

## 증례보고의 임상 경로 탐색을 위한 전이행렬 기반 시각적 분석 시스템 설계

전혜리<sup>1</sup> · 권찬영<sup>2</sup> · 김성희<sup>3\*</sup>

<sup>1</sup>동의대학교 인공지능학과 학사과정

<sup>2</sup>동의대학교 한의과대학 한방신경정신과 조교수

<sup>3</sup>동의대학교 산업ICT기술공학과 부교수

## Design of a Visual Analytics System Utilizing a Transition Matrix to Explore Clinical Pathways in Case Reports

Hye-Li Jeon<sup>1</sup> · Chan-Young Kwon<sup>2</sup> · Sung-Hee Kim<sup>3\*</sup>

<sup>1</sup>Undergraduate Student, Department of Artificial Intelligence, Dong-eui University, Busan 47340, Korea

<sup>2</sup>Assistant Professor, Department of Oriental Neuropsychiatry, College of Korean Medicine, Dong-eui University, Busan 47227, Korea

<sup>3</sup>Associate Professor, Department of ICT Industrial Engineering, Dong-eui University, Busan 47340, Korea

### [요약]

증례보고는 표준화된 작성 지침에도 불구하고 서술의 다양성으로 인해 체계적인 비교와 대규모 분석에 한계가 있다. 본 연구는 증례보고의 임상 경로를 구조화하여 탐색하는 시각적 분석 시스템을 제안한다. CARE guidelines를 기반으로 증상, 검사, 진단, 치료, 결과로 축을 정의하였으며, 2010년 1월부터 2025년 8월까지 게재된 1,320편의 증례보고를 대상으로 대형 언어 모델(LLM; Large Language Model)을 활용해 사건 시퀀스를 추출·정규화하고, 전이행렬 기반의 임상 경로 데이터를 구성하였다. 제안 시스템은 막대그래프, 산점도, Sankey 다이어그램을 통합하고 임상 경로와 개별 증례 원문을 연계 탐색할 수 있도록 지원한다. 그 결과, 임상 경로 데이터의 대부분은 높은 이질성을 보였으나, 일부 군집에서는 공통된 임상 흐름이 시각적으로 명확하게 식별되었다. 본 연구는 증례보고 기반 임상 경로 분석에서 시각적 탐색을 중심으로 한 분석 환경 가능성을 제시한다.

### [Abstract]

Despite guidelines intended to standardize reporting, case reports remain difficult to compare and analyze at scale due to heterogeneous narrative structures. This study presents a visual analytics system for exploring clinical pathway information within case reports. CARE guidelines define a clinical event schema as including symptoms, examinations, diagnoses, treatments, and outcomes. Clinical event sequences were extracted and normalized from 1,320 case reports published between January 2010 and August 2025 using an LLM (Large Language Model), and clinical pathways were modeled using transition matrices. This system integrates bar charts, scatter plots, and Sankey diagrams to support linked exploration of clinical pathways within original texts. While most pathways showed high heterogeneity, distinct and consistent pathway patterns were visually identified within specific clusters. This work supports the viability of visual exploration centered analysis for case report-based clinical pathway studies.

**색인어** : 시각화, 임상 텍스트 분석, 대형 언어 모델, 정보 구조화, 의료 데이터 분석

**Keyword** : Visualization, Clinical Text Analysis, Large Language Model, Information Structuring, Medical Data Analysis

<http://dx.doi.org/10.9728/dcs.2026.27.3.739>



This is an Open Access article distributed under the terms of the Creative Commons Attribution Non-Commercial License (<http://creativecommons.org/licenses/by-nc/3.0/>) which permits unrestricted non-commercial use, distribution, and reproduction in any medium, provided the original work is properly cited.

**Received** 16 January 2026; **Revised** 02 February 2026

**Accepted** 02 February 2026

**\*Corresponding Author; Sung-Hee Kim**

**Tel:** +82-51-890-2366

**E-mail:** sh.kim@deu.ac.kr

## 1. 서론

증례보고는 단일 또는 소수 환자의 증상, 진찰 소견, 진단, 치료, 경과를 상세히 기술하는 관찰 연구의 한 형태로, 기존에 알려지지 않았거나 드물게 발생하는 질병 양상, 진단적 특징, 치료 반응, 부작용 등을 다룬다. 특히 희귀 질환이나 복합 질환처럼 대규모 임상시험이 어려운 분야에서는 의료 현장의 경험과 임상주의 판단 과정을 가장 직접적으로 반영하는 연구 형태로 기능한다[1]. 증례보고는 근거 위계상 낮은 수준으로 분류되며, 제한된 사례 수와 통제되지 않은 임상 환경으로 인해 질병의 원인이나 인과관계를 직접 규명하기는 어렵다. 그럼에도, 실제 임상에서 발생한 사건을 기반으로 새로운 가설이나 임상적 통찰을 제시함으로써, 이후 환자-대조군 연구, 코호트 연구, 무작위 대조 이중맹검 연구 등의 보다 체계적인 연구로 발전할 수 있는 출발점을 제공한다[2]. 이러한 확장을 위해서는 개별 증례에서 기술된 임상 서술이 구조적으로 해석되고 비교할 수 있는 형태로 정리될 필요가 있다.

그러나 증례보고는 대개 자유 서술식으로 작성되기 때문에 서술 방식과 구조가 증례마다 크게 다르다. 사건의 순서, 경과, 판단 근거가 명확히 구분되지 않고, 동일한 임상 상황도 서로 다른 표현과 논리 전개로 기술되는 경우가 많다[2]. CARE guidelines (CAse REport guidelines)와 같은 표준화된 작성 지침이 제시되어 있음에도[3], 실제 임상 현장의 다양한 서술 방식으로 인해 일관된 정보 구조를 확보하기 어렵다. 이러한 특성은 여러 증례에서 공통으로 나타나는 임상 경로나 사건 전개의 차이를 체계적으로 식별하는데 한계를 초래한다.

최근 임상 데이터를 구조화하고 시각적으로 분석하려는 연구가 활발히 이루어지고 있으나, 대부분은 전자 의무기록(EHR; Electronic Health Records)이나 임상 레지스트리 등 정형화된 코드 체계를 중심으로 발전해 왔다[4]-[7]. 하지만 코드 중심의 분석 방식은 환자의 경과가 상세히 기술되는 증례보고 특유의 임상 서술 맥락을 충분히 반영하기 어렵다[8]. 이와 같은 비정형 임상 서술을 구조화하기 위해 자연어 처리(NLP; Natural Language Processing)나 대형 언어 모델(LLM; Large Language Model)을 활용한 시도가 이루어지고 있으나, 여전히 문장 단위의 개체 식별에 머물러 있어 사건 간의 연속성이나 전이 관계를 체계적으로 재구성하는 연구는 부족하다[9]. 특히 기존의 텍스트 임베딩은 개별 사건의 유사성 식별에는 유리하나, 임상 경로의 핵심인 사건 간 순서와 인과관계를 보존하는 데는 취약하다. 따라서 서사 정보를 사건 단위로 재구조화하고 이를 시각적으로 탐색할 수 있는 새로운 접근이 필요하다[10],[11].

이에 본 연구에서는 기존 텍스트 기반 분석이 포착하지 못했던 사건 간 순서와 전이 관계를 임상 경로의 핵심 요소로 보고, 이를 분석하기 위해 LLM 기반 사건 추출과 전이행렬 임베딩을 결합한 방법론을 제안한다. 먼저 CARE guidelines를 참조하여 임상 서사를 ‘증상’, ‘검사’, ‘진단’, ‘치료’, ‘결과’

의 다섯 가지 핵심 축으로 재구성하고[3], LLM을 통해 주요 사건을 추출한다. 이후 의료 텍스트 임베딩 모델을 활용한 군집화를 통해 의미적으로 유사한 임상 사건을 통합함으로써 사건 수준의 정규화를 수행한다. 이렇게 정규화된 임상 경로에 전이행렬 임베딩을 적용하여, 사건 간 전이 패턴을 확률적으로 표현하고 환자별 임상 경로를 구조적으로 비교할 수 있는 형태로 변환한다. 전이행렬 기반 임베딩은 사건 간 관계를 독립적인 의미 벡터가 아닌 전이 확률로 표현함으로써, 텍스트 임베딩 기반 접근보다 임상 경로의 방향성과 흐름을 보다 명확히 반영할 수 있다는 점에서 차별성을 가진다[12]. 제안한 방법론은 임상 경로의 변이가 두드러지는 전통 보완 통합 의학(TCIM; Traditional, Complementary, and Integrative Medicine) 증례보고에 적용하였다[13]. TCIM은 동일 질환에서도 환자의 개인적 특성에 따라 진단과 치료 경로가 다양하게 전개되는 특성이 있어, 사건 전이 구조와 임상 경로의 복잡성을 검증하기에 적합한 사례군으로 판단하였다. 최종 시각화 시스템은 산점도와 Sankey 다이어그램을 결합하여 사건 흐름 탐색과 원문 기반 해석을 동시에 지원한다.

본 연구는 증례보고의 임상 서술을 사건 단위로 구조화하여 분석 가능한 데이터 형태로 전환하고 시각화하는 방법론을 제시한다. 이를 통해 개별 사례의 임상 맥락을 비교하고 사례 간 특성을 분석할 수 있다. 특히 전이행렬 기반 임베딩을 적용함으로써, 임상 사건 간의 순서와 전이 구조를 정량적으로 반영할 수 있도록 하였다. 이러한 접근은 기존 텍스트 기반 분석으로는 파악하기 어려웠던 임상적 의사결정의 과정과 패턴을 시각적으로 탐색할 수 있는 가능성을 제시한다. 또한 실제 임상에서 관찰된 사례를 토대로 새로운 가설과 통찰을 발굴하고, 이후의 체계적 연구로 확장될 수 있는 기반을 마련할 것으로 사료된다.

## II. 관련 연구

### 2-1 의료 데이터 시각화

의료 데이터는 진단, 검사, 치료, 예후 등 시간에 따라 연속적으로 발생하는 사건들로 구성되어 있어, 이 흐름을 시각적으로 표현하는 연구가 활발히 이루어져 왔다. 이러한 시간적 특성을 반영하기 위해 프로세스 마이닝 기법이 주로 활용되었으며, 특히 Sankey 다이어그램이나 흐름 기반 시각화 기법을 통해 환자 집단의 치료 경로나 질병 진행 과정을 파악하고자 하였다[5]. 예를 들어, Perer et al.의 CareFlow는 전자 의무기록(EMR; Electronic Medical Record)을 기반으로 유사 환자군의 치료 경로를 시각화하고, 치료의 전이 흐름과 결과를 비교할 수 있는 Sankey 다이어그램 기반 인터랙티브 도구를 제시하였다[6]. Zhang et al.은 인구 집단 수준의 EHR 데이터를 조건부 확률 기반으로 모델링하여 질병 및 치

료 사건의 시간 축을 기준으로 탐색할 수 있는 시스템을 구축하였고[7], Hjaltelin et al.은 덴마크 전 국민 데이터를 활용해 질병 간 전이 관계를 Sankey 다이어그램으로 표현함으로써 질병의 전개 과정을 시각화하였다[8]. 이처럼 기존 연구들은 정형화된 의료 사건 로그의 시간적 연속성과 전이 구조를 프로세스 마이닝 관점에서 모델링함으로써, 복잡한 임상 패턴을 효과적으로 탐색할 수 있도록 하였다.

다만 이러한 접근들은 모두 EHR이나 EMR 등 정형화된 임상 사건 로그를 전제로 설계되어 있어, 기존 프로세스 마이닝 기법을 서술형 임상 텍스트 자체에 직접 적용하기에는 한계가 있다. 증례보고는 방대한 임상 데이터 중에서도 특이하거나 학문적으로 의미 있는 사례를 임상가가 서사적으로 재구성한 결과물로, 단일 사례의 일반화에는 한계가 있으나 새로운 가설과 통찰을 제시하는 가치를 가진다. 따라서 EHR 및 EMR 기반 시각화가 실제 임상에서 일반화된 경향을 파악하고 의사결정 지원에 활용되는 데 초점을 두었다면, 증례보고의 시각화는 학문적 가치가 있는 개별 사례들이 어떤 공통점이나 패턴을 보이는지 탐색함으로써 상위 수준의 체계적 연구로 확장될 가능성을 보여주는 접근으로 볼 수 있다.

### 2-2 자연어처리 기반 증례보고 분석

한편, 증례보고나 임상 서술과 같은 비정형 임상 텍스트를 구조화하려는 연구도 활발히 이루어져 왔다. NLP 기반 접근에서는 임상 개체명 인식, 용어 정규화, 관계 추출 등을 통해 텍스트 내 증상, 진단, 치료 등의 개체를 자동으로 식별하고 구조화하려는 시도가 이루어졌다[9]. Schulz et al.은 PMC (PubMed Central)의 증례보고를 대상으로 의학 개체 추적 코퍼스를 구축하고, 여러 개체명 인식 모델을 적용하여 증례보고에서 개체를 인식하는 데 나타나는 어려움을 보여주었다[14]. 또한 Zhou et al.은 CREATE (Clinical Report Extraction and Annotation Technology) 시스템을 제안하여, 증례보고 내 임상 사건과 개체를 인식하고 이를 표준화한 뒤, 사건 간의 시간적 관계를 그래프 형태로 표현하였다[15].

그러나 이러한 연구들은 여전히 문장 또는 문서 단위에서의 개체 식별과 관계 추출에 초점을 두고 있어[16], 증례보고에서 사건들이 시간적 흐름 속에서 어떻게 누적되고 전환되는지를 경로 단위로 재구성하는 데에는 한계가 있다[17]. 즉, 개별 사건의 존재와 관계는 포착할 수 있으나, 사건 간 전이와 서사적 전개 구조를 통합적으로 표현하기는 어렵다.

### 2-3 LLM 기반 임상 서술 분석 및 요약

기존 NLP 기반 접근이 문장 또는 문서 단위의 분석에 머무른 반면, 최근 등장한 LLM은 보다 넓은 문맥을 고려하여 임상 서술을 요약하거나 사건 중심으로 재구성할 수 있는 가능성을 보인다[18],[19]. 특히 증례보고처럼 시간과 논리에

따라 결과가 기술되는 임상 서술의 경우, LLM을 활용하면 복잡한 내용을 임상 사건 단위로 명확히 구조화할 수 있다[20].

다만 LLM 기반 분석이 곧바로 임상적으로 신뢰할 수 있는 구조화 결과로 이어진다고 보기는 어렵다. LLM의 추론 과정은 내부적으로 불투명하며, 동일한 입력에 대해서도 일관되지 않은 출력을 생성하거나 사실과 다른 내용을 생성하는 문제(hallucination)가 존재한다. 본 연구는 이러한 한계를 인지한 상태에서 LLM을 증례보고의 사건 추출에 활용하되, 생성된 결과의 신뢰성과 일관성을 강화하기 위한 후속 검증이 필요함을 전제로 한다. 이러한 접근은 기존 NLP의 연장선상에서 LLM을 도입하여, 서술형 임상 텍스트를 사건 단위로 구조화하려는 연구의 발전 방향을 모색하는 시도로 볼 수 있다.

## III. 본 론

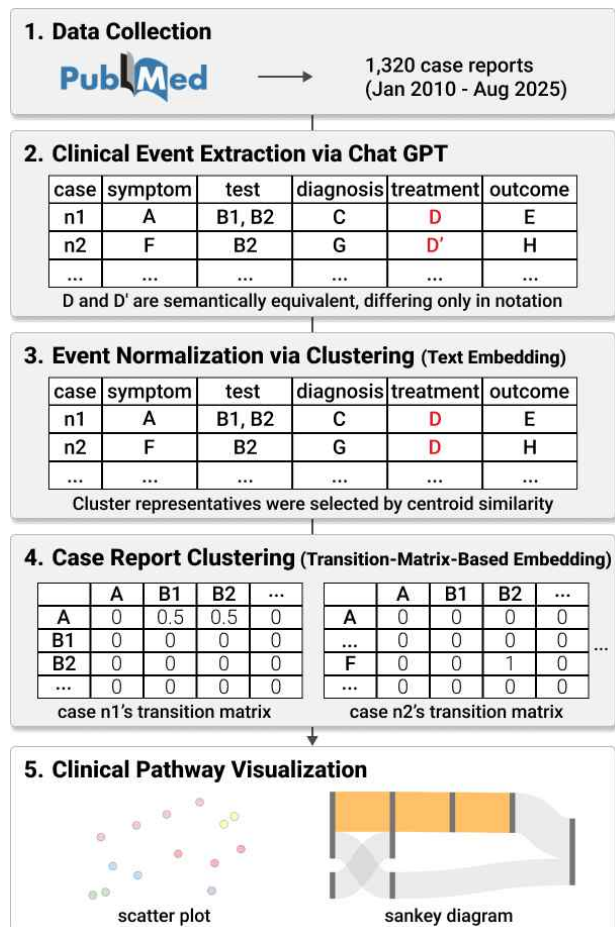


그림 1. 의학 증례보고 기반 임상 사건 구조화 및 시각화 워크플로  
Fig. 1. Structuring and visualization workflow for clinical information from medical case reports

본 연구에서는 증례보고로부터 임상 사건을 구조화하고 이를 시각적으로 탐색할 수 있는 시스템의 설계를 제안하고자

한다. 그림 1은 본 연구의 전체 워크플로를 요약한 것으로, 먼저 2010년 1월부터 2025년 8월까지 PMC에서 증례보고를 수집하여 분석 대상으로 선정하였다. 이후 CARE guideline을 참조하여 ‘증상’, ‘검사’, ‘진단’, ‘치료’, ‘결과’의 다섯 가지 임상 사건 축을 정의하고, ChatGPT를 활용하여 각 증례보고로부터 해당 사건 정보를 추출하였다. 추출된 다양한 표현들을 텍스트 임베딩 기반의 군집화 기법을 통해 정규화하여 유사 사건을 대표 표현으로 통합하였으며, 이어 사건의 순서와 전이 구조를 반영하기 위해 전이행렬 기반 임베딩을 적용하여 각 증례의 임상 경로를 백터로 표현하였다. 마지막으로 이러한 분석 결과를 바탕으로 산점도와 Sankey 다이어그램을 결합한 웹 기반 시각화 인터페이스를 구현하여, 사용자가 사례 간 패턴을 직관적으로 비교하고 개별 증례의 사건 흐름을 탐색할 수 있도록 하였다.

**3-1 데이터 수집 및 검색 전략**

연구 데이터는 PMC에 공개된 TCIM 분야의 증례보고를 대상으로 수집되었다. 수집 범위는 2010년 1월부터 2025년 8월까지 발표된 논문으로 설정되었으며, PMC의 Open Access 자료를 중심으로 데이터가 확보되었다. 데이터 수집은 PMC API를 활용하여 수행되었으며, 본 연구에서는 기존 문헌을 참고하여 표 1과 같이 MeSH (Medical Subject Headings)와 자연어(free-text) 키워드를 결합한 포괄적인 검색 전략을 구성하였다[21].

본 연구의 데이터 선정 및 검색 결과 검토는 단일 연구자에 의해 수행되었다. 1차 검색 결과 도출된 총 1,328편의 문헌에 대한 전수 검토의 물리적 한계를 고려하여, 사전에 정의된 포함 및 배제 기준에 따른 자동화된 필터링을 우선적으로 활용하였다. 데이터 선정의 주된 기준은 증례보고 형식의 여부와 원문 전문 확보 가능성이었으며, 이를 통해 최종 1,320편의 문헌을 분석 대상으로 확정하였다.

**표 1.** 문헌 검색 전략 및 검색어 조합

**Table 1.** Search strategy and query combinations

("Medicine, Traditional"[MH] OR "Complementary Therapies"[MH] OR "Medicine, Korean Traditional"[MH] OR "Medicine, Chinese Traditional"[MH] OR "Medicine, Mongolian Traditional"[MH] OR "Medicine, Tibetan Traditional"[MH] OR "Medicine, Kampo"[MH] OR "Medicine, African Traditional"[MH] OR "Medicine, Ayurvedic"[MH] OR "Medicine, Arabic"[MH] OR "Traditional medicine\*"[TIAB] OR Ayurved\*[TIAB] OR "alternative medicine"[TIAB] OR "alternative therapies"[TIAB] OR "complementary medicine"[TIAB] OR "complementary therapies"[TIAB] OR "integrative medicine"[TIAB] OR "integrated medicine"[TIAB] OR "traditional Chinese medicine"[TIAB] OR "Chinese medicine\*"[TIAB] OR "traditional Korean medicine"[TIAB] OR Kampo[TIAB] OR "Persian medicine"[TIAB] OR "traditional African medicine"[TIAB] OR "Oriental medicine\*"[TIAB] OR "traditional Oriental medicine"[TIAB] OR TCM[TIAB] OR KM[TIAB]) AND "Case Reports"[Publication Type] AND ("2010"[PDAT] : "3000"[PDAT])

본 연구의 핵심 목적은 특정 치료법의 효과 검증이 아닌, 비정형 텍스트의 사건 구조화 방법론 및 시각적 탐색 시스템 설계에 있다. 이에 따라 임상외에 의한 개별 증례의 질적 평가나 임상적 효능에 대한 검증은 연구 범위에서 제외하였으며, 이는 방대한 양의 임상 서술을 사건 단위로 자동 구조화하고 그 전이 패턴을 조망하는 시스템의 범용적 효용성을 확인하기 위한 의도적 설계 선택이다.

**3-2 임상 사건 추출**

본 연구는 임상 서술에서 핵심 사건을 체계적으로 추출하기 위해 CARE guideline[3]을 준용하여 임상 서사를 ‘증상’, ‘검사’, ‘진단’, ‘치료’, ‘결과’의 다섯 개의 축으로 재구성하였다. 또한 기존 연구에서 제시된 증례보고의 보고 목적 분류를 참조하여[22],[23], 증례 유형을 치료 및 관리, 진단 및 발견, 부작용 및 예기치 못한 사건, 기전 또는 기타의 네 가지로 분류하였다.

이를 위해 ChatGPT-4o-mini 모델을 활용하여 증례보고 본문으로부터 환자 단위의 임상 정보와 성별, 연령, 인종 등 메타데이터를 JSON 형식으로 추출하였다. 이때 모델은 정보 손실 방지와 상호 운용성을 동시에 확보하기 위해 ‘원문 표현’, ‘정제 표현’, ‘표준 코드’를 함께 출력하도록 설계되었다. 구체적으로 ‘증상’, ‘검사’, ‘치료’는 SNOMED-CT와 RxNorm을, ‘진단’은 ICD-10을 기준으로 출력하되, 표준화가 어려운 TCIM 개념은 원문 표현을 유지하였다.

특히 ‘결과’ 항목은 표 2와 같이 11개의 정규화된 표현으로 제한하고, 이를 의미적 방향성에 따라 ‘호전’, ‘유지’, ‘악화’의 3가지 범주로 분류하였다[24]. 분석 시 표준 코드를 최우선으로 하여 1차 선별 용어로 선정하되, 코드가 없거나 유효하지 않은 경우 정제 표현을 채택하였다. 그러나 선정된 1차 선별 용어에도 여전히 의미가 같은 용어들이 혼재하는 한계가 있어, 데이터 일관성을 확보하고자 다음 절에서 기술할 유사도 기반 군집화를 추가로 수행하였다.

**표 2.** 결과 항목의 의미적 방향성 분류

**Table 2.** Orientation-based classification of outcome labels

Outcome Label	Outcome Category
complete remission, partial remission, no recurrence, no complication, alive	positive
stable(no change), ongoing follow-up	neutral
recurrence, complication, progressive disease, death	negative

**3-3 군집화를 통한 사건 정규화**

선정된 1차 선별 용어에 남아있는 의미적 중복을 해결하기 위해, 임베딩 기반의 2차 정규화를 수행하였다. 이는 ‘식욕 부진’과 ‘식욕 저하’처럼 표현은 다르나 의미가 동일한 용어들을

하나의 대표 용어로 통합하는 과정이다.

먼저 의료 도메인에 특화된 Gemma-300M 모델[25]을 활용하여 용어 간 의미적 유사도를 벡터화하고, 이를 기반으로 HDBSCAN (Hierarchical Density-Based Spatial Clustering of Applications with Noise) 알고리즘을 적용하였다. HDBSCAN은 군집의 개수 (k)를 사전에 지정할 필요 없이 임의의 형태를 가진 군집을 발견할 수 있고, 노이즈를 효과적으로 식별한다는 장점이 있어 다양한 변이가 존재하는 임상 표현을 정규화하기 위해 선정하였다.

군집화 결과, 의미적으로 유사한 용어들이 하나의 그룹으로 묶였으며, 각 군집의 중심 벡터 (centroid)와 코사인 유사도가 가장 높은 항목을 ‘최종 대표 용어’로 선정하여 일괄 치환하였다. 이렇게 통합된 용어들은 환자별 임상 사건 시퀀스를 구성하고 전이행렬을 생성하는 기본 단위로 활용된다.

### 3-4 전이행렬 기반 임베딩 및 군집화

단순 텍스트 임베딩이 사건의 인과적 흐름을 반영하지 못하는 한계를 보완하기 위해, 본 연구에서는 ‘증상’, ‘검사’, ‘진단’, ‘치료’, ‘결과’의 5개 축으로 정렬된 환자별 임상 사건 시퀀스를 바탕으로 전이행렬을 구성하여 사건의 연속적 관계를 정량화하였다. 이때 특정 축에 동반 진단이나 병용 치료 등 복수의 사건이 존재하는 경우, 독립된 경로로 분리하여 개별 시퀀스로 확장하였다. 데이터의 왜곡을 방지하기 위해 분기된 각 시퀀스에는 가중치를 부여하여, 하나의 증례가 가지는 전체 가중치의 합이 1이 되도록 정규화하였다. 이렇게 확장된 시퀀스를 바탕으로 사건 간 전이 확률을 계산하여 행렬을 구축하였으며, 데이터 희소성 문제를 완화하기 위해 평활화 및 행·열 단위 정규화를 수행하여 빈도 편향을 보정하였다.

구축된 전이행렬의 차원 축소를 위해 주성분 분석 (PCA; Principal Component Analysis)을 선행한 후, 비선형 차원 축소 기법인 UMAP (Uniform Manifold Approximation and Projection)을 적용하였다. 이를 통해 임상 경로의 전역적 분포와 국소적 특징을 균형 있게 보존한 벡터 표현을 획득하였다. 최종적으로 UMAP 벡터에 K-means 군집화를 수행하였으며, 이때 군집 수 (k)는 200으로 설정하였다. 이는 단순한 통계적 최적값을 따르기보다, 시각화 시스템 내에서 사용자가 다양한 임상 패턴을 구체적으로 식별하고 탐색할 수 있도록 군집을 충분히 세분화하기 위함이다. 결과적으로 도출된 200개의 군집은 유사한 전이 패턴을 공유하는 그룹으로서, 시각화 시스템의 핵심 분석 단위로 활용된다.

### 3-5 임상 경로 시각화 시스템 설계

본 연구에서는 증례보고 기반 임상 사건의 전이 패턴을 탐색하기 위한 웹 기반 시각화 시스템을 설계하였다. 앞선 전처리 및 분석 과정을 통해 구축된 최종 데이터는 개별 증례의 메

타데이터, 추출된 임상 사건 시퀀스, 그리고 전이행렬 분석을 통해 도출된 군집 정보와 차원 축소 좌표를 포함한다. 시스템의 데이터 스키마는 표 3과 같이 정의되며, 이는 시각화 인터페이스의 각 구성요소가 유기적으로 연동되는 기반이 된다.

구축된 데이터를 효과적으로 탐색하기 위해, 시스템은 D3.js 라이브러리를 기반으로 구현되었다. 인터페이스는 군집 막대그래프, 산점도 그래프, Sankey 다이어그램, 사례 목록 등 다중 시각화 요소로 구성되며, 각 구성요소는 동일한 선택 상태를 공유하여 상호 연동된다. 전체적인 분석 흐름은 데이터의 전역적 분포를 확인하고, 사례 간 관계와 세부 전이 구조를 단계적으로 탐색한 뒤, 마지막으로 원문 사례에 접근하는 계층적 탐색 구조를 따른다. 또한 필터링 기능을 통해 탐색 범위를 조절할 수 있도록 하여, 사용자가 특정 조건에 맞는 사례를 집중적으로 분석할 수 있도록 지원하였다.

표 3. 시각화 시스템 데이터 스키마

Table 3. Data schema for the visualization system

Category	Variable	Usage	Description
Analysis Info	Cluster	Bar Chart (Group), Scatter Plot (Color)	Cluster ID assigned by K-means
	x_umap, y_umap	Scatter Plot (Position)	2D coordinates for visualization
Metadata	PMID	Interaction Key (Link to PubMed)	PubMed Identifier (Unique Key)
	PubYear	Filter (Bar Chart)	Publication year
	Case Type	Filter (Check Box)	Report type
Demo-graphics	Patients ID	Filter (Check Box)	Individual patients within a case
	Sex	Filter (Check Box)	Male/Female/Unknown
	Age Group	Filter (Check Box)	Child/Adult/Aged/Unknown
	Race	Filter (Check Box)	9 categories
Clinical Path	Symptom	Sankey Diagram (Nodes/Links)	Patient's initial symptoms
	Test	Sankey Diagram (Nodes/Links)	Medical examinations and tests performed
	Diagnosis	Sankey Diagram (Nodes/Links)	Diagnosed diseases or conditions
	Treatment	Sankey Diagram (Nodes/Links)	Therapeutic interventions applied
	Outcome	Sankey Diagram (Nodes/Links)	Final clinical outcome
	Outcome Role	Sankey Diagram (Link Color)	Positive/Negative/Neutral

#### 1) 군집 막대그래프

그림 2의 A에 제시된 가로형 막대그래프는 전이행렬 분석으로 도출된 군집 식별자와 사례 수를 매핑하여, 200개 군집 간의 규모 불균형을 시각적으로 드러낸다. 산점도나 네트워크 그래프가 데이터의 구조적 관계를 보여준다면, 이 막대그래프



그림 2. 전체 시각화 인터페이스 예시  
 Fig. 2. Example of the overall visualization interface

는 각 그룹의 정량적 규모를 직관적으로 비교하게 함으로써 연구자가 분석의 우선순위를 결정하는 데 기여한다. 이를 통해 가장 빈번하게 발생하는 보편적 경로와 빈도는 낮지만, 임상적으로 의미 있는 희귀 변이 경로가 명확히 구분된다. 또한 막대그래프는 상호작용의 기점으로 기능하여, 사용자가 특정 막대를 선택하면 해당 군집 정보가 산점도와 Sankey 다이어그램으로 전달되어 관심 영역을 강조하거나 필터링하는 역할을 수행한다.

2) 산점도 그래프

1,320건 증례 데이터의 전체적인 분포와 군집 간 유사성은 그림 2의 B의 산점도를 통해 나타낼 수 있다. 본 연구에서는 복잡한 임상 경로 벡터를 2차원 평면에 투영하기 위해 차원 축소된 UMAP 좌표를 활용하였으며, 이를 통해 유사한 전이 패턴을 가진 증례들이 공간상에 인접하게 배치되도록 하였다. 이때 각 점의 색상은 군집 식별자에 따라 구분되어 그룹 간의 경계와 밀집도를 명확히 보여준다.

나아가 이 산점도는 거시적 관점과 미시적 탐색을 매개하는 역할을 수행한다. 사용자가 산점도 위의 특정 점을 클릭하면, 시스템은 해당 점에 연결된 논문 고유 식별자를 인식하여 Sankey 다이어그램 상에서 시각적 강조 효과를 적용한다. 구체적으로 해당 군집의 전체 전이 구조는 유지하되 선택된 단일 증례의 경로만이 흰색으로 선명하게 표시되며, 나머지 군집 내 사례들은 회색으로 처리된다. 연구자는 이러한 대비 효과를 통해 집단의 평균적인 흐름 속에서 개별 사례가 갖는 고유한 경로 특성을 명확히 식별하고 비교할 수 있다.

3) Sankey 다이어그램

본 연구의 핵심 분석 도구인 그림 2의 C와 같은 Sankey 다이어그램은 임상 사건의 시간적 흐름과 단계별 인과 관계

를 규명하는 데 최적화되어 있다. 일반적인 그래프가 사건 간의 선후 관계나 분기되는 구조를 표현하기 어려운 반면, 이 시각화는 증상에서 결과로 이어지는 5단계의 임상 경로 데이터를 노드와 링크로 연결하여 사건의 진행 순서를 명확히 시각화한다. 특히 링크의 두께는 전이 빈도를 나타내어, 이를 통해 굵게 표시되는 주류 경로를 통해 보편적인 임상 경향을 직관적으로 파악할 수 있을 뿐만 아니라, 가늘게 연결된 희귀한 변이 경로나 소수의 특이 사례까지 놓치지 않고 포착하여 임상적 다양성을 입체적으로 조망할 수 있다. 특히 결과 노드로 향하는 링크는 예후 분류 정보에 따라 호전은 초록색, 유지는 노란색, 악화는 빨간색으로 구분되어, 복잡한 경로 속에서도 임상 결과의 긍정 및 부정 양상을 즉각적으로 판별할 수 있다. 아울러 단일 증례 내에 복수 환자의 사례가 존재하는 경우, 환자 식별자 기반의 체크박스를 통해 개별 환자의 경로를 선택적으로 활성화함으로써 동일한 증례 내에서 환자 간 임상 경과 차이를 정밀하게 비교할 수 있다.

4) 필터링 및 사례 탐색 인터페이스

단순한 데이터 나열을 넘어 연구자가 수렴한 임상적 가설을 검증할 수 있도록, 그림 2의 D와 같은 동적 필터링 인터페이스를 설계하였다. 전체 데이터를 한 번에 관찰하는 것은 일반적 경향 파악에는 유리하나, 특정 연령대나 질환군에 특화된 패턴을 발견하기는 어렵다. 이에 시스템은 발행 연도, 성별, 연령대, 인종 등 메타데이터 변수를 조작할 수 있는 기능을 제공하여, 사용자가 관심 있는 하위 집합을 추출해 집중적으로 분석할 수 있게 하였다. 조건 변경 시 모든 시각화 화면은 실시간으로 갱신되어 해당 조건에 부합하는 증례들의 분포와 경로만을 다시 그린다.

하단의 논문 목록 또한 시각화 화면과 연동되어 필터링된 증례의 서지 정보를 제공하며, 항목 클릭 시 논문 고유 식별

자를 통해 PubMed 원문 페이지로 연결된다. 이는 시각화에서 발견된 추상적 패턴이 실제 문헌의 텍스트 맥락과 일치하는지 원문을 통해 교차 검증할 수 있도록 지원한다.

#### IV. 연구 결과

##### 4-1 군집 구조의 안정성 검증

본 절에서는 제안된 시스템이 고해상도 임상 경로 탐색을 안정적으로 지원할 수 있는지를 검증하였다. 본 시스템은 임상 경로의 대표적인 흐름뿐만 아니라, 소수 사례에서 나타나는 미세한 변이와 예외적 패턴까지 탐색하는 것을 목표로 한다. 이를 위해 군집 수를 200개로 설정하여 데이터를 고도로 세분화하였으며, 이러한 설정 하에서도 군집 구조가 시각적 분석에 적합한 형태로 유지되는지가 핵심 검증 대상이 된다.

이를 확인하기 위해 전이행렬 임베딩의 품질을 SS (Silhouette Score), CH(Calinski-Harabasz Index), DB(Davies-Bouldin Index)를 통해 정량적으로 측정하였다. 비교 대상으로는 기존 의료 특화 텍스트 임베딩 모델인 Gemma-300M[25] 및 SapBERT[26]를 선정하였으며, 공정한 비교를 위해 세 방식 모두 PCA, UMAP, K-means를 순차적으로 적용하는 동일한 분석 파이프라인을 거치도록 하였다. 분석은 시각화 시스템이 요구하는 군집 수 200개 조건에 주안점을 두되, 공정한 평가를 위해 각 모델의 최적 군집 수 조건에서의 지표도 함께 산출하여 표 4에 제시하였다.

**표 4.** 임베딩 방법과 K 값 변화에 따른 군집화 성능 비교  
**Table 4.** Clustering quality comparison across embedding methods and k values

Method	k	SS	CH	DB
Gemma-300M	17	0.275025	569.242004	3.018630
Gemma-300M	200	-0.197054	91.010368	5.346961
SapBERT	41	-0.091469	95.983482	7.820201
SapBERT	200	-0.402637	34.124847	13.243575
<b>Transition Embedding</b>	<b>200</b>	<b>0.475293</b>	<b>7359.750977</b>	<b>0.645828</b>

분석 결과, 전이행렬 임베딩은 군집 수 200개 조건에서 SS=0.48, CH=7359.75, DB=0.65를 기록하여 가장 우수한 구조적 안정성을 보였다. 이는 단순한 수치적 우수성을 넘어, 세밀하게 분할된 임상 경로들이 시각적 공간 내에서 서로 중첩되거나 붕괴되지 않고, 해석 가능한 구조로 유지됨을 의미한다. 반면 동일한 조건에서 텍스트 임베딩 모델들은 성능이 급격히 저하되는 양상을 보였다. Gemma-300M은 SS=-0.20, CH=91.01, DB=5.35를 기록하였으며, SapBERT 또한 SS=-0.40, CH=34.12, DB=13.24로 나타나 군집 간 경계가

붕괴됨을 보였다.

주목할 점은 텍스트 모델이 잠재적 성능을 최대한 발휘할 수 있는 최적 군집 수 조건과 비교하였을 때에도 전이행렬 방식의 적합성이 확인된다는 것이다. 텍스트 모델 중 가장 양호한 성능을 보인 Gemma-300M은 최적 조건인 17개 군집에서 SS=0.28, CH=569.24, DB=3.02로, 전이행렬 방식이 k=200의 고해상도 조건에서 달성한 성과에 미치지 못했다. 이는 전이행렬 방식이 사건의 흐름을 반영하여 데이터 세분화에 최적화되었음을 방증하며, 결과적으로 본 시스템이 통계적으로 신뢰할 수 있는 탐색 환경을 갖췄음을 확인시켜 준다.

##### 4-2 임상 경로의 패턴 유형화 분석

안정적으로 구축된 군집 구조를 기반으로, 각 군집 내에서 가장 빈번하게 발생하는 주 경로가 차지하는 비율을 경로 집중도로 정의하고, 이를 식 (1)과 같이 산출하였다.

$$C_{path}(i) = \frac{\max_j(n_{ij})}{N_i} \tag{1}$$

여기서  $C_{path}(i)$ 는  $i$ 번째 군집의 경로 집중도,  $n_{ij}$ 는  $i$ 번째 군집 내에서 가장 많이 관측된 단일 경로 유형( $j$ )의 빈도수를,  $N_i$ 는  $i$ 번째 군집에 포함된 전체 증례 수를 의미한다. 이 값은 해당 군집 내 환자들의 임상 경로가 동일한 흐름을 공유하는지, 혹은 서로 다른 경로로 분산되는지를 나타내는 척도로, 1에 가까울수록 경로의 일치도가 높음을 의미한다.

분석 결과, 표 5와 같이 전체 평균 경로 집중도는 0.048로 나타났으며, 전체 군집은 일관된 패턴 (Coherent Pattern)과 이질적 패턴 (Heterogeneous Pattern)으로 분류되었다.

**표 5.** 경로 집중도 기반 군집 유형화 분석 결과  
**Table 5.** Analysis results of cluster categorization based on path concentration

Pattern Type	Concentration Criteria	Count (N)	Percentage
Coherent	$C_{path} \geq 0.4$	4	2.0%
Intermediate	$0.1 \leq C_{path} < 0.4$	19	9.5%
Heterogeneous	$C_{path} < 0.1$	177	88.5%
Total	Mean = 0.048	200	100%

세부 분포를 살펴보면, 전체의 약 2%에 해당하는 4개 군집 (N=4)은 경로 집중도가 0.4 이상인 일관된 패턴으로 식별되었다. 이들은 전체 데이터 대비 낮은 비중임에도 불구하고 군집 내에서 뚜렷한 공통 경로가 형성된 집단이다. 반면, 전체의 88.5%를 차지하는 177개 군집은 집중도가 0.1 미만인 이질적 패턴으로 나타나, 대다수 사례가 단일 흐름 없이 환자별로 다양한 변이를 보임을 시사한다.

결과적으로 본 분석은 임상 데이터의 대다수가 높은 복잡



상적으로 중요한 예외 사례들을 주류 경로에 병합시키는 경향이 있다. 그러나 임상 연구에서는 표준적인 치료 과정뿐만 아니라, 합병증이나 특이 반응과 같은 변이를 식별하는 것이 매우 중요하다. 따라서 시스템은 군집의 해상도를 충분히 높게 설정하여, 시각화 화면상에서 주류 경로와 미세한 변이 경로가 명확히 분리되어 표현되도록 설계해야 한다.

둘째, 거시적 관점과 미시적 관점을 오갈 수 있는 시각화 간의 상호 연동이 필수적이다. Sankey 다이어그램은 전체 경로의 패턴을 파악하는 데 적합한 반면, 산점도는 개별 증례의 분포와 이상치를 탐색하는 데 유리하다. 효과적인 분석을 위해서는 이 두 가지가 독립적으로 존재하지 않고 유기적으로 연결되어야 한다. 즉, 사용자가 전체 흐름을 보다가 특정 영역을 선택했을 때, 즉시 개별 사례의 분포와 세부 경로를 확인할 수 있는 단계적인 탐색 구조가 제공되어야 한다.

셋째, 연구자가 가설을 수립하고 원문을 통해 이를 검증할 수 있는 분석 환경을 지원해야 한다. 시스템의 역할은 고정된 분석 결과를 보여주는 것에 그치지 않는다. 연구자가 시각화된 패턴을 보며 의문을 제기하고, 메타데이터 필터링을 통해 조건을 좁혀가며, 최종적으로 원문 텍스트를 확인하여 사실관계를 검증하는 일련의 과정이 끊김 없이 이어져야 한다. 이러한 데이터와 텍스트 간의 교차 검증 기능은 시각화된 통계 정보가 실제 임상 맥락과 일치하는지를 확인하는 데 핵심적인 역할을 수행한다. 종합하면, 임상 경로 분석 시스템은 데이터의 복잡성을 인위적으로 단순화하기보다는, 그 복잡한 구조 안에서 연구자가 유의미한 질서를 발견하고 해석할 수 있도록 해상도, 연동성, 검증 가능성을 중심으로 설계되어야 한다.

### 5-3 한계점 및 향후 연구 방향

본 연구는 시각적 분석 시스템의 구조적 가능성과 설계적 시사점을 제시하였으나, 다음과 같은 한계와 과제를 남긴다.

첫째, 세분화된 군집 설정에 따른 전반적 경향 파악의 어려움이다. 본 연구는 개별 증례의 구체적인 임상적 변이를 보존하기 위해 군집의 수를 높게 설정하여 데이터를 세분화하였다. 이는 각 군집 내에서 명확한 패턴을 도출하고 특이 사례를 발견하는 데에는 매우 효과적이거나, 반대로 상위 범주의 포괄적인 흐름을 파악하기 위해서는 분산된 다수의 군집을 연구자가 개별적으로 확인하고 종합해야 하는 과제를 남긴다. 향후 연구에서는 세부적인 패턴 탐색을 유지하면서도, 필요에 따라 상위 수준에서 데이터를 통합적으로 조망할 수 있는 계층적 시각화 또는 다중 스케일 탐색 기능의 도입이 필요하다.

둘째, 실제 연구자를 대상으로 한 실증적 사용자 평가의 부재다. 활용 시나리오를 통해 시스템의 잠재력을 보였으나, 이것이 실제 연구 현장에서 가설 수립의 효율성이나 해석의 깊이를 얼마나 증진하는지는 검증되지 않았다. 시스템의 실무적 효용성을 입증하기 위해서는 향후 의료 전문가가 참여하는 정량적·정성적 사용자 실험을 수행하여, 시각화가 연구 성과에 미치는 영향을 객관적으로 측정해야 한다.

셋째, 데이터 수집 대상 및 도메인의 국한성이다. 본 연구는 문헌검색을 PubMed로 한정하였으나, 단일 데이터베이스의 활용은 검색의 포괄성을 제한할 수 있다[27]. 또한, 본 시스템은 TCIM 증례보고 데이터를 기반으로 설계되었기에, 질환 특성이나 서술 구조가 다른 타 임상 데이터에도 동일한 설계 원칙이 유효한지 확인이 필요하다. 향후 연구에서는 검색 데이터베이스의 범위를 확장하고, 타 임상 도메인에 적용 가능성을 검토함으로써 시스템의 범용성을 높일 필요가 있다.

넷째, 자동화된 문헌 수집 과정에서 발생할 수 있는 데이터의 노이즈 문제이다. 본 연구는 검색식에 "Case Reports"[Publication Type] 태그를 사용하여 분석 대상을 한정하였으나, 전문가에 의한 개별 원문 검토 과정을 거치지 않아 데이터베이스 색인 오류에 따른 리뷰 논문이나 임상 시험 데이터가 일부 혼입되었을 가능성을 배제할 수 없다. 이러한 데이터셋의 불확실성은 분석 결과에서 나타난 높은 이질성(88.5%)에 영향을 미쳤을 수 있으며, 향후 연구에서는 PRISMA (Preferred Reporting Items for Systematic Reviews and Meta-Analyses) 가이드라인에 준하는 단계적 선별 절차를 도입하여 데이터의 정밀도를 확보해야 할 것이다[28].

## VI. 결 론

본 연구는 임상 경로 데이터의 복잡성과 이질성을 탐색적으로 이해하기 위한 시각적 분석 시스템을 제안하고, 그 설계 가능성을 논의하였다. 제안된 시스템은 세분화된 임상 경로 수준에서의 시각적 탐색을 전제로 대표적인 임상 경로와 그로부터 파생되는 다양한 변이 양상을 동시에 드러내도록 설계되었다. 또한 Sankey 다이어그램과 산점도 기반 시각화를 연계함으로써 연구자가 서로 다른 시각 수준을 오가며 분석할 수 있는 환경을 제공한다. 이러한 접근은 임상 경로를 단일한 요약 결과로 환원하기보다 구조적 패턴과 개별 사례의 차이를 함께 탐색할 수 있도록 지원한다는 데 의의가 있다. 본 연구는 시각화 시스템이 단순한 결과 표현 도구를 넘어, 연구자의 가설 형성과 해석 과정을 보조하는 분석 환경으로 기능할 수 있음을 시사하며, 향후 임상 데이터 분석을 위한 시각적 탐색 시스템 설계에 참고 가능한 방향성을 제시한다.

## 감사의 글

이 논문은 정부(과학기술정보통신부)의 재원으로 정보통신기획평가원의 지원을 받아 수행된 지역지능화혁신인재양성사업(IITP-2026-RS-2020-II201791, 50%), 정부(과학기술정보통신부)의 재원으로 한국연구재단(RS-2025-00519817, 50%)과 부산광역시 및 (재)부산테크노파크의 BB21plus 사업으로 지원된 연구임.

## 참고문헌

- [1] T. Nissen, and R. Wynn, "The Clinical Case Report: A Review of Its Merits and Limitations," *BMC Research Notes*, Vol. 7, No. 1, 264, April 2014, <https://doi.org/10.1186/1756-0500-7-264>
- [2] S. M. Baek, J. H. Park, S. H. Lee, S. G. Kim, J. H. Lee, B. Y. Kim, and S. M. Choi, "Traditional Korean Medicine Doctors' Awareness and Utilization of the Case Report," *Korean Journal of Acupuncture*, Vol. 29, No. 1, pp. 57-70, March 2012. <https://kiss.kstudy.com/Detail/Ar?key=3035211>
- [3] J. J. Gagnier, G. Kienle, D. G. Altman, D. Moher, H. Sox, and D. Riley, "The CARE Guidelines: Consensus-based Clinical Case Reporting Guideline Development," *Global Advances in Integrative Medicine and Health*, Vol. 2, No. 5, pp. 38-43, September 2013. <https://doi.org/10.7453/gahmj.2013.008>
- [4] A. Perer and D. Gotz, "Data-Driven Exploration of Care Plans for Patients," in *Proceedings of the CHI '13 Extended Abstracts on Human Factors in Computing Systems (CHI EA '13)*, Paris: France, pp. 439-444, April 2013. <https://doi.org/10.1145/2468356.2468434>
- [5] A. Ledesma, N. Bidargaddi, J. Strobel, G. Schrader, H. Nieminen, I. Korhonen, and M. Ermes, "Health Timeline: An Insight-based Study of a Timeline Visualization of Clinical Data," *BMC Medical Informatics and Decision Making*, Vol. 19, 170, August 2019, <https://doi.org/10.1186/s12911-019-0885-x>
- [6] T. Zhang, T. H. McCoy, R. H. Perlis, F. Doshi-Velez, and E. Glassman, "Interactive Cohort Analysis and Hypothesis Discovery by Exploring Temporal Patterns in Population-Level Health Records," in *Proceedings of the 2021 IEEE Workshop on Visual Analytics in Healthcare (VAHC)*, New Orleans: LA, pp. 14-18, October 2021. <https://doi.org/10.1109/VAHC53616.2021.00007>
- [7] M. Pajusalu, K. Mooses, M. Oja, S. Tamm, M. Haug, and R. Kolde, "TrajectoryViz: Interactive Visualization of Treatment Trajectories," *Informatics in Medicine Unlocked*, Vol. 49, 101558, 2024. <https://doi.org/10.1016/j.imu.2024.101558>
- [8] N. Menachemi and T. H. Collum, "Benefits and Drawbacks of Electronic Health Record Systems," *Risk Management and Healthcare Policy*, Vol. 4, pp. 47-55, May 2011. <https://doi.org/10.2147/RMHP.S12985>
- [9] I. Spasic and G. Nenadic, "Clinical Text Data in Machine Learning: Systematic Review," *JMIR Med Inform*, Vol. 8, No. 3, e17984, March 2020. <https://doi.org/10.2196/17984>
- [10] N. A. Abudiyab and A. T. Alanazi, "Visualization Techniques in Healthcare Applications: A Narrative Review," *Cureus*, Vol. 14, No. 11, e31355, November 2022. <https://doi.org/10.7759/cureus.31355>
- [11] J. Kenei and E. Opiyo, "Modeling and Visualization of Clinical Texts to Enhance Meaningful and User-Friendly Information Retrieval," *Medical Sciences Forum*, Vol. 10, No. 1, 9, 2022. <https://doi.org/10.3390/IECH2022-12294>
- [12] W. Plumb, A. Bottle, G. Casale, and A. Liddle, "Clinical Pathway Clustering Using Surrogate Likelihoods and Replayability Validation," in *Proceedings of the 2023 Winter Simulation Conference (WSC)*, San Antonio: TX, pp. 1220-1231, December 2023. <https://doi.org/10.1109/WSC60868.2023.10407280>
- [13] J. Y. Ng, H. Cramer, M. S. Lee, and D. Moher, "Traditional, complementary, and integrative medicine and artificial intelligence: Novel opportunities in healthcare," *Integrative Medicine Research*, Vol. 13, No. 1, 101024, March 2024, <https://doi.org/10.1016/j.imr.2024.101024>
- [14] S. Schulz, J. Ševa, S. Rodriguez, M. Ostendorff, and G. Rehm, "Named Entities in Medical Case Reports: Corpus and Experiments," in *Proceedings of the Twelfth Language Resources and Evaluation Conference*, Marseille: France, pp. 4495-4500, May 2020. <https://aclanthology.org/2020.lrec-1.553/>
- [15] Y. Zhou, W.-T. Chen, B. Zhang, D. Lee, J. H. Caufield, K.-W. Chang, ... and W. Wang, "CREATe: Clinical Report Extraction and Annotation Technology," in *Proceedings of the 2021 IEEE 37th International Conference on Data Engineering (ICDE)*, Chania, Greece, pp. 2677-2680, April 2021. <https://doi.org/10.1109/ICDE51399.2021.00302>
- [16] D. F. Navarro, K. Ijaz, D. Rezazadegan, H. Rahimi-Ardabili, M. Dras, E. Coiera, and S. Berkovsky, "Clinical Named Entity Recognition and Relation Extraction Using Natural Language Processing of Medical Free Text: A Systematic Review," *International Journal of Medical Informatics*, Vol. 177, 105122, September 2023. <https://doi.org/10.1016/j.ijmedinf.2023.105122>
- [17] S. C. Fanni, L. Tumminello, V. Formica, F. P. Caputo, G. Aghakhanyan, I. Ambrosini, ... and E. Neri, "The Journey from Natural Language Processing to Large Language Models: Key Insights for Radiologists," *Journal of Medical Imaging and Interventional Radiology*, Vol. 11, 43, December 2024. <https://doi.org/10.1007/s44326-024-00043-w>
- [18] Y. Wang, L. Wang, M. Rastegar-Mojarad, S. Moon, F. Shen, N. Afzal, ... and H. Liu, "Clinical Information Extraction Applications: A Literature Review," *Journal of Biomedical Informatics*, Vol. 77, pp. 34-49, January 2018.

<https://doi.org/10.1016/j.jbi.2017.11.011>

- [19] S. Wu, K. Roberts, S. Datta, J. Du, Z. Ji, Y. Si, ... and H. Xu, "Deep Learning in Clinical Natural Language Processing: A Methodical Review," *Journal of the American Medical Informatics Association: JAMIA*, Vol. 27, No. 3, pp. 457-470, March 2020. <https://doi.org/10.1093/jamia/ocz200>
- [20] S. Nerella, S. Bandyopadhyay, J. Zhang, M. Contreras, S. Siegel, A. Bumin, ... and P. Rashidi, "Transformers and large language models in healthcare: A review," *Artificial Intelligence in Medicine*, Vol. 154, 102900, 2024. <https://doi.org/10.1016/j.artmed.2024.102900>
- [21] C. Y. Kwon, "Using Artificial Intelligence for the Development of a Living Evidence Map: The Pharmacopuncture Example," *Integrative Medicine Research*, Vol. 14, No. 4, 101217, December 2025. <https://doi.org/10.1016/j.imr.2025.101217>
- [22] A. Erol, "Basics of Writing Case Reports," *Noro Psikiyatri Arsivi*, Vol. 60, No. 1, pp. 1-2, February 2023. <https://doi.org/10.29399/npa.28403>
- [23] Heart Views Editorial Board, "Guidelines to Writing a Clinical Case Report," *Heart Views*, Vol. 18, No. 3, pp. 104-105, September 2017. <https://doi.org/10.4103/1995-705X.217857>
- [24] Clinical Care Classification. Nursing Outcomes Framework [Internet]. Available: <https://clinicalcareclassification.org/framework/nursing-outcomes/>.
- [25] SentenceTransformers. EmbeddingGemma-300m Finetuned on the Medical Instruction and Retrieval Dataset (MIRIAD) [Internet]. Available: <https://huggingface.co/sentence-transformers/embeddinggemma-300m-medical>.
- [26] Cambridgeltl. SapBERT [Internet]. Available: <https://huggingface.co/cambridgeltl/SapBERT-from-PubMedBERT-full-text>.
- [27] T. F. Frandsen, C. Moos, C. I. L. H. Marino, and M. B. Eriksen, "Supplementary Databases Increased Literature Search Coverage Beyond PubMed and Embase," *Journal of Clinical Epidemiology*, Vol. 181, 111704, May 2025, <https://doi.org/10.1016/j.jclinepi.2025.111704>
- [28] M. J. Page, J. E. McKenzie, P. M. Bossuyt, I. Boutron, T. C. Hoffmann, C. D. Mulrow, ... and D. Moher, "The PRISMA 2020 Statement: An Updated Guideline for Reporting Systematic Reviews," *BMJ*, Vol. 372, n71, March 2021. <https://doi.org/10.1136/bmj.n71>



**전혜리(Hye-Li Jeon)**

2022년~현 재: 동의대학교 인공지능학과 학사과정

※ 관심분야: Artificial Intelligence, HCI(Human-Computer Interaction), Data Visualization, Healthcare 등



**권찬영(Chan-Young Kwon)**

2013년: 동의대학교 B.S., KMD

2016년: 경희대학교 M.S.

2020년: 경희대학교 Ph.D

2020년~현 재: 동의료원 한방신경정신과 과장

2020년~현 재: 동의대학교 한의과대학 한방신경정신과 조교수

※ 관심분야: Mind-body Medicine, Mental Health, Digital Therapeutics 등



**김성희(Sung-Hee Kim)**

2006년: 이화여자대학교 B.S.

2008년: 이화여자대학교 M.S.

2014년: Purdue University Ph.D

2015년~2017년: 삼성전자 소프트웨어센터 책임연구원

2019년~2023년: 동의대학교 빅데이터인공지능센터 소장

2017년~현 재: 동의대학교 산업ICT기술공학과 부교수

2024년~현 재: 동의대학교 인공지능 대학원 주임교수

※ 관심분야: HCI(Human-Computer Interaction), User-centered Artificial Intelligence, Data Visualization 등