# Patient Journey: Temporal Disease Pattern Identification Using Machine Learning and Data Mining

Yahui Tian[1], Nan Shao[1], Zheng Zhu[1], Ziwei Yin[1], Linwei Li[1] │ [1]Boehringer Ingelheim

Abstract:

Patient journey is a collection of topics including understanding patient experience in the healthcare ecosystem in order to improve patient care, as well as understanding the disease development paths and their relation to clinical outcomes to potentially facilitate clinical decision making. Electronic health records (EHR) and health insurance administrative databases have provided rich and inexpensive sources of information for patient journey research. We have looked into temporal disease pattern identification in such databases for a rare disease to improve disease understanding, and we have applied sequential pattern mining and network analysis in the exploratory phase. Many results are in alignment with clinical understanding of the disease and we hope it can also facilitate hypothesis generation for further research. We will discuss methods and limitations and share experience for analyzing such large-scale health databases.

## BACKGROUND & STUDY OBJECTIVES

- Machine learning and Data mining in healthcare
  - Advantage: discover unknown features through data driven process
  - Applications: natural disease history, comorbidity, cohort identification, risk prediction, biomarker discovery, etc.
- Topics in patient journey research includes
  - Understanding patient experience in the healthcare ecosystem in order to improve patient care
  - **Understanding the disease development paths and their relation to clinical outcomes to potentially facilitate clinical decision making**
- Potential rich and inexpensive source of information for such research:
  - Electronic health records (EHR)
  - Health insurance administrative databases
  - Characteristics
    - ➤ Primary goal: document patients' care; reimbursement
    - ➤ Data format: structured, semi-structured and unstructured
    - ➤ Advantages: large sample size, long follow-up, cost-effective and time-saving source of research
    - ➤ Challenges: censoring, irregular time series data, completeness, correctness and confounding effects
- **Objective**: identify temporal disease pattern leading to the initial diagnosis of a rare disease using claims and EHR data to improve disease understanding.

## METHODS

Workflow:
- Initial disease cohort selection
  - Based on disease definition and corresponding ICD codes, patient age and study period
- Analysis cohort
  - Further selection of cohort based on length of information coverage in the database, e.g. length of continuous enrollment in health insurance program
  - A trade-off between length of history and # of patients included in the analysis
- Data pre-processing
  - Remove codes which are too general or not relevant to the disease of interest (e.g. codes correspond to general office visit)
- Apply two methods
  - Sequential pattern mining (Figure 1)
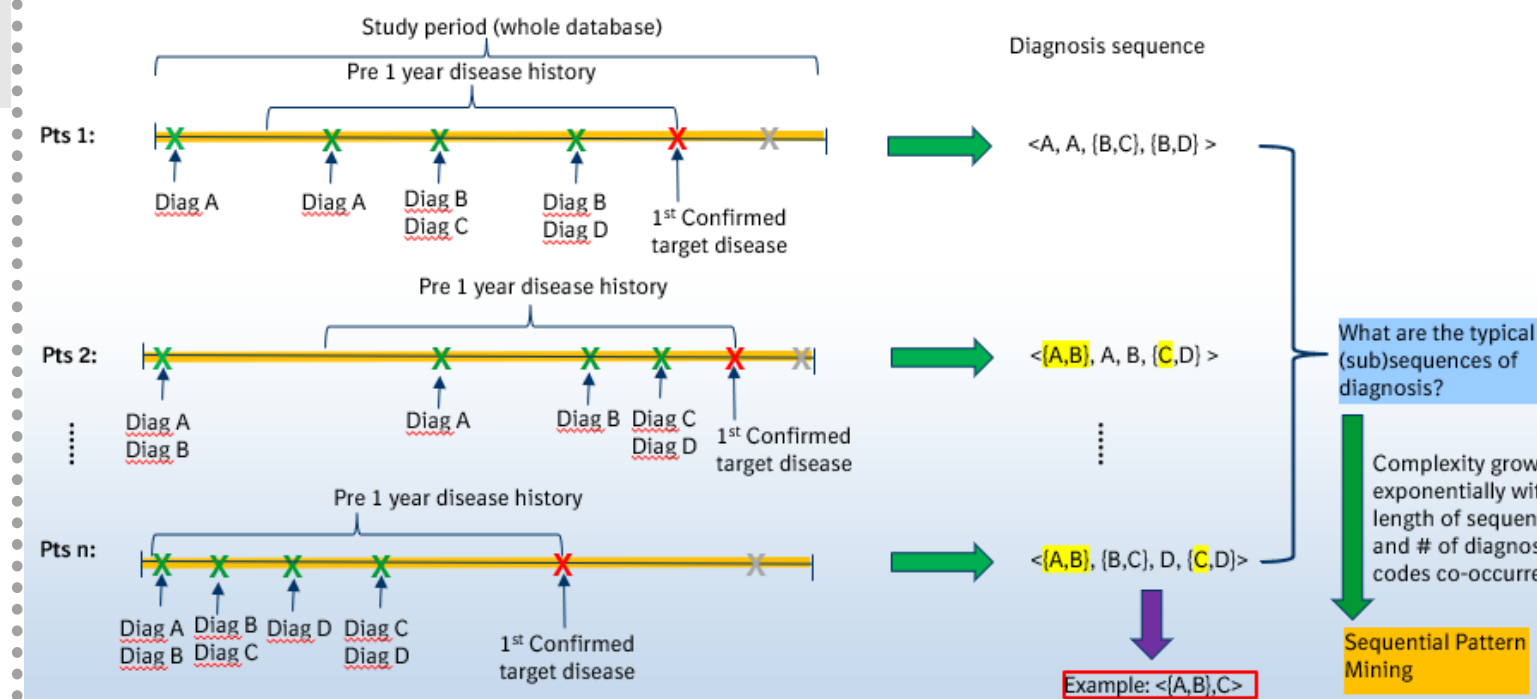  - Network analysis (Figure 2)



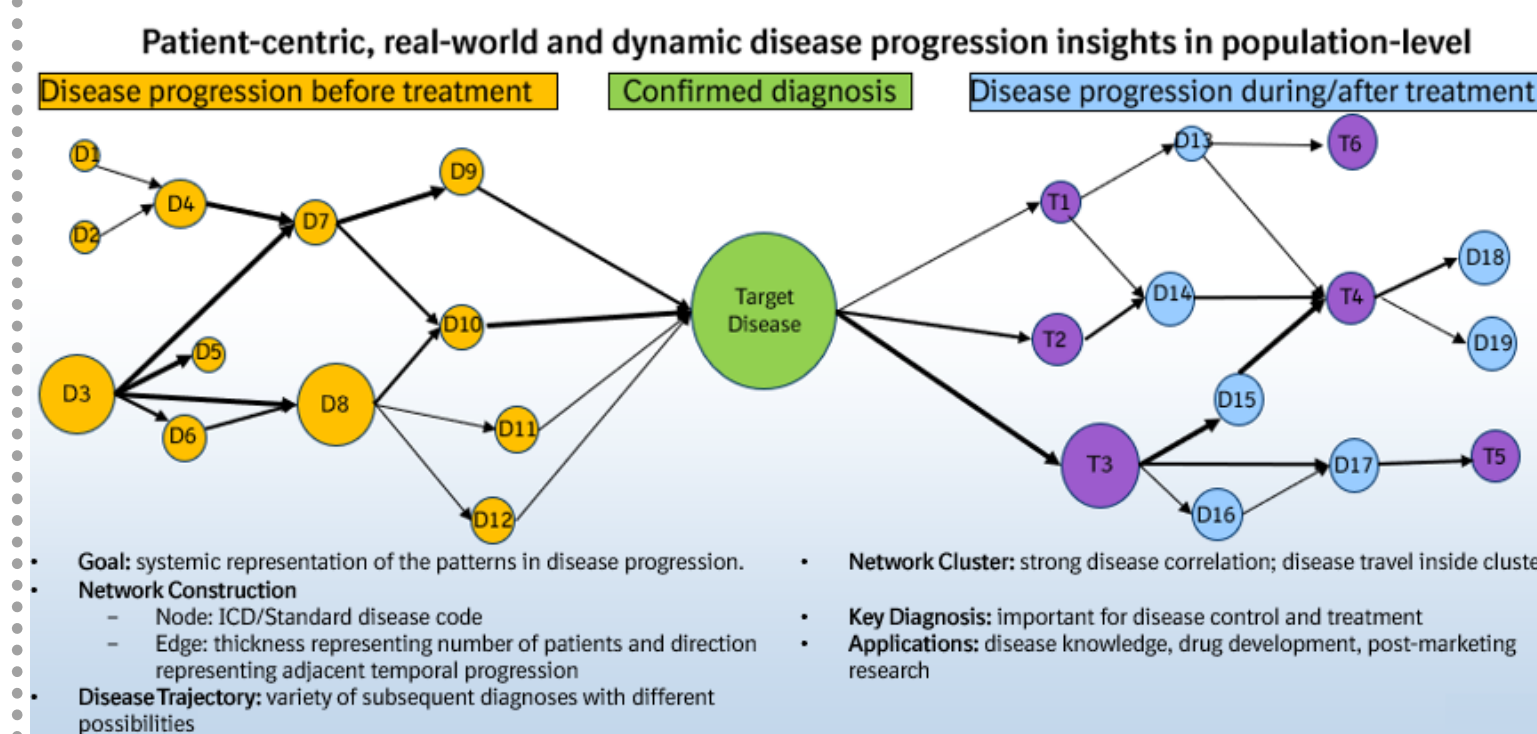Figure 1: Illustration of Sequential Pattern Mining



Figure 2: Illustration of Network Analysis

## RESULTS

### Sequential Pattern Mining
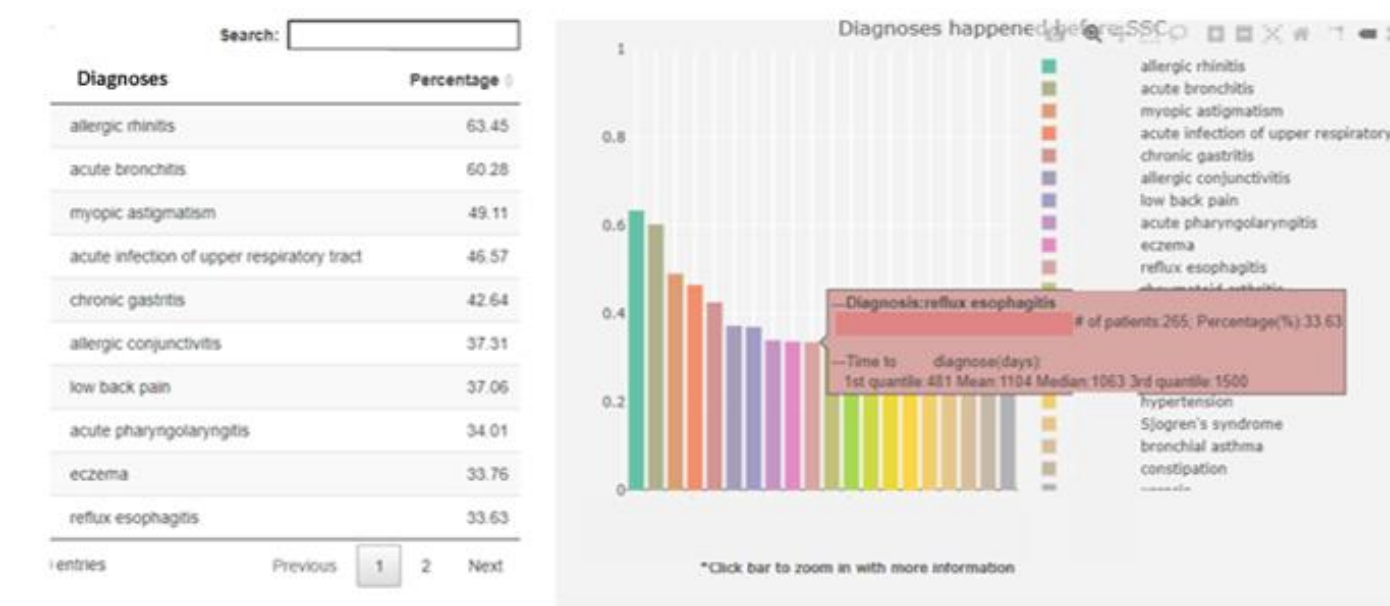
- Representative diagnoses before target disease



Figure 3: Interactive digital platform of representative diagnoses. Table on the left shows percentage. Histogram on the right shows detailed information same as Figure 3.
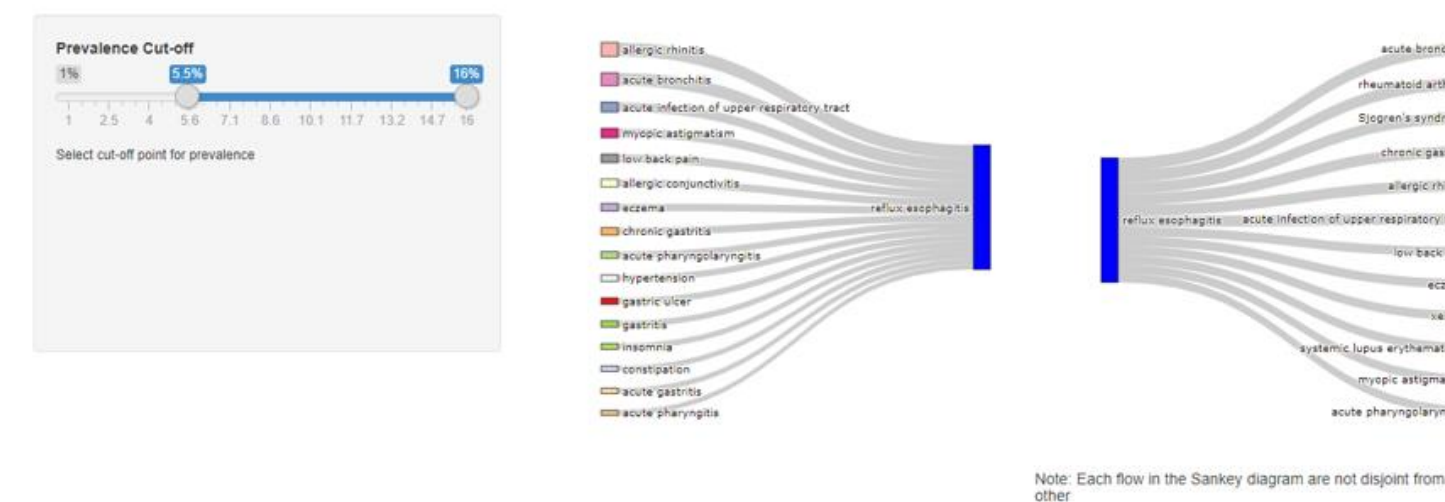
- Zoom in Reflux Esophagitis



Figure 4: Interactive digital platform of diagnosis patterns. Adjustment plate on the left enables change of prevalence cut-off to show more or less patterns. Sankey diagram on the right shows diagnoses which happen before and after reflux esophagitis.

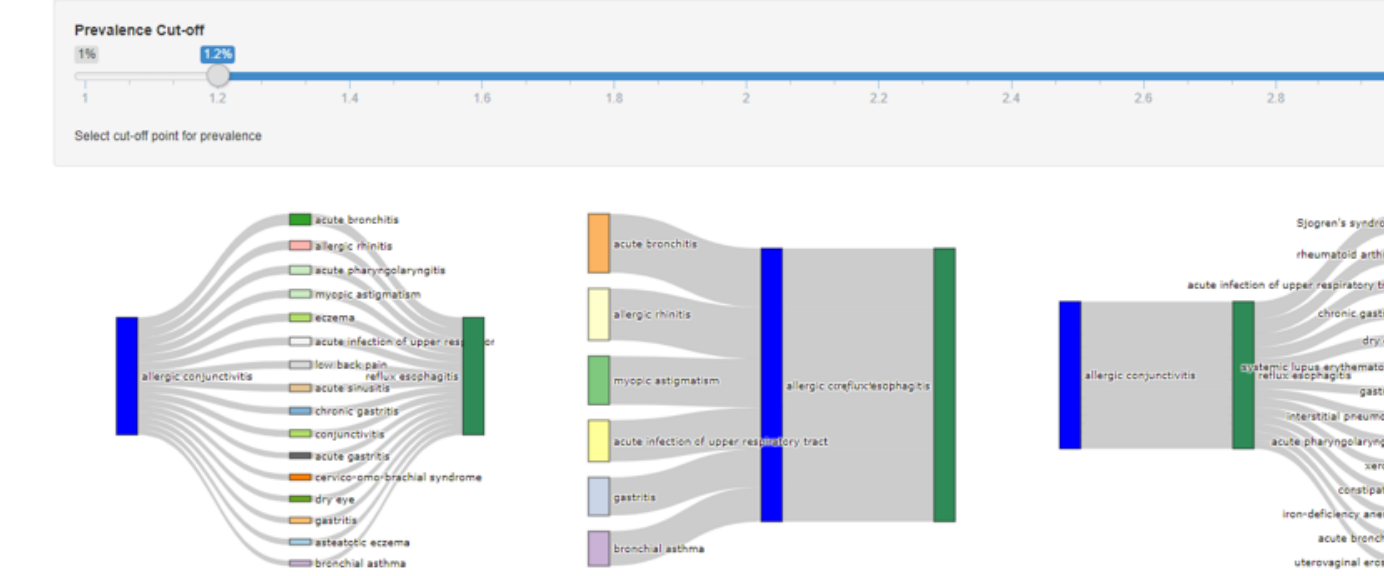- Further Zoom in allergic conjunctivitis to reflux esophagitis



Figure 5: Further zoom in includes more details but also fewer patients. Adjustment plate on the top enables change of prevalence cut-off to show more or less patterns. Sankey diagram on the bottom shows diagnoses which happen before, in-between, and after allergic conjunctivitis and reflux esophagitis.

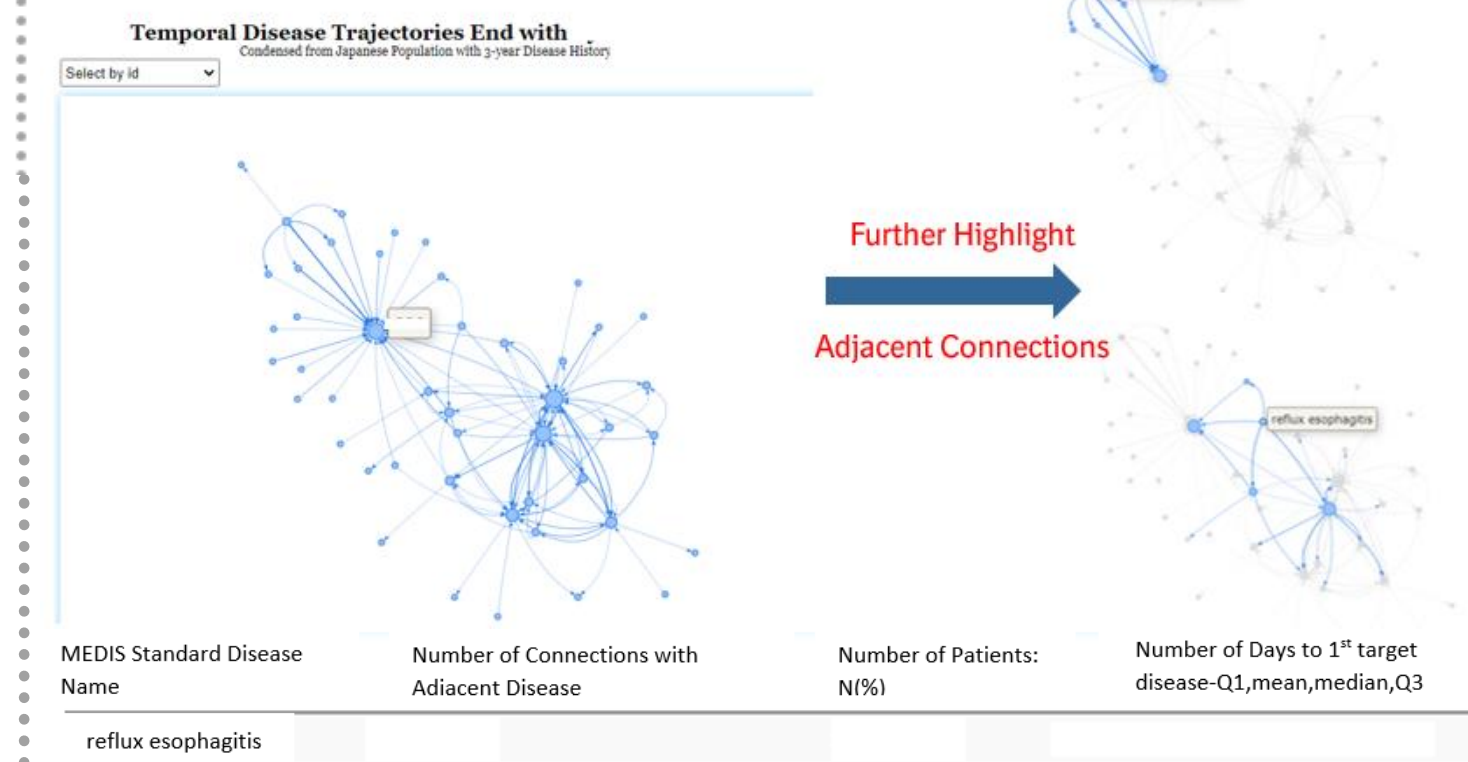### Network Analysis

- Diagnoses network before target disease



Figure 6: Interactive digital platform of diagnosis network. Table on the bottom shows detailed information of selected diagnosis, including standard disease name, number of neighbor diagnoses, prevalence and progression time needed from selected diagnosis to target disease.

## DISCUSSIONS

- The effort is exploratory
  - The results largely confirm existing clinical knowledge of the disease, e.g. common conditions leading to disease diagnosis
  - Identify patterns which may not be well understood and facilitate hypothesis generation for further research
- Data limitation needs to be well understood and communicated
  - Not all conditions are captured in the ICD codes
  - Claims data may capture all encounters during enrollment period, but EHR (especially single EHR provider data) may not
  - Lab results are usually not captured in claims data but can be found in EHR
  - Validation across different databases or different time periods may be desired

## ACKNOWLEDGEMENTS

## REFERENCES

1. Perer, Adam, and Fei Wang. "Frequence: Interactive mining and visualization of temporal frequent event sequences." Proceedings of the 19th international conference on Intelligent User Interfaces. 2014.
2. Perer, Adam, Fei Wang, and Jianying Hu. "Mining and exploring care pathways from electronic medical records with visual analytics." Journal of biomedical informatics 56 (2015): 369-378.
3. Beck, Mette K., et al. "Diagnosis trajectories of prior multi-morbidity predict sepsis mortality." Scientific reports 6.1 (2016): 1-9.