

Background

- Since March 2017 in Korea, the cancer panel test for patient with cancer or rare genetic disease started to be applied to medical insurance.
- Until recent practices in medical research field, clinical and genomic data have been harvested and stored by individual disease or institution. In a Distributed Research Network (DRN), each data holder who uses the OMOP Common Data Model can conduct a statistically powerful study without violation of private information by using the same code and integrating results of the analysis.
- We designed and present a highly scalable standard for clinical - genome convergence database.

Methods

- The International Organization for Standardization (ISO) document ISO/TS20428 defines the required/optional data fields for structured reporting the next-generation sequencing (NGS) result to clinicians or patients. We compared ISO/TS20428 document with the OMOP-CDM v5.0 to identify the relationship between their fields, and design a new database schema.
- New tables for genomic data were created in order to minimize redundancy and include the maximum information required for analysis.
- After converting both clinical and genomic data into OMOP-CDM extension model, descriptive statistics was conducted to compare variants between target and comparator cohorts. As a pilot study, we set advanced stage of lung cancer (stage IIA or more) as a target cohort and non-advanced stage of lung cancer as a comparator cohort.
- The PatientLevelPrediction package was used to predict the stage of lung cancer by adding genomic features to the clinical features through FeatureExtraction package.

Results

- In the existing OMOP-CDM v5.0, 7 clinical data, 3 health system data, and 3 derived elements tables were designated as clinical information tables to be linked to the genome information table. The table for the genomic data was expanded by creating three tables; Sequencing tables, Variant_occurrence table and Variant_annotation table (Figure 1).
- Specification and sample data of each table is uploaded at <http://bit.ly/2hA1ioj>.

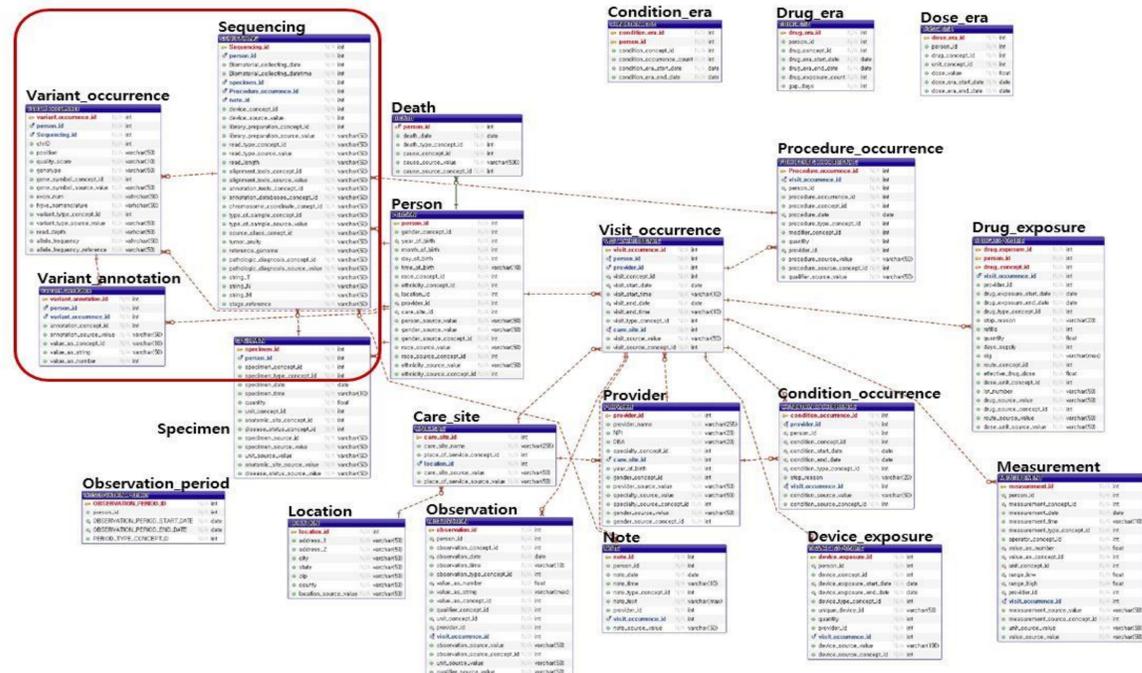


Figure 1. OMOP-CDM extension model for genomic data

- Several results of descriptive statistics of three tables in the extended genomic model are shown in Figure 2~3.

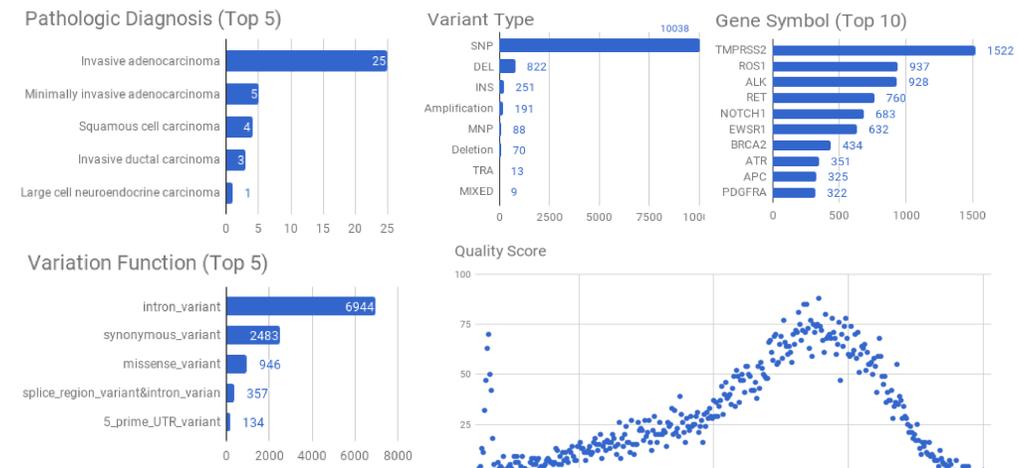


Figure 2. Sequencing and Variant_annotation Table Descriptive Statistics; Pathologic diagnosis and Variant Function

Figure 3. Variant_occurrence Table Descriptive Statistics; Variant Type, Gene Symbol (Top 10), and Quality Score

- The most different mutation occurrence rate ($p=0.002$) between lung cancer stage was PDGFRB, which is already well-known as related to growth of cancer cell (Fig. 4).

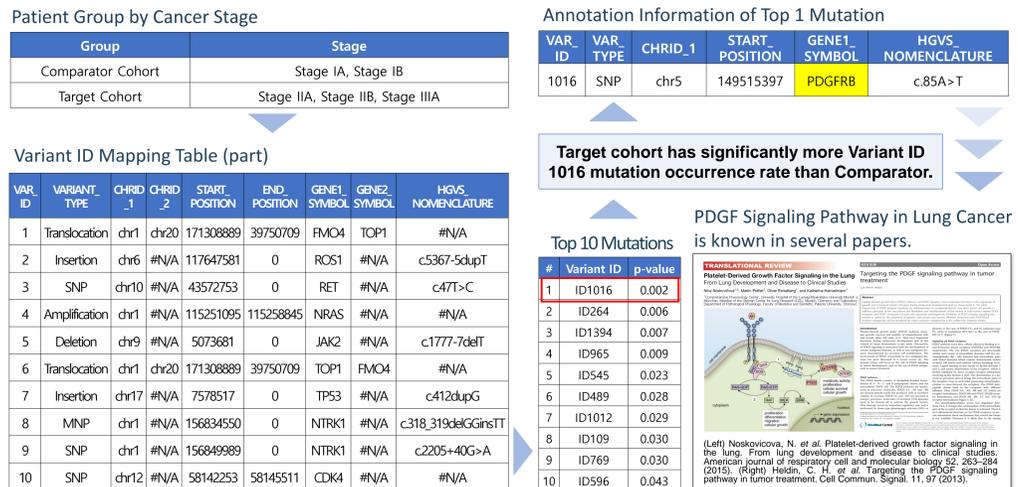


Figure 4. Calculating difference of specific mutant occurrence rate between two cohorts of lung cancer patients.

Conclusions

- We have constructed a database that can systematically manage the genetic variation data of patients who have been examined by NGS. It consists of three tables newly designed to contain genomic information.
- Through this study, we demonstrated that this database can be used not only to manage data quality but also to analyze genomic data integrated with clinical data.
- It is expected that this genomic extension model can be suggested as a prototype database standard which can be easily extended to domestic and overseas medical institutions.

- Acknowledgement : This research was supported by a grant of the Korea Health Technology R&D Project through the Korea Health Industry Development Institute (KHIDI), funded by the Ministry of Health & Welfare, Republic of Korea (grant number : H16C0992). Also, this research was supported by a grant of the Korea company SKT.
- Conflict of interest : None

- With application of PatientLevelPrediction package, predictive model was constructed corresponding to existing OHDSI research pipeline.
- For the pilot study, data of lung cancer patients who have had NGS test was used to predict whether their cancer stage is advanced or not.
- When we use both clinical and genomic data, the prediction accuracy was higher (AUROC 0.87) compared to when we use only clinical data (AUROC 0.82) (Fig. 5).

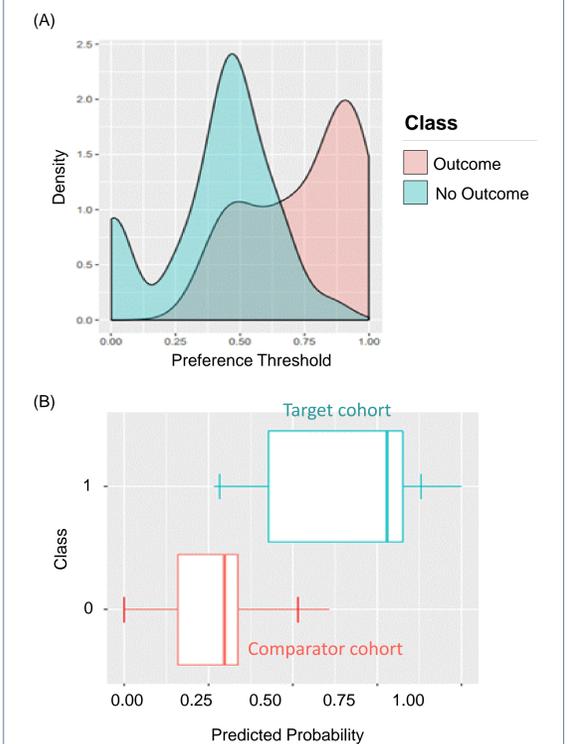


Figure 5. Predicting advanced stage lung cancer patient using both clinical and genomic data. (A) Density plot over preference threshold. (B) Box plot shows that the two class can be easily distinguished among the lung cancer patients.