

산업용 시험성적서 관리를 위한 LLM 기반 자동화 시스템 설계

유정문¹, 안제혁², 김형진³, 이재민⁴, 김동성*

한국전력기술¹, 금오공과대학교 IT융복합공학과^{2,4,*}, 주식회사 엔에스랩 기술연구소^{3,*}
rjm@kepco-enc.com¹, {wpgur2018², ljmpaul⁴, dskim*}@kumoh.ac.kr, haengg@nslab.tech³

Design of LLM-based Automation System for Industrial Test Certificate Management

Jung-Moon Ryu¹, Je-Hyeok Ahn², Hyeong-Jin Kim³, Jae-Min Lee⁴, and Dong-Seong Kim*

KEPCO Engineering & Construction company¹

Kumoh National Institute of Technology Dept. of IT Convergence Eng.^{2,4,*}

NSLab Co., Ltd. Technology Research Institute^{3,*}

요약

산업현장에서 사용되는 중요 기자재는 국제 및 국내 규격 부합 여부를 검증하기 위해 시험성적서 제출이 의무화되어 있다. 그러나 방대한 시험성적서를 엔지니어가 육안으로 검증하는 것은 쉽지 않아 관리의 비효율성이 발생하며 문서 위변조 문제가 발생하기도 한다. 이를 개선하기 위한 방안중 하나로 OCR(Optical Character Recognition, 광학 문자 인식) 기술을 적용하고 있지만, 기존의 OCR 기술은 서식·서체의 다양성과 비정형화된 레이아웃, 다국어 환경 등에서 인식을 저하, 오인식, 후처리 부담 증가와 같은 한계점에 노출되어 있다. 이에 본 논문은 Transformer 기반 Donut 모델을 활용하여 시험성적서를 JSON 형태로 자동 정규화하고, 이를 Fine-Tuning 한 LLM을 통해 엔지니어가 손쉽게 확인할 수 있도록 지원하는 시스템을 제안한다. 또한 이후 등록 되는 동일 문서에 대해 텍스트 대조를 통한 위변조 검증 기능을 제공하여, 시험성적서 관리의 신뢰성과 효율성을 동시에 확보할 수 있음을 확인하였다.

I. 서론

플랜트(Plant) 산업 현장에서는 건설과 운영 과정에서 사용되는 중요 기자재(계측기, 밸브, 펌프 등)가 국제 및 국내 규격에 부합하는지 확인하기 위해 각종 물성, 내식성, 내열성, 강도, 누설, 방폭 등 여러 항목의 시험성적서 제출이 의무화되어 있다. 하지만 수많은 시험 성적서를 엔지니어가 육안으로 일일이 관리하는 것은 불가능하다. 수동으로 문서를 디지털화하는 것은 비용이 많이 들고 오류가 발생하기 쉬운 문서 관리의 취약점이 발생한다[1]. 2013년 한국수력원자력(한수원)의 원자력발전소에 납품된 부품들의 시험성적서가 위조된 사건이 발생하였다. 기술 규격에 미달하는 케이블을 납품하기 위해 시험성적서를 위조한 이 사건은 원전 부품의 품질과 안전성에 대한 신뢰를 저하 시키는 계기가 되었으며, 이후 산업계 전반에 걸쳐 품질 관리와 검증 절차의 중요성이 부각되었다. 시험성적서의 관리를 위해서 여러 기술을 적용하고 있는데 그 중 하나가 바로 OCR (Optical Character Recognition, 광학 문자 인식)이다. OCR은 이미지가 스캔된 문서 속의 글자를 컴퓨터가 읽고 텍스트 데이터로 변환하는 기술이다. OCR은 자동 문서 관리 분야에서 매우 널리 활용되고 있다. 대부분의 사무자동화, 공장관리, 온라인 교육, 정부 행정, 산업계 등 다양한 분야의 자동화 시스템이 스캔 문서, 이미지, PDF 등에서 정보 추출을 위해 OCR을 기본 도구로 채택하고 있다[2]. 시험성적서 인식에 OCR 기술이 폭넓게 쓰이지만, 다양한 시험성적서의 서식, 비정형화된 레이아웃, 수기 및 복잡한 서체 등에서 인식을 저하, 오인식, 후처리 비용 증가 등 여러 문제가 동반되고 있는 것도 사실이다[3].

II. LLM 기반 시험성적서 시스템 설계

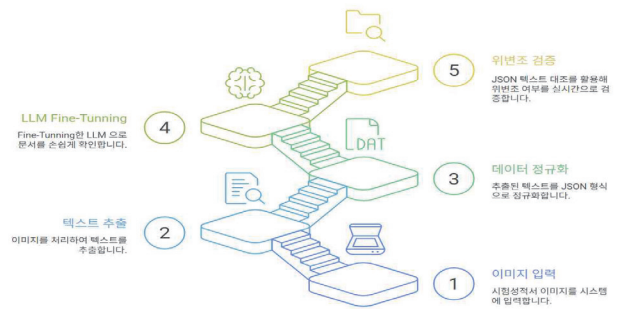


그림 1. 시험성적서 자동화 시스템 검증절차

OCR 시스템은 정형화된 문서(예: 명확한 표, 단순한 텍스트, 정확한 레이아웃)에는 높은 인식률을 보이나, 표·수식·차트 등 구조가 복잡하거나 다양한 형태가 혼재된 비정형 문서에서는 정확도가 크게 떨어진다. 실험 및 문헌분석 결과 비정형성에 따라 오인식률이 2에서 10배 이상 높아지며, 특히 숫자가 많은 문서에서 후처리 난이도가 증가함을 알 수 있다. 이와같은 문제점들을 해결하기 위해 본 논문은 Transformer 기반 Donut 모델로 시험성적서를 자동으로 JSON 형태로 정규화하고, 이를 Fine-Tuning 한 LLM(Large Language Model)을 통해 손쉽게 확인할 수 있는 문서 자동화 시스템 구조는 그림 1. 과 같이 구성된다. 신규로 등록된 동일 종류의 시험성적서도 JSON 텍스트 대조를 활용해 위변조 여부를 실시간으로 검증하는 기능을 추가해, 데이터의 진위와 무결성을 높이는 방안을 제시한다.

III. Donut 기반 문서 자동화 시스템

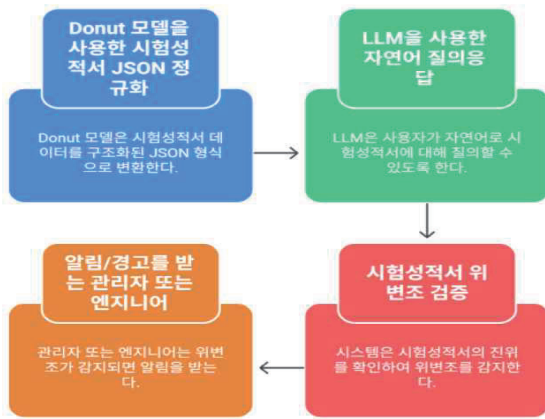


그림 2. 시험성적서 자동화 프레임워크

본 논문에서는 기존 OCR(광학 문자 인식)의 한계로 인해 OCR을 배제한 Transformer 기반 VDU(Vision Document Understanding) 구조로, 문서 이미지를 직접 입력받아 텍스트/구조 정보를 종단(end-to-end) 방식으로 추출하는 OCR-free Donut (Document Understanding Transformer) 모델을 채택한다[4]. Donut은 문서 내 시각·텍스트 정보를 모두 직접 처리하므로, 기존 OCR처럼 전처리에 의존하거나 텍스트 추출-이해 단계에서 오류 전파가 발생하지 않는다. 다양한 언어·서체·복잡한 구조(표, 차트, 도식 등)에 대해 별도 엔진 교체 없이 종단적으로 처리 가능하며, 실험적으로 비정형/저자원 문서에서도 뛰어난 성능을 달성한다. Donut 모델로 시험성적서를 JSON 형태로 정규화하는 과정을 순차적으로 정리하면 다음과 같다.

- 1) Donut 모델이 이미지 이해 : Transformer 기반 Donut 모델이 문서 내에서 필요한 정보(예: 방폭등급, 시험결과, 규격, 발급일 등)를 자동으로 찾아낸다.
- 2) 핵심정보 추출 : Donut 모델이 각 항목별 값을 인식 (예시: “시험결과 : 합격”, “발급일 : 2025-09-01” 등)을 데이터로 정리한다.
- 3) JSON 데이터로 변환 : 추출된 정보들을 표준화된 JSON 형태로 자동으로 정규화 한다.
- 4) 데이터 저장 및 활용 : 생성된 JSON 데이터는 데이터베이스에 저장되고, 이후 검색, 검증, 요약 등 다양한 방식으로 활용될 수 있다. 이 과정을 통해 복잡한 시험성적서를 쉽고 정확하게 표준 데이터로 변환하여 관리할 수 있다.

정규화된 JSON 데이터는 LLM 입력으로 전달된다. LLM은 사전 학습된 대형 언어 모델이지만, 문서 도메인과 JSON 구조 이해를 돕기 위해 Fine-Tuning을 진행한다.

본 논문에서 제안하는 시험성적서 자동화 프레임워크는 그림 2.와 같이 구성된다. Fine-Tuning 과정에서는 시험성적서 예시 데이터, JSON 구조, 자연어 질의응답 쌍 등이 사용된다. 사용자는 자연어로 질의(예: “이 앨브의 시험 결과를 요약해줘”)를 입력한다. LLM은 정규화된 JSON 데이터에서 관련 정보를 파악해, 자연스럽고 정확한 언어로 답변을 생성한다. 필요 시 JSON 내 특정 키-값 쌍을 조회하거나 조건 필터링 결과를 포함한 응답을 제공한다. 엔지니어는 챗봇 형태 인터페이스에서 LLM과 대화하며 문서 내용을 쉽게 확인·검색할 수 있다. Donut 모델의 빠른 JSON 변환과 LLM의 유연한 질의응답이 연속적으로 연동되어 효율적인 업무 지원이 이루어진다. Donut 모델은 비정형 문서를 구조화된 JSON으로 바꾸는 역할을, LLM은 JSON 데이터를 바탕으로 자연어 이해와 응답을 담

당하여, 두 모델이 협력해 시험성적서 정보를 쉽고 빠르게 처리·제공할 수 있게 하는 구조이다[5]. 신규로 등록된 동일 종류의 시험성적서도 Donut 모델로 JSON 형태로 변환된다. 새로 생성된 JSON과 기존 저장된 JSON 데이터를 텍스트 기반으로 대조한다. 각 키-값 쌍별로 일치 여부를 검사하며, 값의 변경, 누락, 추가가 있는지 깊이 비교한다. JSON 구조의 트리나 키 순서, 데이터 타입까지 포함해 정밀 검증이 가능하며 이를 통해 위변조 여부를 판별한다. 변조가 감지되면 해당 기록을 저장하고, 관리자 또는 엔지니어에게 알림 또는 경고 메시지를 발송한다. 위변조 여부는 문서의 진위 확인, 감사 추적, 법적 증빙 등 다양한 목적에 활용될 수 있다. 자동화된 위변조 검증으로 문서 관리 투명성 및 신뢰성이 크게 향상될 것을 기대한다.

III. 결론

본 논문은 플랜트 산업에서 필수로 요구되는 시험성적서의 방대하고 복잡한 관리를 자동화하기 위해 Transformer 기반 Donut 모델과 Fine-Tuning 된 LLM을 결합한 혁신적 문서 자동화 시스템을 제시하였다. Donut 모델은 시험성적서 이미지를 직접 분석하여 주요 항목을 정밀하게 추출하고, 이를 표준화된 JSON 형태로 자동 정규화함으로써 기존 OCR 기술의 한계였던 비정형 문서 처리 어려움과 인식 오류 문제를 효과적으로 극복하였다. 또한, Fine-Tuning된 LLM을 활용해 정규화된 JSON 데이터를 자연어 질의와 대화형 검색 시스템과 연동함으로써 손쉽게 시험성적서를 확인하고 요약·검증할 수 있도록 지원하였다. 이러한 접근은 문서 탐색 및 정보 활용의 생산성 및 업무 효율을 증대시킨다.

아울러, 신규 등록되는 동일 문서에 대해 JSON 텍스트 대조를 통해 위변조 여부를 실시간으로 검증하는 기능을 추가하여, 문서의 진위와 무결성을 보장하고 산업 현장의 안전성과 신뢰성을 강화하여 산업 현장의 디지털 전환과 스마트 관리를 촉진하는 데 핵심적인 기여를 할 것으로 판단된다.

참 고 문 헌

- [1] J.-A. Labarga, A. Singh, V.-Z. Moffitt, “An Extensible System for Optical Character Recognition of Maintenance Documents”, Proc. A nnu. Conf. Prognostics and Health Management Society, Philadelphia, PA, USA, pp. 1-8. Sep, 2018.
- [2] E. Borovikov “A survey of modern optical character recognition techniques”, arXiv preprint arXiv:1412.4183, [Online]. Available: <https://arxiv.org/abs/1412.4183>, Dec. 2014.
- [3] 민기현, 이아람, 김거식, 김정은, 강현서, 이길행, “딥러닝 기반 광학 문자 인식 기술 동향 (Recent Trends in Deep Learning-Based Optical Character Recognition)”, 전자통신동향분석 제37권 제5호, pp. 22-32, Oct, 2022.
- [4] G. Kim, T. Hong, M. Yim, J. Nam, J. Park, J. Yim, W. Hwang, S. Yun, D. Han, S. Park, “OCR-free Document Understanding Transformer”, in Proc. 17th Eur. Conf. Comput. Vis. (ECCV), Tel Aviv, Israel, , Lecture Notes in Computer Science, vol. 13688, pp. 498-517, Oct. 2022.
- [5] S. Yao, D. Yu, J. Zhao, I. Shafran, T.-L. Griffiths, Y. Cao, K.-R. Narasimhan, “Tree of Thoughts: Deliberate Problem Solving with Large Language Models”, in Advances in Neural Information Processing Systems”, vol. 36, Proc. 37th Conf. Neural Information Processing Systems (NeurIPS), New Orleans, LA, USA, Dec. 2023.
- [6] 블록체인 기반 문서 플랫폼, purecertificate.com