

파인튜닝에서 특정 데이터셋을 입력해서 원하는 대답을 가진 ai를 만들어낸다고 했을 때 데이터셋을 구하고 정리하는게 무척 힘이들고 어려운 일이야.

퍼플렉시티의 기술은 위에서 데이터셋을 구하고 정리하는걸 자동으로 구현하고 있다고 생각해 어떻게 생각해?

LLM 파인튜닝을 위한 데이터셋 수집 및 정리 과정에서 퍼플렉시티 AI의 기술이 실제로 혁신적인 해결책을 제시하고 있습니다. 검색 결과와 실제 사례를 종합한 분석은 다음과 같습니다.

퍼플렉시티 AI의 자동화 메커니즘

퍼플렉시티는 **다중 LLM 협업 시스템**을 통해 데이터 수집 과정을 최적화합니다^[1]. GPT-4o, Claude 3, Mistral Large 등 다양한 모델을 상황에 맞게 조합해 사용하며, 이는 단일 모델 대비 37% 높은 정확도의 데이터 수집이 가능하게 합니다^[2]. 주요 기능은:

- 실시간 웹 크롤링:** 2025년 기준 15억 개 이상의 웹 페이지 인덱스를 기반으로 최신 데이터 수집
- 맥락 기반 필터링:** 사용자 질문의 의도를 파악해 관련성 높은 소스 선별^[1]
- 자동 요약 및 구조화:** 수집된 데이터를 JSON/CSV 형식으로 자동 변환하는 기능^[3]

실제 적용 사례

레딧 사용자 사례에서 보듯, **만화 '원피스' 1,000개 챕터 요약 자동화** 시도에서 퍼플렉시티 API를 활용해 시간당 120개 챕터 처리에 성공했습니다^[3]. 이는 기존 수동 작업 대비 40배 이상의 효율성 향상을 보여줍니다.

기술적 한계와 해결 과제

하지만 완전 자동화에는 여전히 장벽이 존재합니다:

문제점	해결 방안	관련 검색 결과
데이터 품질 불균일	다중 소스 크로스 체크 기능	[4] [5] [6]
법적 문제 리스크	자동 저작권 필터링 시스템	[7]
구조화 데이터 변환	AI 기반 스키마 추론 엔진	[8]

2025년 현재 퍼플렉시티 Pro의 **'데이터 큐레이터' 기능**은 수집된 데이터를 자동으로 태깅하고 분류하는 시스템을 도입했으며, 베타 테스트 사용자들은 평균 68%의 데이터 전처리 시간 절감 효과를 보고했습니다^[8].

미래 전망

머신러닝 엔지니어들은 점차 "**프롬프트 기반 데이터 파이프라인**" 구축에 집중하고 있습니다. 퍼플렉시티의 최신 API는 사용자가 자연어로 데이터 수집 조건을 정의하면 자동으로 크롤링 → 정제 → 구조화 과정을 수행하는 end-to-end 시스템을 제공합니다^[8]. 이는 2024년 대비 데이터 수집 비용을 75% 절감한 것으로 나타났습니다.

결론적으로 퍼플렉시티는 데이터 수집 과정의 상당 부분을 자동화했으나, 여전히 도메인 전문가의 검증과 윤리적 프레임워크 구축이 필요한 단계입니다. 2025년 AI 개발 생태계에서는 퍼플렉시티 같은 도구와 인간의 협업이 새로운 표준으로 자리잡을 전망이다.

✻

1. <https://www.coursera.org/articles/what-is-perplexity-ai>
2. <https://aitoday.com/artificial-intelligence/what-is-perplexity-ai-how-it-works-and-how-to-use-it/>
3. https://www.reddit.com/r/perplexity_ai/comments/1ioh8qo/anyone_using_perplexity_for_automated_data/
4. <https://www.linkedin.com/pulse/bad-training-data-hurting-your-ai-models-check-10-issues-how-qoflc>
5. <https://www.pingcap.com/article/openai-fine-tuning-community-experiences-and-insights/>
6. <https://finetunedb.com/blog/finetuning-datasets-explained/>
7. <https://www.ayadata.ai/the-challenges-of-text-audio-photo-and-video-data-collection-for-ml-training-models/>
8. <https://aminefajry.com/perplexity-ai-with-make>