

딥시크가 불러온 AI전쟁 2.0시대 소버린 AI 전략

2025. 3

하정우, PhD: jungwoo.ha@navercorp.com

네이버클라우드 AI Innovation 센터장

네이버 Future AI 센터장

과실연 공동대표 겸 AI 미래포럼 공동의장

한국공학한림원 컴퓨팅분과 정회원

<AI전쟁>, <2025 AI 대전환> 작가

NAVER Cloud

생성형 AI의 패러다임의 변화

LLM 1단계: Knowledge AI의 시대
베이스 LLM 모델: 글 잘쓰는 AI

Pre-training + Post-training

GPT4o, Claude 3.5, Gemini, DeepSeek v3,
Qwen2.5-Max

HyperCLOVA X, EXAONE 3.5

강화학습

LLM 2단계: Thinking AI의 시대
추론 강화 LLM 모델:
아주 긴 논리/수리적 추론 능력 보유 AI

강화학습 (아주 긴 CoT 데이터)

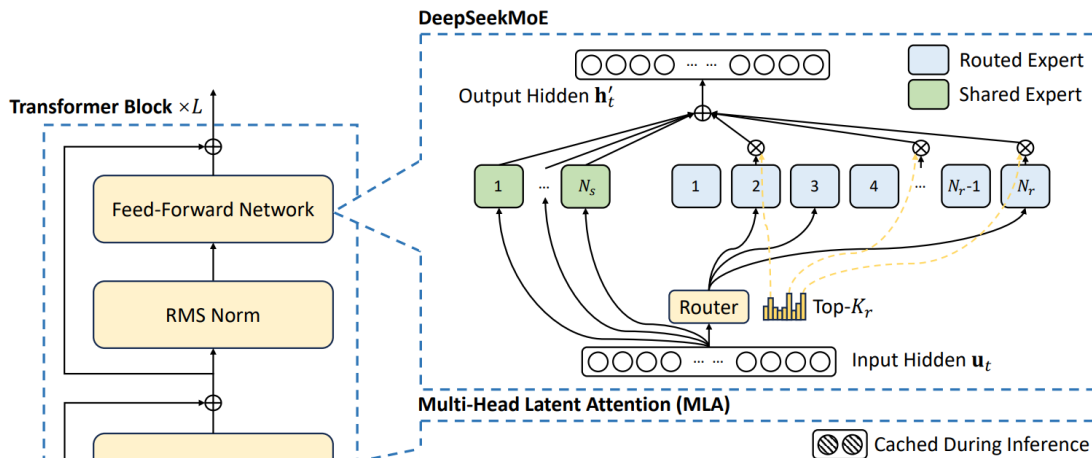
강력한 베이스 LLM 모델 필수

o1, o3, Deepseek-R1

DeepSeek 주요 시 현황

DeepSeek V1 (24. 1)

새로운 오픈소스 LLM



DeepSeek R1 (25. 1)

2세대 Thinking AI

그룹 활용 (GRPO)

향상된 성능

Training Costs	Pre-Training	Context Extension	Post-Training	Total
in H800 GPU Hours	2664K	119K	5K	2788K
in USD	\$5.328M	\$0.238M	\$0.01M	\$5.576M

80억원

H100 환산
5500장
15일 학습

기술 노하우
및 반복횟수
고려해야

Table 1 | Training costs of DeepSeek-V3, assuming the rental price of H800 is \$2 per GPU hour.

Input Hidden h_t ○○○○ ... ○○○○

딥시크 충격: 잘못 알려진 오해와 아마도(?) 진실 (자세한 내용은 QR통해 유튜브를 보시면 됩니다)



- 개발 비용 80억원? No: 1회 학습비용일 뿐 그 동안 실패, 인건비, 데이터구축비 제외
- 그래서 아무나 할 수 있다? No: 1세대 Knowledge 기반의 전문 AI만 가능
- 정말 저렴한가? 운영비는 수배 비
- 이제 고가의 GPU, 고사양 GPU를 넣
- 메모리 사업은? HBM도 고효율
- 개인정보 위험? AI모델을 다운로
- 편향되어 있나?
- 글로벌 AI 경쟁이

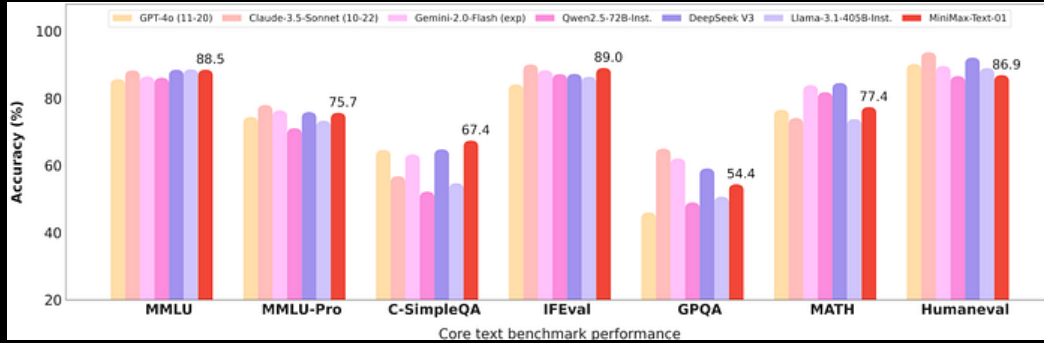
Providers for DeepSeek R1

OpenRouter routes requests to the top-ranked providers able to handle your prompts. ⓘ

Provider	Context	Max Output	Input	Output	Latency	Throughput
DeepSeek	64K	8K	\$0.55	\$2.19	27.18s	14.60t/s
DeepInfra	16K	16K	\$0.85	\$2.5	34.24s	3.07t/s
Together	164K	164K	\$7	\$7	22.69s	1.94t/s
Fireworks	164K	164K	\$8	\$8	4.27s	1.54t/s

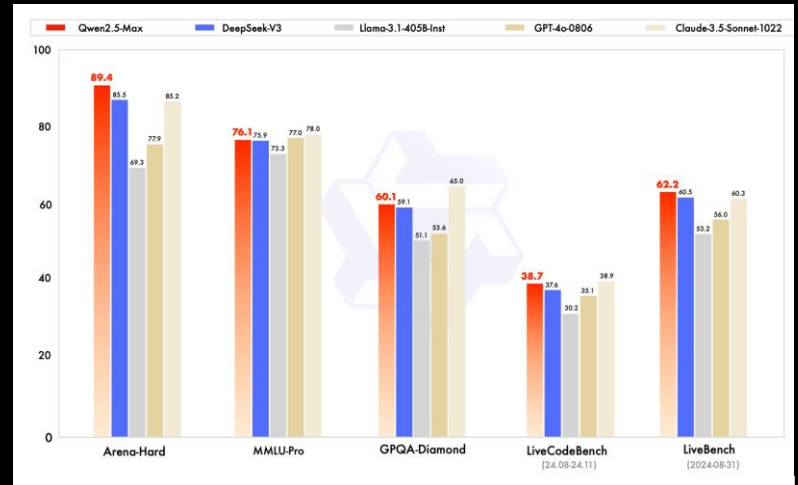
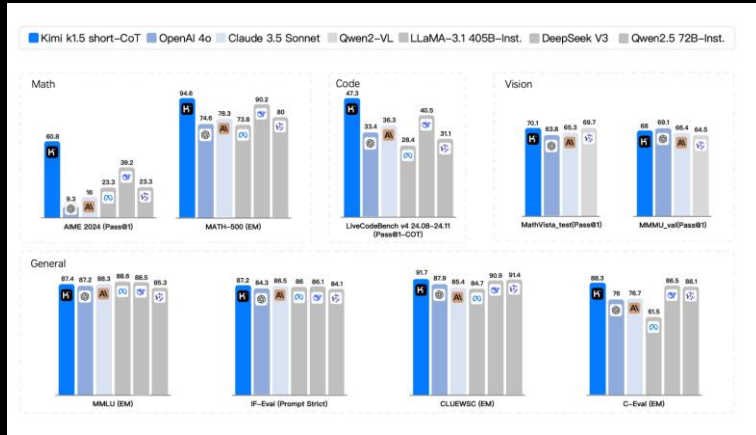
DeepSeek 뿐이 아니다. 강력한 다른 중국 AI들, GPT4o나 o1이 동네북

Minimax-Text-01



Qwen2.5-Max

Moonshot AI – Kimi-k1.5



딥시크 충격: 잘못 알려진 오해와 아마도(?) 진실 (자세한 내용은 QR통해 유튜브를 보시면 됩니다)



- 개발 비용 80억원? No: 1회 학습비용일 뿐 그 동안 실패, 인건비, 데이터구축비 제외
- 그래서 아무나 할 수 있다? No: 1세대 Knowledge 기반의 전문 AI만 가능

Providers for DeepSeek R1

OpenRouter routes requests to the top-ranked providers able to handle your prompts. ⓘ

Provider	Context	Max Output	Input	Output	Latency	Throughput
DeepSeek	64K	8K	\$0.55	\$2.19	27.18s	14.60t/s
DeepInfra	16K	16K	\$0.85	\$2.5	34.24s	3.07t/s
Together	164K	164K	\$7	\$7	22.69s	1.94t/s
Fireworks	164K	164K	\$8	\$8	4.27s	1.54t/s

- 정말 저렴한 고성능 AI 운영비는 수배 비
- 이제 고가의 GPU를 넘어서고, 고사양 GPU를 넘어서고
- 메모리 사업은? HBM도 고효율
- 개인정보 위험? AI모델을 다운로드
- 편향되어 있나?
- 글로벌 AI 경쟁이

param 모델보다
화학습에 더 많은
상
합재 & 확산. 고성능
어에 저장. 하지만

Second, the rule institutes new controls on the model weights of the most advanced closed-weight AI models. These controls will initially apply to the weights of models trained with 10^{26} computational operations or more, and authorizations will be required to export, reexport, or transfer (in-country) such weights to a broad set of countries. Additionally, the rule creates a new foreign direct product rule that applies these controls to certain model weights produced abroad using advanced computing chips made with U.S. technology or equipment. As with advanced computing chips, however, this rule includes several license exceptions for model weights:

GPT4급 이상, Gemini Ultra 이상, Claude 3.5 Sonnet 이상

- **Exception for deployments by U.S., ally and partner-headquartered entities:** New License Exception Artificial Intelligence Authorization (AIA) allows for the export, reexport, or transfer (in-country) of otherwise controlled closed AI model weights, without an authorization, by companies headquartered in the United States and certain allies and partners, except to an arms-embargoed country.
- **Exception for open models:** Models with widely available model weights (*i.e.*, open-weight models) are not subject to controls. Additionally, the model weights of closed models that are less powerful than the most advanced open-weight models, even if they exceed the 10^{26} threshold, are not controlled.

DeepSeek 이후 미국의 대응

ARTIFICIAL INTELLIGENCE

DeepSeek fallout: GOP Sen Josh Hawley seeks to cut off all US-China collaboration on AI development

DeepSeek's release of a new AI model that costs less to run than existing versions sent a chill through US markets

By **Morgan Phillips** · Fox News

Published January 29, 2025 10:20am EST



More From Fox News



Tension builds around Tulsi Gabbard's confirmation with key GOP senators undecided



미중의 AI 안전성 약화와 EU의 AI 레이스 참가 (300조 투자)

뉴스시스 구독

PICK

미국·영국, '포용·지속가능 AI' 공동성명 불참...美 "규제가 AI 죽일것"

입력 2025.02.12. 오전 11:11 · 수정 2025.02.12. 오후 12:39 [기사원문](#)

김승민 기자



댓글



카



영국 "거버넌스·안보 불충분...美와 무관"
EU도 '규제 완화'..."유럽, 세계와 동기화"



디지털타임스 구독

美中 경쟁 속 EU도 진흥 무계... 'AI정상회의'서 AI안전성 뒷전으로

입력 2025.02.12. 오후 6:45 · 수정 2025.02.12. 오후 7:01 [기사원문](#)

팽동현 기자



댓글



카



미국 압박에 따른 EU 규제 완화

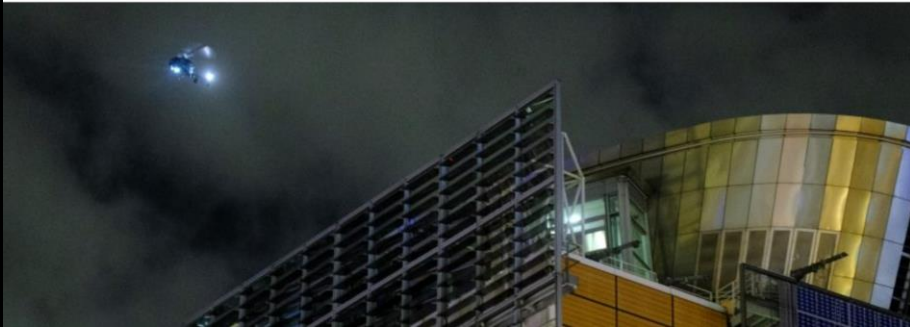
Tech

Commission withdraws AI liability directive after Vance attack on regulation

The Commission justified its removal by saying it could see "no foreseeable agreement" on the directive.

This article is exceptionally available for free! Want access to more exclusive content like this?
Discover all the benefits of Euractiv Pro.

Request a trial



Tech

Commission to withdraw proposed ePrivacy Regulation

The Commission's updated draft work programme confirms the withdrawal of the ePrivacy Regulation proposal



<https://www.euractiv.com/section/tech/news/commission-withdraws-ai-liability-directive-after-vance-attack-on-regulation/>
<https://www.euractiv.com/section/tech/news/commission-to-withdraw-proposed-eprivacy-regulation/>

글로벌 주요 국가 AI 투자 현황

미국 (25. 1): 스타게이트 (민간 중심 720조원)

EU (25. 2): 300조원 규모 AI 투자 (컴퓨팅 인프라와 CERN for AI에 30조원)

영국 (25. 1): AI Opportunities Action Plan. 컴퓨팅 인프라 정부 40조원, 민간 포함 71조원 투자

프랑스 (25. 2): 167조원 규모 투자 (캐나다와 UAE로 부터 700억 유로 투자 유치)

일본 (24. 11): AI와 반도체에 공적자금 2030년까지 91조원 규모 지원

소버린 AI의 범위? 생태계 전체 그리고 전략자산 영역 중요성

모든 산업 생성형 AI 확산 및 산업 생태계: 산업 경쟁력이 국가 AI 기술 경쟁력에 제한 (80점 vs. 90점)

생성형 AI 파인튜닝, 배포 및 활용 (미국 클라우드 위에서만)

범용생성형 AI
(미국 100점 vs. 미국 제공 80점 vs. 소버린 90점)

학습 데이터 및 서비스 데이터 - 미국의 Cloud Act
(데이터 유출가능성)

반도체, AI 데이터센터, 클라우드 (IaaS) (GPU 전략자산화)

에너지 (전력 공급 - 에너지 전략자산화, SMR but 나중에 탄소?)

전략자산화 진행 중

소버린 Foundation Model의 필요성: 공급망 리스크

반드시 경쟁력 있는 자체 기술 필요

강력한 범용 베이스

Reasoning foundation model
(선생님 모델)

- 전략자산화, 접근 제한 가능성
- AGI 보유국 or 종속국
- 제한된 라이선스 (국방, 에너지, CBRN 등)

Distillation

제조 특화 sLM

Distillation

공공 특화 sLM

Distillation

금융 특화 sLM

...

Distillation

법률 특화 sLM

Model	mixed	cont'	en	zh (/en)	others (/en)
LLaMA	N	N	50.36	15.98 (0.32)	17.95 (0.36)
Chinese-LLaMA	N	Y	29.38	21.12 (0.72)	5.88 (0.20)
LLaMA2	Y	N	73.29	43.57 (0.59)	35.78 (0.49)
Chinese-LLaMA2	Y	Y	60.93	32.32 (0.53)	22.59 (0.37)
LLaMA	N	N	50.36	15.98 (0.32)	17.95 (0.36)
LLaMA2	Y	N	73.29	43.57 (0.59)	35.78 (0.49)
Baichuan2-base	Y	N	87.58	59.30 (0.68)	38.61 (0.44)
Chinese-LLaMA	N	Y	29.38	21.12 (0.72)	5.88 (0.20)
Chinese-LLaMA2	Y	Y	60.93	32.32 (0.53)	22.59 (0.37)

Table 5: The comparison of the selected models' RA scores on the *Basic* knowledge (the mean of xCSQA and xCOPA scores), where "mixed" and "cont'" means having Chinese mixed or continued pretraining, "/en" means the ratio to the English scores, and "others" refers to the mean scores in the other 8 languages. The first lines in each division (LLaMA, LLaMA2, LLaMA and Chinese-LLaMA) are the baseline values (in black). The green values are higher than baseline, and the red ones are lower than baseline. (Same notations in below.)

Model	mixed	cont'	en	zh (/en)	others (/en)
LLaMA	N	N	79.28	7.49 (0.09)	42.19 (0.53)
Chinese-LLaMA	N	Y	44.23	13.01 (0.29)	15.44 (0.35)
LLaMA2	Y	N	91.58	40.39 (0.44)	56.00 (0.61)
Chinese-LLaMA2	Y	Y	79.84	45.92 (0.58)	43.70 (0.55)
LLaMA	N	N	79.28	7.49 (0.09)	42.19 (0.53)
LLaMA2	Y	N	91.58	40.39 (0.44)	56.00 (0.61)
Baichuan2-base	Y	N	86.34	65.81 (0.76)	50.51 (0.59)
Chinese-LLaMA	N	Y	44.23	13.01 (0.29)	15.44 (0.35)
Chinese-LLaMA2	Y	Y	79.84	45.92 (0.58)	43.70 (0.55)

Table 6: The comparison of the selected models' RA scores on the *Factual* knowledge (the mean of xGeo and xPeo scores).

소버린 시는 한국형 시를 넘어 다문화 포용적 시로 진화하기 위한 첫출발

