

현대어 변환 프롬프트톤

[Track 2 : 업스테이지] SOLAR Pro2 API 활용 프롬프트 엔지니어링

BIBIBIG 김호성 박준희 이예지

현대어 변환 프롬프트

CONTENTS

01

주제 분석 및 EDA

02

전략 및 접근 방법

03

사용 방법

04

최종 성능

05

기타 시도들

주제 분석 및 EDA

데이터

- 변수 파악
 - id, section, publication_date
 - original/answer_title
 - original/answer_sentence
- 원본 기사
 - 한자, 기호, 영어 혼합
 - 띄어쓰기 X
- 변환 기사
 - 한자, 한자어 → 문맥 따라 변환 후 삭제
 - 기호 □ → 문맥 따라 복원
 - 그 외 기호들 (『 』) → 유지
 - 조사 변경(‘에’ → ‘의’)

평가지표

- 네 가지 평가 영역으로 구성
 - Omission
 - Restoration
 - Naturalness
 - Accuracy
- 네 가지 값들의 평균으로 최종 점수 산정

전략 및 접근 방법

• 데이터, 평가지표에 따라 프롬프트 구성

데이터

- 한문 처리 2단계 전략
 - 직역 → 자연스러운 한국어 기사체로 재작성하는 중간 변환 절차 정의
 - 의미 왜곡 없이 현대 문체 정착 유도
- 현대 한국어 맞춤법·띄어쓰기 통제
 - 조사·접미사·명사 결합 규칙 명확화
 - 2025년 표준 맞춤법 기준 반영
- 예시 기반 Few-shot 프롬프팅
 - 복잡한 규칙을 단순 지침이 아니라 예시 기반으로 쉽게 인식하도록 유도
 - '에/의', 한자 번역, □ 복원 등 핵심 규칙 사례로 명확히 전달

평가지표

- 정보 손실 방지 규칙화
 - 문장·구조·수치·조건 등 모든 의미 요소 유지를 명시
 - 변환 과정에서 임의 삭제나 축약 금지
- □(결손 텍스트) 복원 규칙 설정
 - 문맥 기반 최소 단위로 복원
 - 과도한 해석·추가 정보 생성 금지

사용 방법 : 프롬프트만 사용 vs 멀티턴 추가

0.5200

baseline_prompt = (
""

아래 문장을 현대 한국어 기사체 느낌으로 자연스럽게 바꿔줘.
필요한 설명 없이 변환된 문장만 한 줄로 출력해!

1. 필수 변환 원칙

① 내용 누락 방지

- 원문에 존재하는 정보는 절대로 빠뜨리지 않는다.
- 문장, 절, 핵심 서술, 조건, 수식, 수치 정보까지 모두 보존한다.
- 한문 문장은 전체 의미가 동일하게 유지되도록 완전 번역한다.

② □ 복원 원칙

- □는 문맥상 자연스러운 **최소 단위(단어 또는 짧은 구)**로 복원한다.
- 과도한 의미 확장 금지(단 하나의 □를 전체 문장으로 만들지 않는다).
- 문맥과 맞지 않거나 임의로 유추한 복원 금지.

③ 띄어쓰기 원칙

- 2025년 표준국어대사전 기준의 현대 한국어 띄어쓰기 원칙을 따른다.
- 조사 ‘이/가/은/는/을/를/에/에서/으로/에게/뿐’ 등은 항상 붙여 쓴다.
- 명사+명사 결합은 의미 단위에 따라 표준에 맞게 판단한다. (예: 국가 정책, 정보 처리, 데이터 분석)
- 한문어 표현을 직역했을 때 띄어쓰기 오류가 나는 경우, 의미 단위를 기준으로 자연스럽게 풀어서 바르게 띄어쓴다.
- 접속어·부사 어절은 관용적 결합을 고려해 자연스럽게 띄어쓴다. (예: 또한, 다만, 그러나, 결국 → 관용적 형태 유지)
- ‘-되다/-하다’류의 파생어는 현대 표기 기준으로 붙여 쓰기 (예: 개선되다, 적용하다, 분석하다)

④ 조사 원칙 (‘에’, ‘의’ 포함)

- ‘의’: 소유·귀속·설명 관계일 때 사용 (예: 군주의 명령, 인민의 생활)
- ‘에’: 장소·시간·상태·기준을 나타낼 때 사용 (예: 성 안에, 그때에, 규정에 따라)
- 한문 구조를 그대로 옮기지 않고, 문맥에 맞게 조사 재배치
- ‘으로, 에게, 에서, 마다, 뿐, 따위’ 등도 문맥에 맞춰 자연스럽게 적용
- 조사 선택으로 인해 의미 왜곡이 생기지 않도록 정확히 판단한다.

--- 종략 ---

프롬프트

0.7190

3턴 구조: 의미보존 -> 자연스러움 -> 평가기준 기반 최종 교정

1턴: 의미 보존

first_system_prompt = (
"당신은 한자, 고어, 한문, 영어 혼용문을 현대 기사체로 변환하는 전문가입니다.\n"
"- 핵심 기준: *Omission 방지*, *Accuracy 유지*\n"
"- 단 하나도 의미를 빼거나 추가하지 마세요.\n"
"출력은 변환된 문장만."
)

2턴: 자연스러움

second_system_prompt = (
"당신은 한국어 문장 다듬기 전문가입니다.\n"
"- *Naturalness* 향상\n"
"- 그러나 *Restoration & Omission & Accuracy*는 유지\n"
"출력은 개선된 문장만."
)

second_user_prompt = f"원문:\n{text}\n\n1차 변환:\n{first_result}\n\n위 문장을 더 자연스럽게 다듬되 의미 보존."

3턴: 평가 기준 기반 최종 검수

third_system_prompt = (
"최종 검수 단계입니다.\n"
"다음 기준으로 문장을 다시 점검하세요:\n"
"1. Omission(누락 X)\n"
"2. Restoration(□ 복원 시 정확성 유지)\n"
"3. Naturalness(읽기 자연스러움)\n"
"4. Accuracy(원문 의미 왜곡·추가 금지)\n\n"
"가능한 최소 수정만 적용하세요. 출력은 최종 문장만."
)

third_prompt = f"원문:\n{text}\n\n2차 문장:\n{second_result}\n\n위 문장이 기준을 충족하는지 확인하고 필요하면 수정."

프롬프트 + 멀티턴

현대어 변환 프롬프트

사용 방법 : 멀티턴 프롬프트 수정

0.7190

3턴 구조: 의미보존 -> 자연스러움 -> 평가기준 기반 최종 교정

```
### 1턴: 의미 보존
first_system_prompt = (
    "당신은 한자, 고어, 한문, 영어 혼용문을 현대 기사체로 변환하는 전문가입니다.\n"
    "- 핵심 기준: *Omission 방지*, *Accuracy 유지*\n"
    "- 단 하나도 의미를 빼거나 추가하지 마세요.\n"
    "출력은 변환된 문장만."
)

### 2턴: 자연스러움
second_system_prompt = (
    "당신은 한국어 문장 다듬기 전문가입니다.\n"
    "- *Naturalness* 향상\n"
    "- 그러나 *Restoration & Omission & Accuracy*는 유지\n"
    "출력은 개선된 문장만."
)

second_user_prompt = f"원문:\n{text}\n\n1차 변환:\n{first_result}\n\n위 문장을 더 자연스럽게 다듬되 의미 보존."

### 3턴: 평가 기준 기반 최종 검수
third_system_prompt = (
    "최종 검수 단계입니다.\n"
    "다음 기준으로 문장을 다시 점검하세요:\n"
    "1. Omission(누락 X)\n"
    "2. Restoration(□ 복원 시 정확성 유지)\n"
    "3. Naturalness(읽기 자연스러움)\n"
    "4. Accuracy(원문 의미 왜곡·추가 금지)\n"
    "가능한 최소 수정만 적용하세요. 출력은 최종 문장만."
)

third_prompt = f"원문:\n{text}\n\n2차 문장:\n{second_result}\n\n위 문장이 기준을 충족하는지 확인하고 필요하면 수정."
```



0.8200

3턴 API 호출 함수:

- 1턴 - 의미 보존,
- 2턴 - 자연스러움 + 의미 보존,
- 3턴 - 원문과 비교하여 내용 보존 확인 및 보강

```
# 첫 번째 턴: 의미 보존도가 중요
first_system_prompt = (
    "당신은 한자, 고어, 한문, 영어 혼용문을 현대 기사체로 변환하는 전문가입니다. "
    "가장 중요한 것은 원문의 모든 의미와 정보를 빠짐없이 보존하는 것입니다. "
    "원문에 있는 모든 내용을 누락 없이 변환하고, 의미를 왜곡하거나 축약하지 마세요. "
    "변환된 텍스트만 출력하세요.")

# 두 번째 턴: 자연스러움과 의미 보존도가 중요
second_system_prompt = (
    "당신은 현대 한국어 기사체 문장을 다듬는 전문가입니다. "
    "주어진 문장을 더 자연스럽게 읽기 쉽게 개선하되, "
    "원문의 의미와 정보는 절대 누락하거나 왜곡하지 마세요. "
    "자연스러운 현대 한국어 표현으로 다듬으면서도 의미 보존도를 최대한 유지하세요. "
    "개선된 텍스트만 출력하세요. "
    "단, 문장을 다듬은 이유나 판단 과정은 절대 출력하지 마세요.")

second_user_prompt = (
    f"원문:\n{text}\n\n"
    f"첫 번째 변환 결과:\n{first_result}\n\n"
    "위 변환 결과를 더 자연스럽게 읽기 쉽게 개선해주세요. "
    "단, 원문의 모든 의미와 정보는 반드시 보존해야 합니다.")

### 3턴: 평가 기준 기반 최종 검수
# 세 번째 턴: 원문과 두 번째 결과 비교하여 내용 보존 확인 및 보강
third_system_prompt = (
    "당신은 원문과 변환 결과를 비교하여 내용 보존을 검증하는 전문가입니다. "
    "원문의 모든 의미, 정보, 세부사항이 변환 결과에 제대로 포함되어 있는지 철저히 확인하세요. "
    "누락된 내용이 있다면 반드시 보강하고, 왜곡된 내용이 있다면 수정하세요. "
    "원문의 모든 정보를 빠짐없이 포함하면서도 자연스러운 현대 한국어로 표현하세요. "
    "최종 보강된 텍스트만 출력하세요.")

third_user_prompt = (
    f"원문:\n{text}\n\n"
    f"두 번째 변환 결과:\n{second_result}\n\n"
    "위 변환 결과를 원문과 비교하여 다음을 확인하고 보강해주세요:\n"
    "1. 원문의 모든 의미와 정보가 포함되어 있는지 확인\n"
    "2. 누락된 내용이 있다면 반드시 추가\n"
    "3. 왜곡되거나 잘못된 내용이 있다면 수정\n"
    "4. 자연스러운 현대 한국어 표현 유지\n"
    "원문의 모든 내용을 빠짐없이 보존하면서도 자연스러운 최종 결과를 출력하세요.")
```

사용 방법 : 프롬프트 원본 vs 간소화

0.8200

baseline_prompt = (
""

필요한 설명 없이 변환된 문장만 한 줄로 출력해.

0. 목적

- 한문·고전문 + 누락된 글자(□)를 포함한 문장을 두 단계 변환 절차로 처리한다.
- 원문의 모든 정보를 보존하면서 매우 자연스러운 2025년 한국어 기사체로 변환한다.
- multi-turn 검토를 통해 자연스러움과 정확성을 추가 개선한다.

1. 필수 변환 원칙

① 내용 누락 절대 금지

- 원문에 존재하는 정보는 단어 단위까지 모두 보존한다.
- 한문 어휘, 조건, 인과, 대조 구조, 수식(數値), 고유명사, 시간·장소 정보 등을 빠짐없이 유지한다.
- 문장 구조가 아무리 부자연스러워 보여도 삭제하거나 축약하지 않는다.
- 하나의 문장이 너무 길더라도 문장 분리 없이 원문 정보량을 보존한다.

② □ 복원 원칙(강화 + 특수문자 추가)

- □는 문맥상 자연스러운 **최소 단위(단어 또는 짧은 구)**로만 복원한다.
- 과도한 의미 확장 금지, ‘임의 추측 단어’ 금지.
- 『』 「」 《》 〈〉 【】 ▲★● 등 특수문자·기호는 원문 그대로 유지하고 변경하지 않는다.
- 기호 옆이나 안에 있는 □는 기호 구조를 유지한 채 복원한다. (예: 『□ 정책』 → 『신규 정책』)

③ 띄어쓰기 원칙

- 2025년 표준국어대사전 띄어쓰기 기준 적용 + 다음 세부 규칙을 반드시 따른다.

가. 조사

- 조사(이/가, 은/는, 을/를, 에, 에서, 으로, 에게, 뿐 등)는 항상 붙여쓴다.
- 한문 구조를 그대로 베끼지 않고 현대 한국어 기준으로 조사 위치를 조정한다.

나. 명사 결합

- 명사 + 명사 결합이 의미 단위에 따라 띄어쓰기 여부를 결정한다. (예: 국가 정책, 기업 진단 결과 검토, 정보 처리, 데이터 분석)
- 특히 원문에서 모두 붙어 있는 경우에도, 의미 단위별로 정확히 분할해 표준적으로 띄어쓴다.

다. 의존명사

- 것, 수, 점, 차, 때, 바, 데 등은 앞말과 반드시 띄어쓴다.

라. 한자어·전문어

- 관례적으로 띄어쓰는 표현은 그대로 따른다. (예: 기술 개선, 정책 시행, 업무 보고, 성 안, 산 위)
- 한문 투 표현을 해석할 때 띄어쓰기가 불명확하면 의미 단위 기준으로 재구성한다.

--- 종락 ---

목표1

0.8330

baseline_prompt = (
""

아래 문장을 현대 한국어 기사체로 변환해줘.
설명 없이 변환된 문장만 한 줄로 출력해.

[핵심 원칙 — 반드시 모두 준수]

1) 내용 누락 절대 금지

- 원문의 정보(문장·절·구·조건·대조·수치·고유명사)를 단어 단위까지 모두 보존한다.
- 어떤 정보도 축약·삭제·생략하지 않는다.

2) 문장 구조 보존

- 원문의 의미 단위와 논리 흐름을 그대로 유지한다.
- 문장이 길어도 분리하지 않고 정보량을 유지한다.

3) □ 복원

- □는 문맥상 필요한 최소 단어만 복원한다.
- 과도한 의미 확장·임의 추측 금지.

4) 표현 현대화

- 고어·한문투를 모두 현대 한국어 표현으로 바꾼다.
- 조사·어순·띄어쓰기를 바로잡되 의미는 절대 변형하지 않는다.
- 자연스럽게 다듬되 정보 삭제 없이 재배치만 한다.

5) 한문 처리(2단계 내부 변환)

- (1) 한문을 의미 단위별 직역으로 해석한 뒤
- (2) 그 내용을 그대로 유지하며 현대 기사체로 재구성한다.

[출력 규칙]

- 자연스럽게 원문의 모든 의미·정보가 완전히 포함된 한 줄 문장만 출력.

원문:

{text}
""

.strip()
)

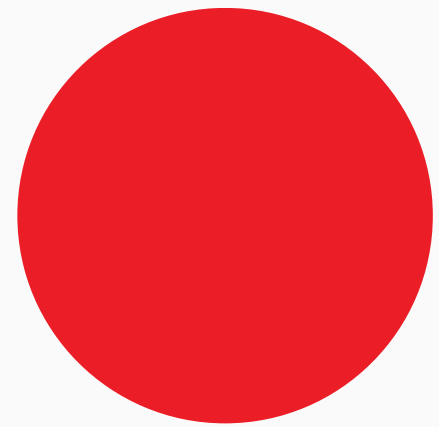
현대어 변환 프롬프톤

Temperature 실험

prompt	1st temp	2nd temp	3rd temp	omission	restoration	naturalness	accuracy	evaluate	leaderboard	final
full	0.0	0.0	0.0	0.932	0.995	0.969	0.961	0.964	0.8200	0.7990
full	0.0	0.1	0.1	0.936	0.994	0.989	0.983	0.976	-	-
full	0.1	0.1	0.1	0.930	0.992	0.958	0.975	0.964	-	-
full	0.0	0.15	0.1	0.951	0.991	0.985	0.978	0.976	-	-
full	0.0	0.2	0.1	0.947	0.993	0.980	0.971	0.973	-	-
full	0.0	0.15	0.05	0.945	0.993	0.976	0.959	0.968	0.8330	0.7960
<u>simple</u>	<u>0.0</u>	<u>0.0</u>	<u>0.0</u>	<u>0.932</u>	<u>0.995</u>	<u>0.969</u>	<u>0.961</u>	<u>0.964</u>	<u>0.8330</u>	<u>0.8140</u>
simple	0.0	0.15	0.1	0.944	0.987	0.974	0.956	0.965	0.8360	0.7900
simple	0.0	0.15	0.05	0.946	0.994	0.966	0.953	0.965	0.8240	0.8000

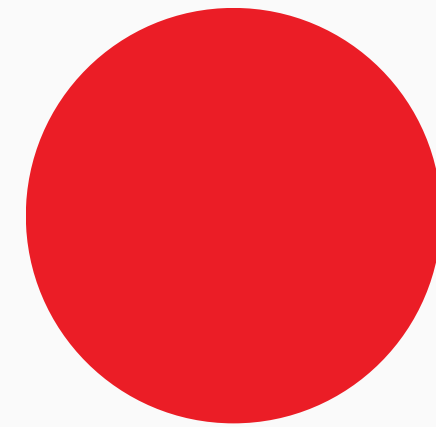
현대어 변환 프롬프트

최종 성능



Overall 0.8330
Final 0.8140

-
- 멀티턴 3회
 - 정확성 → 자연스러움 → 보강
 - 프롬프트 간소화



Overall 0.8360
Final 0.7900

-
- 멀티턴 3회
 - 정확성 → 자연스러움 → 보강
 - 프롬프트 간소화
 - 파라미터 수정

현대어 변환 프롬프트

기타 시도들

✓ 멀티턴 횟수 증가 (2 → 3 → 5)

성능 하락 확인, 멀티턴 3번으로 확정

✓ RAG

- BGE-M3 임베딩
 - 시도 1) metadata : answer_sentence
 - 시도 2) section에 따라 filtering 진행
→ 일반 모델 대비 성능 하락
-

2025.11.23

감사합니다

MixUP_BIBIBIG