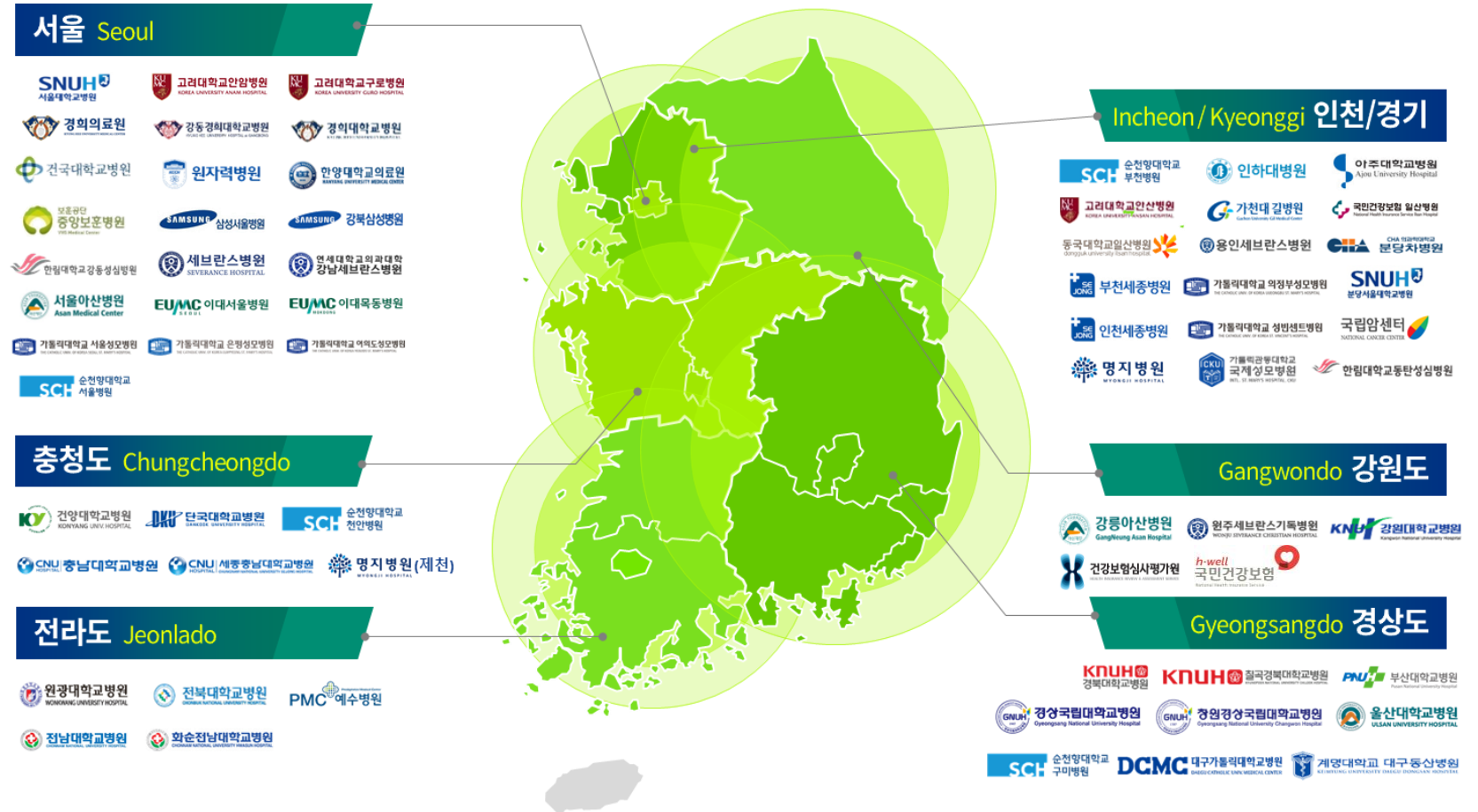scientific **data**

OPEN

ARTICLE

Check for updates

# Scalable Infrastructure Supporting Reproducible Nationwide Healthcare Data Analysis toward FAIR Stewardship

Ji-Woo Kim[1,7], Chungsoo Kim[2,7], Kyoung-Hoon Kim[3], Yujin Lee[3], Dong Han Yu[1], Jeongwon Yun[1], Hyeran Baek[1], Rae Woong Park[2,4,8]✉ & Seng Chan You[5,6,8]✉

# OMOP in Korea

## 65 Data partners

- 63 Hospitals

  (70% of tertiary hospitals)

- **2 National Institutions**

- Coordinating center

  : Ajou University

# National Health Insurance Data in Korea

**National Claims**
- From Clinics, general/tertiary hospital, pharmacy
- Including demographics, income level, diagnosis, prescription, measurement (no value of measures), procedure, device, etc
- Cost
- National Health checkup data (NHIS only)

**Data Linkage**
- Death registry (Statistics Korea)
- Cancer registry (National Cancer Center)
- Vaccine registry (KCDC)

**Health Insurance Review and Assessment Service (HIRA)**

# HIRA CDM projects

**2018 - 2020 Pilot CDM ETL
Study published in JAMA (Chan et al)**

**2020 OpenData4Covid project
World first data open for research
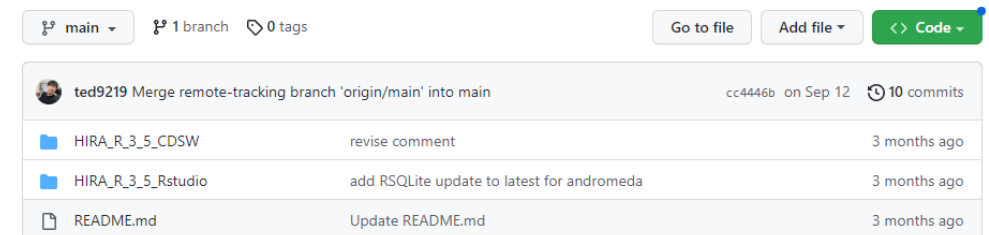(Published in many journals)**

# HIRA CDM projects

## HIRA CDM analytic environment

- ABMI supports HIRA to manage their own analytic environment for future research.
- It was not easy for implementing R into the off-line environment, but we resolved using Docker.
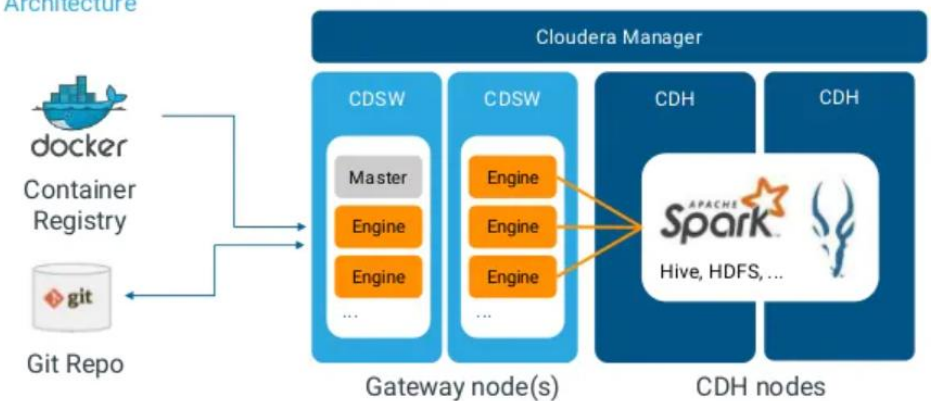
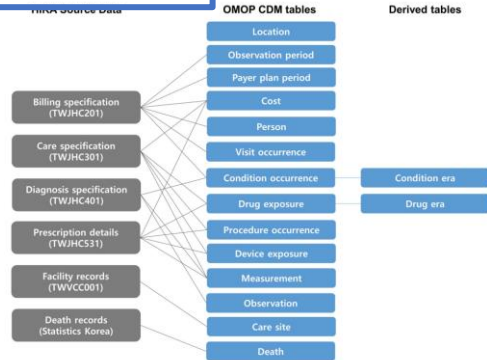## On-premise environment (manually developed)



## Off-the-shelf platform + Docker (flexible)

**Data ETL**

Fig. 3 Data mapping to OMOP-CDM from HIRA source claims database. OMOP: Observational Medical Outcome Partnership; CDM: common data model; HIRA: Health Insurance Review and Assessment Service.
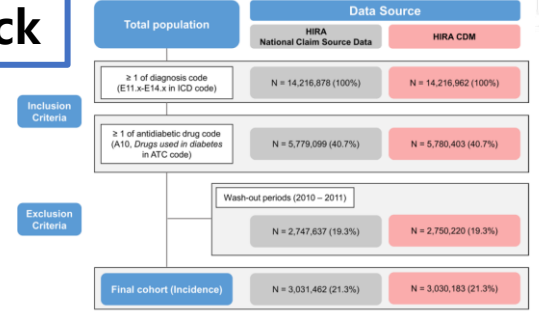
**Mapping (EDI to SNOMED)**

**Quality check**

Fig. 1 Flow chart of type 2 diabetes mellitus phenotype and comparison of incidences from the source and converted CDM databases. CDM: common data model.
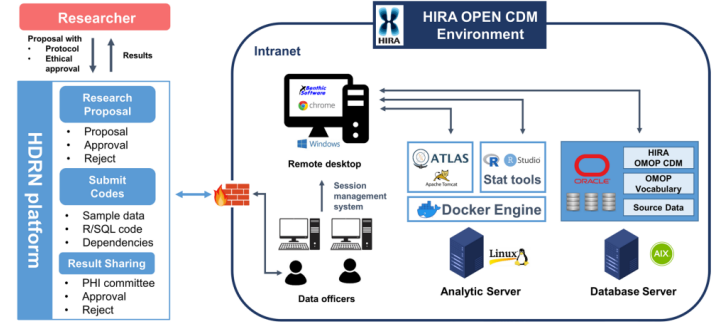
**Analytic infra**

**Open to public**

**Data catalog (ETL rule and sample data)**

Fig. 2 HIRA CDM analytic environment and data open process. Researchers can request the use of the HIRA CDM through the HDRN platform, which is an open public healthcare data platform. HDRN, Healthcare Distributed Research Network; PHI, personal health information; HIRA, Health Insurance Review and Assessment Service; CDM, common data model; OMOP, Observational Medical Outcome Partnership.

# HIRA CDM projects

## HIRA CDM Open Project

### 1st pilot opening with 10M data (2022)



**About 40 studies were done using this data**

### 2nd pilot opening with 50M data (2023)



**About 15 studies are running using this data**

**Currently, HIRA is accepting applications 1-2 times a year, but they are preparing to accept applications more frequently.**

# Thank you

**Big data department**

Ji-Woo Kim
Kyoung-Hoon Kim
Yujin Lee
Dong Han Yu
Jeongwon Yun
Hyeran Baek

Rae Woong Park
Seng Chan You