

# 자동화 팩트체킹 연구 동향: 벤치마크 데이터셋 중심으로

숭실대학교 | 박건우\*

## 1. 서론

정보 통신 기술의 발달로 온라인 공간에서 많은 사람과 연결되어 정보를 빠르고 손쉽게 접할 수 있게 되었다. 이에 따라, 재난 상황 내 위급 정보 전달 등 사회적으로 의미 있는 응용들이 연구 및 개발되었다. 반면, 정보 공유와 전파의 용이성으로 인해, 가짜뉴스 문제가 대두되었다. TV, 신문 등 전통적인 미디어와 다르게, 소셜 미디어, 온라인 포럼 등은 정보 공유 시 진위 검증 절차가 부재하거나 간단하다. 거짓 정보는 자극적인 내용을 담는 경우가 많고 전파 속도도 진실 정보에 비해 빠르므로[1], 거짓 정보가 퍼진 후 이를 바로잡는 것은 어렵다. 더욱이, 최근 대규모 언어 모델(large language models), 확산 모델(diffusion models) 등 생성형 인공지능 기술이 급속도로 발전하고 ChatGPT 등 이를 다루기 위한 인터페이스가 간단해지면서 전문 지식 없이도 거짓 정보를 쉽게 생성할 수 있게 되었다. 온라인 공간 내 공유된 거짓 정보를 탐지할 수 있는 인공지능 기술의 개발이 사회적으로 요구되고 있다.

팩트체킹은 주장의 사실 및 거짓 여부를 판단하는 과업으로 주로 저널리즘에서 전문가 팩트체커에 의해 다루어져 왔다. 보도의 신뢰성을 높이기 위해 언론사는 내부적으로 보도 전 진위성 검증을 수행하며, PolitiFact 등 전문 팩트체크 기관은 정치인의 발언, 소셜 미디어에 공유된 정보 중 검증할 만한 주장에 대해 팩트체크 하여 관련 정보를 공개한다. 수작업 팩트체킹은 많은 시간이 소요되는 고비용 저효율 작업이다. 주장에 대해 관련된 근거를 찾고, 이를 종합하여 거짓 여부를 판단하고, 충분한 설명을 제공하는 작업은 검증에 며칠이 걸리기도 한다[2]. 온라인 채널을



그림 1. 소셜 미디어에 공유된 거짓 정보의 예시 (출처: 페이스북)

통해 공유되고 전파되는 수많은 정보를 검증하기에 수작업 팩트체킹은 불충분하다.

전문가 팩트체킹의 한계를 극복하고 확장 가능한 사실 검증을 위해 자동화 팩트체킹 연구가 수행되었다[3]. 검증할 만한 주장을 찾고, 근거를 검색하고, 확보한 근거를 바탕으로 진위성을 예측하는 팩트체킹 절차는 자연어처리 및 인공지능 응용 기술로써 다뤄질 수 있다. 과거 연구들이 과업을 설정하고 이에 대한 데이터셋을 제공함으로써 관련 기술 개발을 촉진하였다.

이 고는 자동화 팩트체킹을 위해 제안된 주요 벤치마크 데이터셋을 소개한다. 팩트체킹 기술을 개발하고 평가하기 위한 주요 데이터셋의 구축 과정을 알아봄으로써, 사회-기술적 관점에서 자동화 팩트체킹 과업에 대해 이해하고 발전 방안을 모색해보고자 한다.

## 2. 자동화 팩트체킹

전문가 팩트체킹을 모사한 자동화 팩트체킹은 그림 2와 같이 세 가지 세부 과업으로 제안[3]되었다: 주장 탐지(claim detection), 근거 검색(evidence retrieval), 평

\* 정회원

† 이 논문은 정부(과학기술정보통신부)의 재원으로 정보통신기획평가원의 지원을 받아 지역지능화혁신인재양성사업(IITP-2025-RS-2022-00156360)과 메타버스 융합대학원(IITP-2025-RS-2024-00430997)의 연구 결과로 수행되었음

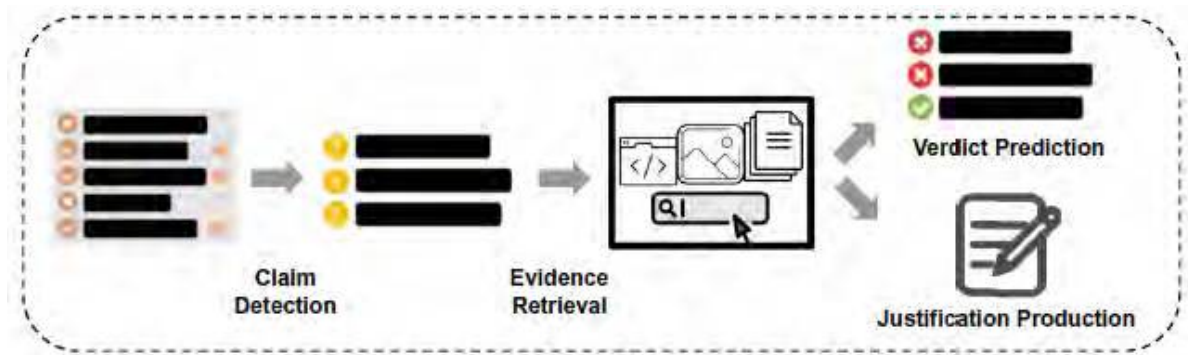


그림 2 자동화 팩트체크 흐름 및 세부 과업 (그림 출처: [3])

결 예측(verdict prediction). 주장 탐지 단계는 소셜 미디어, 뉴스 기사 등 온라인 텍스트로부터 검증이 필요한 텍스트 형태의 발언, 즉, 주장(claim)을 탐지하는 것을 목표로 한다. 대표적인 주장 탐지 연구로 사실이 확인되지 않은 정보의 탐지를 목표로 하는 소문 탐지 (rumor detection) 연구가 있다. 검증할 만한 가치가 있는 정보를 필터링하는 단계로 이해할 수 있으며, 특정 정보에 대한 검증을 목표로 할 경우 이 단계는 생략될 수 있다. 근거 검색은 주어진 주장의 진위성을 판단할 수 있는 추가 정보, 즉 근거를 확보하는 것을 목표로 한다. 평결 예측은 관련된 근거를 바탕으로 주장의 진위성을 예측하는 단계이다. 평결 예측과 함께 예측에 대한 타당한 이유를 생성하는 과정이 추가되기도 한다. 3장에서 소개할 데이터셋은 근거 검색 및 평결 예측 과업을 다루며, 주장 탐지에 관해 관심 있는 독자는 서베이 논문[3]을 참고하기 바란다.

### 3. 주요 벤치마크 데이터셋 소개

#### 3.1 LIAR [5]

자연어처리, 컴퓨터비전 등 인공지능 응용 분야에서 대규모 라벨링 데이터셋을 구축할 때 크라우드소싱이 사용되어 왔다[6]. 크라우드소싱은 상대적으로 저렴한 가격에 많은 수의 작업자를 고용할 수 있어 규모 있는 데이터셋 확보에 용이하지만, 진위성 여부를 판단하는 것은 저널리즘 전문 지식이 필요한 어려운 작업이기 때문에 일반 작업자의 라벨은 품질 문제가 발생할 수 있다. 따라서, 많은 팩트체크 데이터셋은 신뢰할 수 있는 팩트체크 웹사이트에서 공개한 전문가 라벨 정보를 활용하였다. LIAR[5]는 이런 방식으로 구축된 대표적인 데이터셋이다. Wang은 미국의 주요 팩트체크 웹사이트 중 하나인 PolitiFact로부터 12,800 건의 주장 및 검증 정보를 수집하였다.

**Statement:** "The last quarter, it was just announced, our gross domestic product was below zero. Who ever heard of this? Its never below zero."  
**Speaker:** Donald Trump  
**Context:** presidential announcement speech  
**Label:** Pants on Fire  
**Justification:** According to Bureau of Economic Analysis and National Bureau of Economic Research, the growth in the gross domestic product has been below zero 42 times over 68 years. Thats a lot more than "never." We rate his claim Pants on Fire!

그림 3 LIAR 데이터셋 샘플 (그림 출처: [5])

그림 2와 같이, 웹사이트에서 제공하는 평결 기준에 따라 진위성 라벨은 6단계로 나뉜다: *pants-fire*, *false*, *barely-true*, *half-true*, *mostly-true*, *true*.

검증 대상 주장은 보도자료, TV/라디오 인터뷰, 선거 유세 발언, 광고, 소셜 미디어 포스트 등 다양한 맥락을 포괄한다. 이 데이터셋은 별도의 지식 저장소를 제공하지 않아, 주로 평결 예측을 위한 분류기 구축에 활용되었다. 이후 연구에서 설명문을 machine-readable 한 형태로 처리하여 확장한 LIAR-PLUS[7] 데이터셋이 제안되었고, 설명문을 추가 입력으로 고려하였을 때 검증 성능을 향상시킬 수 있음을 보였다.

#### 3.2 FEVER [8]

LIAR 형태의 데이터셋은 지식 저장소를 활용하지 않아 이에 학습된 팩트체크 기술의 활용성이 제한된다. 파라미터화된 지식에만 기반해 진위 여부를 예측할 경우, 알지 못하는 지식과 다양한 도메인에 대한 주장을 다루는 팩트체크를 효과적으로 수행할 수 없다. 따라서, Throne 등[8]은 텍스트 형태의 지식 저장소에

**Claim:** The Rodney King riots took place in the most populous county in the USA.

**[wiki/Los Angeles Riots]**

The 1992 Los Angeles riots, also known as the Rodney King riots were a series of riots, lootings, arsons, and civil disturbances that occurred in Los Angeles County, California in April and May 1992.

**[wiki/Los Angeles County]**

Los Angeles County, officially the County of Los Angeles, is the most populous county in the USA.

**Verdict:** Supported

그림 4 FEVER 데이터셋 샘플 (그림 출처: [8])

기반해 주장의 진위성 여부를 검증하는 근거 기반의 팩트체크 과업과 데이터셋 FEVER(Fact Extraction and Verification)를 제안하였다. 주어진 주장에 대해, 위키피디아 문서를 지식 저장소로 활용하여 관련된 근거를 검색하며, 검색된 근거에 기반해 주장의 진위성 여부를 다음의 세 가지 라벨 중 하나로 예측하는 것을 목표로 한다: *supported*, *refuted*, *not enough information*. 학습 및 평가를 위한 주장과 라벨 정보는 다음 두 단계에 걸쳐 구축되었다.

(1) 주장 생성: 50,000 개 정도의 인기 있는 위키피디아 문서의 소개 섹션에서 문장을 무작위로 추출하여, 사람 작업자에게 제공한다. 작업자는 주어진 문장으로부터 해당 위키피디아 문서가 다루는 엔티티에 관한 주장을 직접 생성한다. 검증이 쉽지 않은 주장을 구축하기 위해 연결된 문서의 키워드 등을 활용하도록 한다. 이어서, 다음의 여섯 가지 방법을 적용해 주장을 복잡하게 만든다: 패러프레이징, 부정, 엔티티 및 엔티티 간 관계를 유사한 것과 다른 것으로 대체, 주장 구체화

(2) 주장 라벨링: 사람 작업자는 1단계에서 생성된 각 주장에 대해, 위키피디아 문서를 직접 참조하여 진위성 라벨링을 수행한다.

위 절차에 따라 구축된 FEVER 데이터셋은 그림 3의 예시와 같은 185,445개의 주장과 근거 쌍으로 구성된다. 주장별 관련 문서와 근거 문장이 라벨링된 형태로 볼 수 있으며, 검색기 학습을 위한 라벨로 사용 가능하다. FEVER의 주장 및 문서 쌍은 정보 검색 벤치마크

BEIR[9]에 포함되어 있다. 또한, 근거 기반의 팩트체크는 대표적인 지식-집약적(knowledge-intensive) 과업으로 간주되며, FEVER는 관련 벤치마크 KILT[10]에 포함되어 검색 증강 생성 기술 학습 및 평가에 활용되고 있다. FEVER와 유사한 방식으로 모조 주장을 다루되 멀티-홉 근거 검색이 필요한 HOVER[11], 테이블 형태의 근거를 포함하는 FEVEROUS[12] 등이 이어 제안되기도 하였다.

### 3.3 MultiFC [13]

FEVER는 주장을 검증할 수 있는 라벨된 문서와 근거 문장 정보를 위키피디아 지식 저장소와 함께 제공하여, 규모 있는 근거 기반의 팩트체크 기술 연구가 가능하게 한 첫 시도로서 의의가 있다. 하지만, 모조 주장을 생성하는 방식으로 데이터셋이 구축됨으로써 FEVER 데이터셋이 다루는 주장은 검증이 필요한 실세상의 주장과 도메인 및 특성 등이 다를 수 있다. 또한, LIAR는 하나의 팩트체크 웹사이트에서 다루는 주장에 집중하여 주장의 도메인이 제한될 수 있다. MultiFC는 위 한계점들을 해결하기 위해 Augenstein 등[13]이 제안한 데이터셋으로, 27개의 전문가 팩트체크 데이터 웹사이트에 공유된 주장 36,534 건으로 구성되어 다양한 도메인의 실제 주장을 포괄한다. 각 팩트체크 웹사이트 별 평결 클래스 개수가 달라 데이터셋 내 총 라벨 수가 165개에 달하나, 연구팀은 라벨을 통합하는 대신 그대로 두고 라벨 임베딩 기반의 다중-과업 학습(multi-task learning) 방식으로 모델링 하는 방법을 택했다. 그림 5는 MultiFC 데이터셋 샘플을

Feature	Value
ClaimID	farg-00004
Claim	Mexico and Canada assemble cars with foreign parts and send them to the U.S. with no tax.
Label	distorts
Claim URL	<a href="https://www.factcheck.org/2018/10/factchecking-trump-on-trade/">https://www.factcheck.org/2018/10/factchecking-trump-on-trade/</a>
Reason	None
Category	the-factcheck-wire
Speaker	Donald Trump
Checker	Eugene Kiely
Tags	North American Free Trade Agreement
Claim Entities	United_States, Canada, Mexico
Article Title	FactChecking Trump on Trade
Publish Date	October 3, 2018
Claim Date	Monday, October 1, 2018

그림 5 MultiFC 데이터셋 샘플 (그림 출처: [13])



**Claim:** *The USA has succeeded in reducing greenhouse emissions in previous years.*  
**Date:** 2020.11.2 **Speaker:** Morgan Griffith **Type:** ...

**Q1:** What were the total gross U.S. greenhouse gas emissions in 2007?  
**A1:** In 2007, total gross U.S. greenhouse gas emissions were 7,371 MMT.  
**Q2:** When did greenhouse gas emissions drop in US?  
**A2:** In 2017, total gross U.S. greenhouse gas emissions were 6,472.3 MMT, or million metric tons, carbon dioxide.  
**Q3:** Did the total gross U.S. greenhouse gas emissions rise after 2017?  
**A3:** Yes. After 3 years of decline, US CO2 emissions rose sharply last year. Based on preliminary power generation, natural gas, and oil consumption data, we estimate emissions increased by 3.4% in 2018.

**Verdict:** **Conflicting Evidence/Cherrypicking.**  
**Justification:** *It is true they did reduce emissions however they have now increased again. It is unknown exactly what years are being referred to.*

그림 6 AVeriTeC 데이터셋 샘플 (그림 출처: [14])

나타낸다. 지식 저장소는 제공하지 않으며, 근거 기반의 팩트체킹을 위해 구글 검색 API를 활용하였다.

### 3.4 AVeriTeC [14]

Ousidhoum 등[15]과 Glokner 등[16]은 기존의 자동화 팩트체킹 데이터셋이 크게 세 가지 관점에서 제한됨을 보고하였다. 첫째, 일부 주장은 맥락화되어 이해하기 어렵다. 예를 들어, “unemployment is rising” 이라는 주장은 지역과 시간에 대한 맥락이 추가되어야 검증이 가능하다. 둘째, 주장 라벨과 근거가 대응되지 않는 경우가 있다. 주장을 검증하기에 주석 처리된 근거가 충분하지 않거나, 주석 처리된 근거와 상충되는 외부 정보가 있는 경우도 있다. 셋째, 미래 정보가 근거로 포함된 경우가 있다. 예를 들어, 2024년 10월의 주장을 검증하는 데 12월의 정보가 근거로 포함되는 비현실적인 상황이다.

위와 같은 세 가지 한계점을 해결하기 위해, Schlichtkrull 등[14]은 AVeriTeC (Automated VERification of TExtual Claims) 데이터셋을 제안하였다. 이 데이터셋은 4,568개의 주장으로 구성된다. 한 주장 당 다섯명의 작업자를 할당하고 근거를 획득하기 위한 절차를 세분화함으로써 기존의 데이터셋이 지니는 문제를 해결하고 고품질의 데이터를 확보할 수 있도록 했다.

연구팀은 LIAR, MultiFC 등을 구축한 방식과 유사하게 50개 팩트체크 웹사이트로부터 데이터셋을 수집하고, 수집된 각 주장에 대해 크라우드소싱을 이용해 근거를 확보하는 주석 작업을 수행하였다. 주어진 주장에 대해 근거를 찾고 평결을 내리는 추론 과정은, 주장에 대한 검증 질문을 생성하고 이에 대한 답을 찾는 과정을 반복하여 수행할 수 있다. 저자들은 수작업자들에게 검증 질문을 생성하도록 하고 웹 검색 엔

진을 사용해 근거를 확보하도록 하여, 주어진 주장의 검증에 필요한 질문-근거 쌍을 확보하고 평결을 내릴 수 있도록 했다. 검색 시 대상 문서를 주장이 발행된 이전 날짜로 제한하여 미래의 근거를 사용하는 것을 방지하였다. 수작업자가 내린 평결 라벨이 전문가 라벨과 일치하면 과업 종료, 불일치하면 이전 과정을 반복하는 방식으로 라벨 품질을 높이고 충분한 근거를 확보할 수 있도록 했다. 평결 라벨 종류는 네 가지로, FEVER에서 사용한 세 가지 분류에 더해 *conflicting evidence/cherrypicking* 분류를 추가하였다. 마지막으로, 어떻게 평결에 이르게 되었는지 검증 과정에 대한 타당한 이유를 수작업자들에게 작성하게 하여 데이터셋에 포함하였다.

그림 6은 AVeriTeC 데이터셋 샘플을 나타낸다. 주장에 대한 메타데이터, 평결 라벨 및 설명문과 함께, 근거에 대한 질문-답변 쌍이 제공된다. 이에 더해, AVeriTeC은 작업자들이 근거 검색 시 참조한 웹 문서를 지식 저장소로 함께 제공한다. 연구팀은 구축된 데이터셋을 바탕으로, 주장에 대해 질문을 생성하고, 근거를 확보하고, 평결 예측 후 설명문을 생성하는 세부 과업을 제안하였다. 사람의 팩트체킹 과정과 유사한 형태로 과업을 설계한 것이다. 또한, 주석 처리된 질문과 생성 질문을 비교하여 유사한 경우에만 평결 예측 정확도를 측정하는 평가 지표 AVeriTeC score를 함께 제안하였다. 최근, EMNLP 2024에서 열린 팩트체킹 분야 저명 워크숍 FEVER에서 AVeriTeC 데이터셋에 기반한 공유 과업 경진대회[17]가 개최되었다. 각 팀은 gpt, llama 등 다양한 대규모 언어 모델에 기반한 시스템을 제안하였다. 1위 팀 TUDA\_MAI 는 API를 사용해 gpt-4o로 각 절차를 수행하는 InFact를 개발하여 최고 성능(0.63)을 달성하였으며, 필자가 소속된

HUMANE 팀은 llama 3.1 기반 HerO를 개발하여 베이스라인 모델 대비 5배 이상의 성능 향상(0.57)을 기록했다.

#### 4. 결론

이 고는 자연어처리 분야에서 자동화 팩트체킹을 위해 제안된 핵심 데이터셋 네 가지를 살펴보았다. 벤치마크의 역할을 하고 있는 주요 데이터셋을 살펴봄으로써 사회-기술적 관점에서 자동화 팩트체킹 기술의 현주소를 가늠하고자 하였다.

팩트체킹 데이터셋 구축의 어려움은 진위성 평결의 어려움에서 기인한다. 일반 작업자가 라벨링을 수행할 수 있는 간단한 과업과 달리, 진위성 평결은 전문 지식을 지닌 팩트체커가 수행하지 않을 경우 라벨 품질을 보장할 수 없기 때문에 LIAR, MultiFC처럼 전문가 팩트체커 웹사이트에 공개된 정보를 수집하여 데이터셋을 구성하는 방법이 주를 이루었다. 이 방법은 사실 여부에 대한 고품질의 전문가 라벨을 제공하지만, 전문가 팩트체커가 사실을 검증하는 과정을 자동화하기 위해 관련 근거를 검색해 주장의 사실 여부를 검증할 수 있도록 근거를 포함하는 지식 저장소를 함께 제공하는 것이 바람직하다. 따라서, 위키피디아를 지식 저장소로 활용하는 FEVER 데이터셋이 제안되었다. 하지만, 진위성 라벨 할당을 위해 모조 주장을 생성하는 방식을 채택하여 실제 주장의 도메인 및 형태와 차이가 발생하게 되었다. 마지막으로 소개한 AVeriTeC은 전문가 평결 예측 라벨을 사용하되 검증 질문 기반으로 웹 문서를 검색해 근거를 확보하는 창의적인 방법을 제안함으로써 두 방법의 장점을 결합할 수 있게 되었다. AVeriTeC은 웹 스케일 근거 검색에 기반한 팩트체킹 데이터셋이라는 점에 그 의의가 있으나, 팩트체킹 과정에 대한 전문 지식이 없는 크라우드소싱 작업자가 생성한 질문과 이에 대한 근거 데이터셋을 구성하여 전문가 팩트체킹 과정에 대한 대표성이 떨어질 수 있다는 한계가 있다.

자동화 팩트체킹 기술이 전문가 팩트체킹 과정을 충실히 반영하는 방향으로 발전되기 위해, 미래에 요구되는 데이터셋 및 연구 방향은 크게 세 가지가 있다.

(1) 다양한 종류의 근거를 포괄하는 데이터셋: 실제 세상에서 검증이 필요한 주장과 근거는 텍스트 뿐 아니라 이미지, 동영상 등 다양한 형태로 나타난다. 멀티모달 정보를 포괄하는 데이터셋 구축(예. [18]) 및 검증 기술 연구가 필요하다.

(2) 한국 맥락을 반영하는 데이터셋: 국가와 문화에 따라 검증이 필요한 주장의 도메인과 형태가 다를 수 있다. 국내 인물과 사건에 관한 팩트체킹 데이터셋 구축 및 검증 기술 연구가 필요하다.

(3) 전문가 지식 기반 데이터셋: 전문가 팩트체커의 지식과 원칙 등을 반영한 데이터셋을 구축할 경우 보다 효과적인 검증이 가능할 것을 기대할 수 있다.

위와 같은 데이터셋 및 인공지능 기술 연구가 활성화되어, 실제 세상의 주장 검증을 위해 활용 가능한 자동화 팩트체킹 기술이 개발될 수 있기를 바란다.

#### 참고문헌

- [ 1 ] S. Vosoughi, D. Roy, and S. Aral, "The spread of true and false news online," *science*, vol. 359, no. 6380, pp. 1146-1151, 2018.
- [ 2 ] A. Zhao and M. Naaman, "Insights from a Comparative Study on the Variety, Velocity, Veracity, and Viability of Crowdsourced and Professional Fact-Checking Services," *Journal of Online Trust and Safety*, vol. 2, no. 1, 2023.
- [ 3 ] Z. Guo, M. Schlichtkrull, and A. Vlachos, "A survey on automated fact-checking," *Transactions of the Association for Computational Linguistics*, vol. 10, pp. 178-206, 2022.
- [ 4 ] A. Vlachos and S. Riedel, "Fact checking: Task definition and dataset construction," *ACL 2014 workshop on language technologies and computational social science*, pp. 18-22, 2014.
- [ 5 ] W. Y. Wang, "'Liar, Liar Pants on Fire': A New Benchmark Dataset for Fake News Detection," *55th Annual Meeting of the Association for Computational Linguistics, Volume 2: Short Papers*, pp. 422-426, 2017.
- [ 6 ] J. W. Vaughan, "Making better use of the crowd: How crowdsourcing can advance machine learning research," *Journal of Machine Learning Research*, vol. 18, no. 193, pp. 1-46, 2018.
- [ 7 ] T. Alhindi, S. Petridis, and S. Muresan, "Where is your evidence: Improving fact-checking by justification modeling," *the first workshop on fact extraction and verification(FEVER)*, pp. 85-90, 2018.
- [ 8 ] J. Thorne, A. Vlachos, C. Christodoulopoulos, and A. Mittal, "FEVER: a Large-scale Dataset for Fact Extraction and VERification," *2018 Conference of the North American Chapter of the Association for*

---

Computational Linguistics: Human Language Technologies, Volume 1(Long Papers), pp. 809-819, 2018.

- [ 9 ] N. Thakur, N. Reimers, A. Rücklé, A. Srivastava, and I. Gurevych, “BEIR: A Heterogeneous Benchmark for Zero-shot Evaluation of Information Retrieval Models,” Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track(Round 2), 2021.
- [10] F. Petroni et al., “KILT: a Benchmark for Knowledge Intensive Language Tasks,” 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pp. 2523-2544, 2021.
- [11] Y. Jiang, S. Bordia, Z. Zhong, C. Dognin, M. Singh, and M. Bansal, “HoVer: A Dataset for Many-Hop Fact Extraction And Claim Verification,” Findings of the Association for Computational Linguistics: EMNLP 2020, pp. 3441-3460, 2020.
- [12] R. Aly et al., “FEVEROUS: Fact Extraction and VERification Over Unstructured and Structured information,” Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track(Round 1), 2021.
- [13] I. Augenstein et al., “MultiFC: A Real-World Multi-Domain Dataset for Evidence-Based Fact Checking of Claims,” 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP), pp. 4685-4697, 2019.
- [14] M. Schlichtkrull, Z. Guo, and A. Vlachos, “Averitec: A dataset for real-world claim verification with evidence

from the web,” Advances in Neural Information Processing Systems, vol. 36, 2024.

- [15] N. Ousidhoum, Z. Yuan, and A. Vlachos, “Varifocal Question Generation for Fact-checking,” 2022 Conference on Empirical Methods in Natural Language Processing, pp. 2532-2544, 2022.
- [16] M. Glockner, Y. Hou, and I. Gurevych, “Missing Counter-Evidence Renders NLP Fact-Checking Unrealistic for Misinformation,” 2022 Conference on Empirical Methods in Natural Language Processing, pp. 5916-5936, 2022.
- [17] M. Schlichtkrull et al., “The Automated Verification of Textual Claims(AVeriTeC) Shared Task,” Seventh Fact Extraction and VERification Workshop(FEVER), pp. 1-26, 2024.
- [18] B. M. Yao, A. Shah, L. Sun, J.-H. Cho, and L. Huang, “End-to-end multimodal fact-checking and explanation generation: A challenging dataset and models,” 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, pp. 2733-2743, 2023.

## 약 력



### 박 건 우

2021~현재 숭실대학교 AI융합학부 교수  
2020~2021 미국 University of California Los Angeles,  
Postdoctoral researcher  
2018~2020 카타르 Qatar Computing Research  
Institute, Postdoctoral researcher  
2018 KAIST 전산학부 박사

관심분야: 자연어처리, 멀티모달, 데이터과학  
Email : kunwoo.park@ssu.ac.kr