

# 국어 문항 개발에 최적화된 모델

국어 문항 개발 및 분석(특히 수능/모의고사 수준)은 인공지능에게 가장 까다로운 영역 중 하나입니다. 글자 그대로의 번역을 넘어, 맥락(Context), 함축적 의미, 시대적 배경, 그리고 출제자의 의도(출제 오류 가능성 배제)까지 모두 계산해야 하기 때문입니다.

현존 최강의 프론티어 모델인 **OpenAI GPT-5.5**와 **Anthropic Claude 4.8 Opus**를 번갈아 가며 국어 문항을 개발하고, 기존 문항을 분석하고, 해설하는 프로젝트를 진행하고자 할 때, 특히 난해한 현대시나 한자와 고어가 들어 가는 고전 문항을 다룰 때, 어떤 모델을 이용한 것이 효율적이며, 두 모델을 번갈아 사용했을 때의 장단점을 문항 개발자 관점에서 철저히 분석해 보겠습니다.

## 🏛️ 문학 문항 개발 시 두 모델의 특성 및 유효성

### 1. Claude 4.8 Opus: "문학적 맥락과 텍스트 해석의 마스터" (현대시/고어 분석 우위)

- **현대시 (정서와 은유 해석):** 시어의 다의적 의미와 시적 화자의 정서적 태도를 읽어내는 데 독보적입니다. 단순히 사전적 정의를 나열하는 게 아니라, 작품 전체를 관통하는 '주제적 긴장감'이나 '반어·역설'의 미묘한 결을 인간 평론가에 준하는 수준으로 포착합니다.
- **고전 문학 (한자어/고어 맥락 추론):** 훈민정음 표기나 17~18세기 가사 문학의 난해한 고어가 나왔을 때, 주변 맥락과 당시 시대 상(예: 강호가도, 연군지정)을 결합하여 **가장 정교한 현대어 플래시백(번역) 및 해설**을 제공합니다. 억지스러운 직역보다 문맥에 맞는 의역에 강합니다.

### 2. GPT-5.5: "출제 메커니즘과 매력적인 오답 설계의 마스터" (문항/선지 개발 우위)

- **출제 의도 및 메커니즘 부합:** GPT-5.5는 교육과정 성취기준(화법과 작문, 독서, 문학 등)과 한국교육과정평가원의 출제 매뉴얼을 코드화하여 이해하는 능력이 뛰어납니다. 문항의 구조를 정밀하게 짚습니다.
- **오답 선지(Distractor) 제작:** 국어 시험의 핵심은 '그럴싸한 오답'을 만드는 것입니다. GPT-5.5는 '수험생이 가장 많이 낚이는 논리적 함정(예: 인과관계 뒤틀기, 주객전도)'을 정교하게 설계하여 매력적인 선택지를 만들어내는 능력이 Opus보다 뛰어납니다.

## 🔄 두 모델을 번갈아 사용할 때의 장점 (시너지)

### 1. '해석의 오류'가 없는 무결한 해설지 작성

- **워크플로우:** [4.8 Opus] 작품 독해 및 핵심 코드 추출 → [GPT-5.5] 이를 바탕으로 해설 및 보기 ( <보기> ) 지문 작성
- 현대시의 특정 구절을 두고 AI가 자의적으로 해석하는 '환각(Hallucination)' 현상이 국어 문항에서는 치명적인 출제 오류로 이어집니다. Opus가 1차로 작품을 엄밀하게 플롯 분석하고 주석을 달아주면, GPT-5.5가 이를 수능형 문항 레이아웃으로 가공합니다. 이 과정에서 **텍스트 해석의 오류가 완벽히 스크리닝**됩니다.

### 2. 고전 시가의 다의성을 활용한 고난도 변형 문제 개발

- 고전문학은 해석의 관점에 따라 이중적인 의미를 지니는 경우가 많습니다. Opus에게 "이 고전 가사에서 학계의 이설(異說)이 있거나 다의적으로 해석될 수 있는 시어를 찾아줘"라고 요청해 소스를 얻은 뒤, GPT-5.5에게 "이 다의성을 바탕으로 수험생이 착각하기 쉬운 <보기> 적용형 3점짜리 고난도 문항을 설계해줘"라고 넘기면, 평가원 기출에 필적하는 퀄리티의 변형 문항이 나옵니다.

## ⚠️ 번갈아 사용할 때의 단점 및 주의점 (리스크)

### 1. 텍스트 톤앤매너(Tone & Manner)의 괴리

- **문제점:** Opus는 문학을 대할 때 다소 진중하고 깊이 있는 비평가적 어조를 사용하는 반면, GPT-5.5는 철저히 시험 지침서나 EBS 해설지 같은 딱딱하고 구조화된 어조를 씁니다.
- **결과:** 두 모델의 결과물을 가이드라인 없이 섞으면 해설지의 앞부분(시어의 의미)은 굉장히 감성적이고 깊이 있는데, 뒷부분(선지 분석)은 기계적인 논리 구조로 쪼개져 있어 **전체 문항의 텍스트 일관성이 깨집니다**. 검토하는 인간 연구원의 피로도가 증가할 수 있습니다.

## 2. '기준선(Ground Truth)'의 붕괴

- 고어 해석이나 현대시의 상징은 100% 정답이 없는 경우가 있습니다. GPT-5.5가 A라는 관점으로 선지를 구성해 나가는 와중에, 중간에 컨텍스트를 이어받은 Opus가 "이 시어는 B로 보는 것이 더 타당하다"며 중간에 논리를 틀어버릴 수 있습니다. 이 경우 **질문(발문)과 정답 선지 간의 논리적 연결고리가 끊어진 불량 문항**이 탄생합니다.

### 🔗 추천하는 프로젝트 세부 운용 가이드

국어 문항 개발 프로젝트에서는 철저하게 "**Opus를 싱크탱크(내용 전문가)로, GPT-5.5를 아키텍트(출제 전문가)로**" 포지셔닝해야 성공합니다.

단계	작업 내용	주력 모델	이유
1단계: 원전 분석	고어/한자어 현대어역, 현대시 시어 상징성 추출, 화자의 정서 및 태도 정리	Claude 4.8 Opus	텍스트의 미묘한 결을 가장 왜곡 없이 인간답게 읽어냄.
2단계: 문항 설계	발문(질문) 작성, <보기> 지문 구성, 정답 선지 및 핵심 오답 메커니즘 빌딩	GPT-5.5	평가원 특유의 출제 프레임과 함정 설계에 압도적 강세.
3단계: 교차 검증	개발된 문항을 입력하여 "출제 오류나 이의제기 가능성"이 없는지 최종 검토	Claude 4.8 Opus	논리적 비약이나 지나치게 기계적인 오답 설계의 모순을 잡아냄.

💡 **핵심 팁:** 모델을 교대할 때, 이전 모델이 낸 결과물만 톡 던져주면 안 됩니다. 반드시 프롬프트 서두에 **[목표: 2027학년도 수능 국어 영역 수준의 문항 개발 / 기준: 1단계에서 확정된 작품 해석 가이드를 절대적으로 따를 것]**과 같은 **앵커(Anchor, 닷)** 역할을 하는 제약 조건을 명시해 주어야 두 괴물 모델 사이의 컨텍스트 충돌을 막을 수 있습니다.

## 📊 모델 사양·벤치마크 비교 (검증, 2026-06-06 기준)

### 🕒 웹 검증 결과

두 모델 모두 현행 모델로 확인됨. **공식 한국어 벤치마크는 두 모델 모두 부재**하며, 한자·고어 해석에 대한 정량 비교 자료도 현재 존재하지 않음 → 사실 검증은 반드시 1차 출처 대조로 보완할 것.

항목	OpenAI GPT-5.5	Anthropic Claude Opus 4.8
출시일	2026-04-23	2026-05-28
출력 토큰 단가	\$30 / 1M	\$25 / 1M
다분야·인문 추론 (HLE)	기준	<b>+7~8점 우위</b>
장문 맥락 (1M 토큰, GraphWalks)	45.4%	<b>68.1%</b>
사실 정합성 / 환각	—	<b>결함 통과율 전세대비 1/4, 불확실성 명시</b>
터미널·에이전트 코딩	<b>Terminal-Bench 우위</b>	SWE-bench Pro 우위
한국어 유창성	<b>새 사전학습, GPT-5 계열 최상</b>	자연스러우나 한국어 토큰 소모 다소 ↑
출제 관점 강점	발문·오답 설계·문항 구조화	원전 해석·정합성·교차검증

**해석:** 국어 출제에 직결되는 **인문 추론·장문 맥락 낮은 환각**은 Opus 4.8이 우세(→ 원전 분석·검증 앵커), **한국어 표현 유창성·출제 프레임**은 GPT-5.5가 우세(→ 초안·윤문·문항 설계). 위 1·2·3단계 역할 분담과 정확히 일치한다.

## 출제 루브릭 템플릿 (두 모델에 매번 주입)

### 사용법

모델 교대 시 컨텍스트 충돌·기준선 붕괴를 막기 위해, 아래 블록을 **모든 프롬프트 서두에 고정 주입**한다. 프로젝트별로 [ ] 항목만 채워 재사용.

[프로젝트 앵커 - 절대 준수]

- 목표 수준: [예: 2027학년도 수능 국어 / 6월 모평 수준]
- 영역/성취기준: [문학 / 독서 / 화작 / 연애] - [관련 성취기준 코드]
- 확정 해석 가이드: [1단계에서 Opus가 확정된 작품 해석을 절대 기준으로 삼을 것. 임의 재해석 금지]
- 정답 유일성: 정답은 1개, 나머지 선지는 명백히 오답이어야 함(이의제기 차단)
- 오답 메커니즘: 각 오답이 '왜 매력적인지' + '왜 틀렸는지'를 1줄씩 명시
- 근거 표기: 모든 해석은 지문 내 근거(행/구절) 인용. 외부 사실(한자 훈음·출전)은 [확실/불확실] 라벨링
- 문체 통일: 해설 어조 = [EBS 해설지체 / 평이한 설명체] 로 고정
- 금지: 환각성 단정, 학계 이설을 정설처럼 서술, 발문과 무관한 선지

### 품질 체크리스트(최종 검토용):

- 정답이 유일하고, 지문 근거로 100% 방어되는가
- 오답 4개가 각각 다른 함정 유형(인과뒤틀기·주객전도·범위확대·시제오류 등)인가
- 한자 훈·음, 고어 어형, 출전이 1차 출처로 검증되었는가
- 발문-<보기>-선지의 논리 연결이 끊김 없이 일관되는가
- 해설 전반의 어조·용어가 통일되어 있는가
- 이의제기·복수정답 소지가 없는가

## 용어 통일 글로서리 템플릿

### 목적

GPT-5.5(EBS 지침서체)와 Opus 4.8(비평가체)의 **용어 표류**를 막는 단일 사전. 두 모델 모두 아래 '표준 용어'만 사용하도록 프롬프트에 첨부한다. 프로젝트 진행 중 새 용어가 나오면 이 표에 누적.

개념	표준 용어 (채택)	혼용 금지 변형	비고
시 속 발화 주체	시적 화자	서정적 자아, 화자, 시인	'시인'과 구분
작품의 분위기	정서 / 분위기	무드, 톤	
표현상 특징	표현상의 특징	수사, 기법	선지 발문 통일
자연 친화 주제	강호가도(江湖歌道)	자연예찬, 전원생활	고전 한정
임금 향한 충정	연군지정(戀君之情)	충성심, 임 그리움	고전 한정
반대로 말하기	반어(反語)	아이러니, 비꼼	
모순 표현	역설(逆說)	패러독스, 모순어법	
보조 자료	<보기>	지문2, 참고자료	꺾쇠 표기 고정
점수 표기	3점 / 2점	[고난도], (배점)	
현대어 풀이	현대어역	번역, 의역	고어·한문 한정

구분	표준 표기 규칙
한자 병기	한글(漢字) 형식, 최초 1회만 병기
고어 인용	원문 그대로 + 괄호 안 현대어역
출전 표기	[작품명], 『문헌명』 (홀낫표/겹낫표 구분)
선지 번호	① ② ③ ④ ⑤ (원문자 고정)

### 📖 한 줄 운용 원칙

**Opus 4.8 = 내용·검증 앵커, GPT-5.5 = 설계·윤문 엔진.** 교대 시 [루브릭 앵커] + [글로서리]를 항상 동봉하고, 한자·고어·출전은 모델 신뢰 금지·1차 출처 대조 필수.