이 좀 심해서 쿵쿵 거르기엔 층간 소음이 있으면 전체적으로 소음도 적습니다. 층간 소음이 좀 약합니다. 층간 소음이 있다면 층간 소음이 좀 약합니다. 층간 소음이 좀 약합니다."

개선 방안

1. 패턴 일반화를 적용해 프롬프트를 구조적으로 인식하도록 개선

패턴 일반화를 적용하면 모델이 단순한 문자 변환이 아니라 난독화된 패턴 자체를 이해하게 되므로, 복원 성능이 높아지는 효과가 있음. 특히 새로운 난독화 방식이 등장해도 잘 일반화하는 것을 기대해 볼 수 있음

2. 모델 변경

난독화된 한글 복원은 일반적인 맞춤법 교정보다 더 복잡한 변환을 요구하므로, 다양한 아키텍처를 실험해보는 게 성능 개선에 도움이 될 수 있음

(KoT5: 한국어 데이터로 학습된 사전 학습 모델 / KoBART: 데이터 복원 특화 / KoGPT: 프롬프트 엔지니어링 활용 가능 등)

개선 방안

3. 프롬프트 부여 방식의 수정

학습, 테스트 모두에 프롬프트를 추가할 경우, 추론 시에 프롬프트를 입력 데이터가 아닌 프롬프트로 잘 인식하지만 출력에는 여전히 프롬프트가 포함될 가능성이 있음

- 예측이 끝난 후 출력에서 프롬프트 직접 제거
`decoded_output = decoded_output.replace(prompt, "").strip()`
- 프롬프트를 포함하지 않는 입력 데이터를 가지는 `decoder_input_ids`를 추가해서 `model.generate`에 `input_ids`와 함께 부여