

켈리저널 키워드 뉴스 발행 시스템 분석 및 고도화 제안

1. 시스템 개요

켈리저널 시스템은 하루에 하나의 키워드를 정하여 관련 콘텐츠를 자동으로 수집 · 큐레이션하고 뉴스 형태로 발행하는 것을 목표로 한다. 편집자가 키워드를 입력하면 공식 뉴스, 블로그, 학술 자료, 커뮤니티 글 등 여러 출처의 콘텐츠를 수집하고, 그중 승인된 기사 15 건 이상을 선택하여 특집 페이지를 발행한다 【86275622104348†L2-L14】 . 모든 수집 결과는

selected_keyword_articles.json 파일에 저장되고, 편집자가 승인한 기사는

selected_articles.json 에 정리되어 Markdown/HTML 뉴스 페이지로 출력된다

【86275622104348†L18-L27】 . 발행을 위해서는 최소 15 개의 기사가 승인되어야 하며, 이 기준은 설정 옵션 require_editor_approval 을 통해 강제된다 【86275622104348†L162-L170】 .

2. 현재 구현과 워크플로우 분석

2.1 디렉터리 및 데이터 구조

- **작업본(selected_keyword_articles.json)** - 하루의 키워드에 대한 모든 수집 기사와 편집 여부를 저장하는 JSON 파일이다. 기사별로 제목, 요약, 출처, 날짜, 점수, approved 플래그, editor_note 등이 포함된다 【86275622104348†L18-L23】 .
- **발행본(selected_articles.json)** - 작업본에서 approved=True 인 기사만 추린 최종 발행 리스트이다. 최종 뉴스 페이지 생성 시 이 파일을 참조한다 【86275622104348†L24-L27】 .
- **archive/YYYY-MM-DD_KEYWORD.{json|html|md}** - 발행 완료 후 결과물을 날짜별로 저장하는 아카이브다. 동일 내용의 JSON · Markdown · HTML 파일을 모두 저장하여 후일 재검토가 가능하다 【86275622104348†L29-L35】 .
- **tools/** - 파이프라인 지원 스크립트들이 위치한 폴더로, 재구성, 보강, 승인 처리, 동기화, 포맷 교정 등에 사용하는 다양한 Python 스크립트를 포함한다 【86275622104348†L40-L88】 .

2.2 주요 스크립트 기능

보고서는 도구 폴더에 있는 중요한 스크립트를 다음과 같이 설명한다 【86275622104348†L40-L88】 .

스크립트	주요 역할
rebuild_kw_from_today.py	당일 수집된 원본 JSON 들을 표준 구조(딕셔너리 + articles 배열)로 재구성하여 selected_keyword_articles.json 작업본을 만든다 【86275622104348†L41-L45】 .
augment_from_official_pool.py	작업본의 기사 수나 승인된 기사가 부족할 경우 미발행 상태의 공식 기사 풀에서 추가 기사들을 채워 기사 총량을 20 개 수준으로 보정하고, 승인 기사도 최소 15 개 이상 확보하도록 돕

스크립트	주요 역할
force_approve_top20.py	는다 【86275622104348†L46-L49】 . 점수 순 등 정렬 기준으로 상위 20 건의 approved 필드를 강제로 True 로 설정하여 일괄 승인 처리를 지원한다 【86275622104348†L73-L77】 .
sync_selected_for_publish.py	selected_keyword_articles.json 에서 approved=True 인 기사만 추출하여 selected_articles.json 을 갱신하고, 기존 기사 메타데이터를 병합한다 【86275622104348†L79-L83】 .
repair_selection_files.py	selected_articles.json 등의 구조가 어긋난 경우 표준 딕셔너리 구조({"date": ..., "articles": [...]})로 교정해 데이터 포맷 오류를 예방한다 【86275622104348†L84-L88】 .

2.3 PowerShell 자동화 워크플로우

일일 발행 프로세스는 run_quali_today.ps1 스크립트에 의해 자동 실행된다. 주요 단계는 다음과 같다 【86275622104348†L89-L129】 .

1. **수집(Collect)** - 공식 소스와 커뮤니티 소스, 필요하면 키워드 기반 외부 RSS 를 모두 수집한다. orchestrator.py 의 --collect, --collect-community, --collect-keyword --use-external-rss 명령을 차례로 호출한다 【86275622104348†L94-L99】 .
2. **재구성(Rebuild)** - rebuild_kw_from_today.py 를 실행하여 수집된 원본 JSON 을 하나의 작업본으로 합치고 selected_keyword_articles.json 을 생성한다
【86275622104348†L100-L104】 .
3. **보충(Augment)** - 승인된 기사가 15 개 미만일 때 augment_from_official_pool.py 로 기사 수를 20 개 수준으로 보충한다 【86275622104348†L105-L109】 .
4. **승인 수 점검** - 파이썬 원라이너를 사용해 approved=True 인 기사 수를 집계하여 발행 조건 충족 여부를 확인한다 【86275622104348†L110-L114】 .
5. **동기화(Sync)** - repair_selection_files.py 로 포맷을 교정한 뒤 sync_selected_for_publish.py 를 실행해 승인 기사만 selected_articles.json 에 갱신한다 【86275622104348†L115-L120】 .
6. **발행(Publish)** - orchestrator.py --publish-keyword <키워드>를 실행하여 Markdown/HTML 페이지를 생성하고 아카이브에 저장한다. 승인 기사가 15 개 이상이면 정상 발행하고, 부족하면 require_editor_approval 을 임시로 false 로 바꿔 응급 발행을 수행한 후 즉시 복구한다 【86275622104348†L121-L129】 .

2.4 운영 정책과 설정

- **품질 게이트(QG)** - 발행 한 호마다 최소 15 개의 기사가 승인돼야 한다는 조건이 적용된다. 이 기준을 만족하지 못하면 발행 명령이 오류로 종료되거나 UI 에서 발행 버튼이 비활성화된다 【86275622104348†L162-L167】 . 권장 선정 기사 수는 15~20 개이다
【86275622104348†L162-L168】 .

- **require_editor_approval 설정** – config.json 의 features.require_editor_approval 값으로 편집자 승인 없이 발행이 가능할지 제어한다. 기본값은 True 이며, 응급 상황에서만 임시로 False 로 변경한다 【86275622104348†L169-L173】 .
- **커뮤니티 필터** – 커뮤니티 수집은 잡음을 줄이기 위해 필터링을 적용한다. 키워드 포함 여부, 제목 길이, 금지 패턴, 최소 추천수/댓글수, 점수 임계값 등을 설정하며, 현재 score_threshold 는 약 3.5~4.2 수준에서 조정된다 【86275622104348†L174-L181】 .
- **공식 소스 관리** – feeds/official_sources.json 의 RSS/뉴스 피드 URL 은 정기적으로 점검해 404 오류나 구조 변경으로 인한 문제를 해결해야 한다 【86275622104348†L183-L188】 .
- **에러 대응** – PowerShell 에서 f-string 을 사용하면 중괄호가 PowerShell 포매팅으로 처리되어 오류가 발생할 수 있으므로 heredoc 방식으로 Python 코드를 실행하도록 권장한다 【86275622104348†L193-L200】 . 또한 파이썬 원라인 호출 시 경로 · 인코딩 이슈를 피하기 위해 python -c 또는 here-string 을 사용하는 것을 권장한다 【86275622104348†L201-L205】 .
- **데이터 일관성 및 백업** – selected_articles.json 과 selected_community.json 은 항상 {"date": "YYYY-MM-DD", "articles": [...]} 구조를 유지해야 하며, 수동 편집 후에는 sync_selected_for_publish.py 를 다시 실행해 동기화를 보장해야 한다 【86275622104348†L238-L257】 . 아카이브 폴더는 백업 · 압축 등을 통해 용량 관리를 해야 하고, 동일 키워드로 여러 번 발행할 경우 시간까지 포함한 파일명을 사용한다 【86275622104348†L258-L261】 .

2.5 보고서에서 제시한 개선사항

보고서 말미에는 향후 발전 방향으로 키워드 자동 회전/추천, 스케줄러 연동, 관리자용 편집 UI, Elasticsearch 연동, Docker 이미지화 및 CI/CD 파이프라인, 코드 리팩터링, 로그 관리, 단위 테스트, SEO 최적화 등이 제안되었다 【86275622104348†L264-L305】 . 이러한 제안들은 고도화 로드맵의 출발점으로 활용할 수 있다.

3. 시스템 강점과 문제점 분석

3.1 강점

1. **모듈화된 파이프라인** – 수집 – 재구성 – 보충 – 승인 – 발행으로 이어지는 단계가 명확히 구분되어 있으며, 각 단계를 별도의 스크립트로 분리하여 유지보수가 용이하다.
2. **품질 게이트 도입** – 최소 승인 기사 수(≥ 15)를 강제하는 품질 게이트가 있어 뉴스 특집의 내용과 분량을 일정 수준 이상으로 유지할 수 있다 【86275622104348†L162-L170】 .
3. **공식/커뮤니티 출처 분리** – 커뮤니티 글은 필터링 기준과 점수 임계값을 적용해 잡음을 최소화하고, 공식 소스는 별도의 보강 스크립트로 품질을 유지한다 【86275622104348†L174-L181】 .
4. **자동화 스크립트 제공** – run_quali_today.ps1 를 통해 원클릭으로 전체 워크플로우를 실행할 수 있으며, 응급 발행 모드를 포함해 예외 처리까지 자동화했다 【86275622104348†L121-L129】 .
5. **아카이브 구조** – 모든 발행 결과를 날짜별로 저장하여 히스토리를 보존하고, JSON/Markdown/HTML 형태로 남겨 후속 분석과 재가공이 가능하다 【86275622104348†L29-L35】 .

3.2 문제점 및 개선 필요 영역

1. **편집자 의존성** - 최소 15 개 승인이라는 기준 때문에 편집자의 수동 검수 및 승인 작업이 필수적이다. 긴급 모드가 존재하지만 이는 임시 방편이며, 근본적으로 수집 단계에서 높은 품질의 기사 선별을 자동화할 필요가 있다.
2. **정적 임계값** - 커뮤니티 점수 임계값과 최소 기사 수 등이 정적으로 설정돼 있어 키워드에 따라 지나치게 엄격하거나 느슨할 수 있다. 상황에 따라 동적으로 조정하거나 학습 기반 모델이 필요하다.
3. **RSS 소스 불안정** - 공식 소스 URL 이 변경될 때 수집이 실패할 수 있으므로, 자동 검증 및 알림 시스템이 필요하다 【86275622104348†L183-L188】 .
4. **PowerShell 기반 자동화** - Windows 환경에 종속된 PowerShell 스크립트는 f-string 오류와 경로 · 인코딩 문제 등 유지보수 부담이 크다 【86275622104348†L193-L205】 .
5. **편집 UI 부재** - 기사 선택과 코멘트 입력을 위해 직접 JSON 을 수정해야 하므로 사용자 경험이 낮다. 보고서에서도 관리자용 웹 UI 개발을 제안하고 있다 【86275622104348†L288-L297】 .
6. **트렌드 분석 및 추천 기능 부족** - 키워드를 운영자가 수동 지정해야 하므로 트렌드 변화에 대한 대응이 늦다 【86275622104348†L264-L268】 .
7. **검색 및 데이터 탐색 기능 미흡** - 현재 JSON 파일 기반으로 목록을 불러오기 때문에 빠른 검색이나 연관 기사 탐색이 어렵다.
8. **코드 품질과 테스트** - 장기적 확장성을 위해 코드 리팩터링과 에러 로그 관리, 단위 테스트 추가가 필요하다 【86275622104348†L299-L304】 .

4. 개선점 및 고도화 방향 제안

4.1 기사 수집·필터링 고도화

LLM 기반 선별 및 요약 도입 - 전통적 키워드 매칭 기반 자동화는 잡음을 걸러내지만 편집자의 뉴스 가치 판단을 대신하기 어렵다. 2025 년 연구에서는 대형 언어모델(LLM)이 시의성, 영향력, 논쟁성, 대중성 등 저널리즘의 뉴스 가치를 프롬프트에 인코딩해 잠재적 기사 리드를 추출 · 평가하는 것이 가능함을 보여주었다 【763569562464242†L198-L214】 . 또한 LLM 은 요약, 감성 분석, 팩트체크 같은 복잡한 작업을 수행할 수 있으며, 프롬프트를 적절히 설계하면 자동 큐레이션 성능을 크게 향상시킬 수 있다 【763569562464242†L198-L214】 .

보고서에서 이미 force_approve_top20.py 와 같은 점수 기반 일괄 승인을 제공하고 있지만, LLM 을 활용해 기사 내용을 이해한 뒤 뉴스 가치가 높은 기사만 추천하도록 개선할 수 있다. 예를 들어 다음과 같은 절차를 도입할 수 있다.

1. **LLM 프롬프트 설계** - 기사 본문을 입력으로 받아 ① 시의성(최근성), ② 산업/표준에 미치는 영향, ③ 논쟁성, ④ 일반 독자의 관심도를 평가하도록 LLM 에게 지시한다. 연구에서는 이러한 네 가지 요소를 기반으로 프롬프트를 반복적으로 개선하여 모델의 판단을 인간 편집자와 유사하게 만들었다 【763569562464242†L257-L261】 .
2. **요약 및 키포인트 생성** - LLM 이 기사 내용을 요약하고 주요 포인트를 Bullet 형식으로 생성하게 한다. 뉴스업계 사례에 따르면 야후 뉴스와 월스트리트저널은 AI 요약을 '편의 기능'으로 제공하여 독자들이 기사 핵심을 빠르게 파악하도록 돕고 있으며, AI 요약은 독자의 참여도 향상에 기여한다 【358394435857084†L97-L115】 . 요약은 원문에 있는 정보만 사용하도록 제약해 오류 발생을 줄이고, 여러 단계의 품질 검수를 거쳐야 한다 【358394435857084†L97-L115】 .

3. **에디터 검수 보조** - LLM 이 생성한 평가와 요약참고하여 편집자는 승인 여부를 결정한다. 월스트리트저널의 사례처럼 AI 가 생성한 요약이나 키포인트는 반드시 편집자가 검토해야 하며, 인간이 루프에서 빠지는 것은 바람직하지 않다 【358394435857084†L166-L175】. 따라서 LLM 은 '첫 필터'로 활용하고 최종 결정은 편집자가 내리는 구조를 유지한다.

4.2 키워드 추천 및 트렌드 분석

보고서에서는 향후 기능으로 자동 키워드 회전 및 추천을 제시하였다 【86275622104348†L264-L268】. 이를 구현하기 위해 다음 방안을 제안한다:

- **트렌드 분석 모듈** - 전일 또는 최근 일주일 동안 수집한 기사 데이터를 분석하여 가장 많이 언급된 표준명이나 산업 키워드를 추출하고, 외부 트렌드 데이터(구글 트렌드 API 등)에 연동해 인기 키워드를 추천한다.
- **LLM 기반 주제 추출** - LLM 에게 “어제 발행된 기사들에서 공통되는 주제나 떠오르는 이슈를 요약하고 다음 주제 키워드를 추천하라”는 프롬프트를 제공할 수 있다. LLM 은 맥락을 파악해 다음날 발행에 적합한 후보 키워드를 제안할 수 있다.
- **키워드 풀 로테이션** - 운영자가 미리 정의한 키워드 풀을 주간 또는 월간 주기로 자동 회전하도록 하여 수동 개입을 줄인다. 추천 결과와 로테이션 풀을 결합해 최종 키워드를 선택한다.

4.3 검색 및 데이터 관리 향상

Elasticsearch(ES) 연동 - 보고서는 장기적으로 검색 최적화와 대용량 데이터 대응을 위해 Elasticsearch 와의 연동을 고려하고 있다 【86275622104348†L299-L305】. 학계 연구에 따르면 Elasticsearch 를 Retrieval-Augmented Generation(RAG) 프레임워크에 적용하면 기존 TF-IDF 기반 검색보다 정확도와 효율이 향상되며, 복잡한 쿼리를 유연하게 처리할 수 있다 【669341352768893†L51-L66】.

따라서 다음과 같은 구조를 제안한다:

- 발행된 기사 및 원본 데이터 전체를 ES 인덱스로 저장한다. 인덱싱 시 키워드, 요약, 출처, 날짜, 점수, 섹션 등을 필드로 지정한다.
- 편집자 UI 에서 검색창을 제공하여 과거 기사, 유사 기사, 특정 키워드를 빠르게 탐색할 수 있게 한다. ES 의 퍼지 검색이나 벡터 검색(semantic search)을 활용하면 표준번호 변형(예: “IPC A610”, “IPC-A-610”)도 쉽게 찾을 수 있다.
- LLM 과 연계한 ES-RAG 를 구현하면 질문 응답이나 요약 생성 시 관련 문서를 효과적으로 검색할 수 있어 향후 챗봇 서비스 확장에 도움이 된다 【669341352768893†L51-L66】.

4.4 자동화 및 운영 개선

- **일정 기반 실행** - 매일 아침 정해진 시간에 자동으로 수집·발행 작업을 실행하도록 윈도우 작업 스케줄러 또는 cron 에 등록하는 것이 권장된다 【86275622104348†L282-L287】. 또한 작업 결과를 이메일이나 슬랙 웹훅으로 통지하여 실패 시 빠르게 대응할 수 있게 한다.
- **플랫폼 독립적인 스크립트로 전환** - PowerShell 특유의 f-string 오류 및 경로 문제를 해소하기 위해 전체 자동화 스크립트를 Python 으로 통합하거나, 오픈소스 워크플로우 도구(n8n 등)로 이동하는 것을 고려한다. 연구 사례에서는 n8n 을 사용해 RSS 모니터링, 크롤

링, LLM 처리, 결과 저장을 자동화했다 【763569562464242†L294-L305】. 이를 활용하면 GUI 없이도 크로스 플랫폼에서 안정적으로 파이프라인을 운영할 수 있다.

- **컨테이너화 및 CI/CD** - Docker 이미지화와 CI/CD 파이프라인을 통해 배포를 일원화하고 의존성 문제를 줄인다 【86275622104348†L299-L305】. GitHub Actions 등을 이용해 테스트와 빌드, 배포 과정을 자동화할 수 있다.
- **모니터링과 로깅 강화** - 현재는 PowerShell 에서 \$ErrorActionPreference = "Stop"으로 중단시키고 로그를 남기는 수준이다 【86275622104348†L230-L235】. 중앙 로깅 서버나 Sentry 같은 서비스와 연동해 에러 발생 시 알림을 받고, 주요 지표(수집 기사 수, 승인 기사 수, 발행 성공 여부)를 대시보드로 시각화하면 운영 신뢰성을 높일 수 있다.

4.5 편집자 인터페이스 및 협업 도구

보고서에서 제안한 관리자용 웹 UI 는 가장 시급한 개선 사항 중 하나이다

【86275622104348†L288-L297】. Flask/FastAPI 백엔드와 React/Vue 기반 대시보드를 구축하여 다음 기능을 제공할 수 있다:

- **기사 목록 표시 및 필터링** - selected_keyword_articles.json 을 테이블로 불러와 기사별 제목, 출처, 점수, 승인 상태를 보여주고, 체크박스로 승인/보류를 선택할 수 있게 한다.
- **편집자 코멘트 입력** - 각 기사에 한 줄짜리 editor_note 를 입력할 수 있는 필드와, 코멘트 미입력 시 경고 메시지를 제공한다.
- **키워드 입력 및 발행 버튼** - 키워드를 입력하면 수집이 자동 트리거되며, 승인 완료 후 발행 버튼을 눌러 결과물을 생성하도록 한다.
- **역사 관리와 협업** - 이전 발행 기록을 조회하고, 여러 편집자가 동시에 작업할 때 충돌을 방지하는 잠금 또는 버전 관리 기능을 제공한다.

이러한 UI 는 편집자의 업무 효율을 높이고, JSON 파일을 직접 다루며 발생하는 오류를 줄여준다.

4.6 코드 품질 및 테스트 강화

- **모듈화와 리팩터링** - engine_core.py 와 orchestrator.py 등 주요 코드 베이스를 모듈화해 테스트 가능성을 높인다. 함수 단위로 책임을 분리하고, 공통 유틸 함수는 별도 모듈에 배치한다.
- **단위 테스트 도입** - 수집, 파싱, 보강, 발행 등 주요 함수에 대해 단위 테스트를 작성하고, GitHub Actions 를 통해 지속적 통합 테스트를 수행한다 【86275622104348†L299-L304】.
- **JSON 스키마 검증** - selected_keyword_articles.json 과 selected_articles.json 을 저장하기 전에 JSON Schema 로 형식을 검증해 포맷 오류를 사전에 방지한다. PowerShell 에서 JSON 을 수정할 때는 반드시 ConvertFrom-Json/ConvertTo-Json 을 사용하라는 보고서의 권장 사항을 자동화 도구에서도 적용한다 【86275622104348†L245-L251】.

4.7 SEO 및 배포 전략

- **SEO 메타데이터** - 발행된 HTML 페이지에 구조화된 데이터(schema.org 기사 스키마)와 메타태그를 추가하여 검색 엔진 친화성을 높인다. Nieman Lab 기사에서는 AI 요약이 검색 인덱싱에 도움을 줄 수 있다는 의견도 제시했다 【358394435857084†L63-L70】.
- **웹사이트 통합** - 메인 웹사이트에 키워드 뉴스 섹션을 추가하고, 최신 발행물을 쉽게 탐색할 수 있도록 개선한다.

4.8 다국어 지원과 번역

시스템은 현재 OpenAI API 를 사용해 번역을 제공하고 있으나, 향후 한국어 외에 다양한 언어로 발행을 확대하려면 번역 품질 관리가 중요하다. 월스트리트저널 사례에서는 AI 번역을 이용해 한국어·일본어 제품을 출시하는 등 다국어 확장 효과를 보고 있다 【358394435857084†L187-L196】. 번역 모델 업데이트 시 품질 테스트를 진행하고, 편집자가 번역 내용을 검수할 수 있는 워크플로우를 마련해야 한다.

4.9 사용자 피드백 및 학습 기반 개선

독자들이 어떤 기사에 관심을 보였는지, 어느 요약이 도움이 되었는지 로그를 수집하고, 추천 모델을 개선하는 데 활용할 수 있다. 야후 뉴스는 AI 요약 기능 출시 후 사용자 참여가 50% 증가하고 1인당 체류 시간이 165% 늘어난 것으로 보고했다 【358394435857084†L97-L115】. 퀄리저널도 발행 페이지에 피드백 기능을 추가해 요약의 품질과 기사 선택 기준을 지속적으로 개선할 수 있다.

5. 결론

퀄리저널 키워드 뉴스 발행 시스템은 명확한 파이프라인과 품질 게이트를 갖춘 견고한 기반을 마련했지만, 변화하는 미디어 환경과 독자 요구에 맞추어 지속적인 고도화가 필요하다. LLM 을 활용한 기사 평가와 요약, 자동 키워드 추천, Elasticsearch 기반 검색, 일정 기반 자동화, 관리자용 UI 구축, 코드 품질 개선, SEO 최적화 등은 시스템의 안정성과 확장성을 크게 높일 것이다. 이러한 개선을 단계적으로 적용하면서도 인간 편집자의 판단과 책임을 중심에 두어야 하며, AI 는 편집자를 보조하는 도구로 활용해야 한다 【358394435857084†L166-L175】. 발행 결과와 사용자 피드백을 기반으로 알고리즘과 프롬프트를 지속적으로 개선하는 선순환을 통해, 퀄리저널 시스템을 더욱 신뢰할 수 있고 효율적인 뉴스 발행 플랫폼으로 발전시킬 수 있을 것이다.