# Development of a SNOMED-CT Mapping Framework for the Interoperability of Biobank Resources in Korea

Hyun Sang Park[1], Hune Cho[1], Lee Sung Hee[2] and Hwa Sun Kim[3*]

[1]*Department of Medical Informatics, Kyungpook National University,
Daegu 700-842, South Korea*
[2]*College of Nursing, Kyungpook National University,
Daegu 700-842, South Korea*
[3]*Department of Medical Information Technology, Daegu Haany University,
Gyungsan 700-715, South Korea*

## Abstract

*The existing Systematized Nomenclature of Medicine - Clinical Terms (SNOMED-CT) toolkit includes only a clinical term search, and its output functions are not linked to a particular database, making it difficult for users to map biological data, and the use of biobanks requires much time and effort. This study developed a SNOMED-CT toolkit using an optimized search, allowing users to map clinical items for efficient utilization and standardization of biobank resources. The toolkit was developed using the Java programming language and Java Database Connectivity to connect to a previously developed research database. The user interface was implemented using Swing. The functions and interfaces of the toolkit were developed after a requirement analysis by SNOMED-CT experts. Using the toolkit, the user can search SNOMED-CT concepts. Using pre-and post-coordinated mapping of clinical items with the toolkit, it is possible to maximize the semantic interoperability of the previously developed database.*

## 1. Introduction

Biobanks worldwide have collaborated to construct a Biobank Information Management System (BIMS) for the efficient management of collected biological data, to enable the development of customized methods for treating and preventing disease. Each BIMS was developed with the aim of providing high-quality information to a large number of researchers, regardless of affiliation. However, the BIMS experienced difficulties with information exchange because incompatibility and database problems between heterogeneous system, such as Electronic Medical Record System (EMRS) [1]. Typically, different institutions and practitioners use various methods to express the same concept of an item or represent it using local code or abbreviations, because BIMS operates based on the range of collected information, content, and format in each biobank. Consequently, this polymorphism of items leads to misinterpretation by researchers and acts as an obstacle to information exchange.

National projects in the United states and Europe have constructed network environments that facilitate the exchange of biological data between biobanks, such as the Cancer Biomedical Informatics Grid (caBIG) [2] and the Biobanking and Biomolecular Resources Research Infrastructure (BBMRI) [3]. Therefore, we studied database development in order to integrate biological data managed by 15 university hospital biobanks in the Republic of Korea, together with the National Biobank of Korea (NBK).

---

* Corresponding Author: Hwa Sun Kim (Daegu Haany University.) Tel: +82-53-819-1591, email:
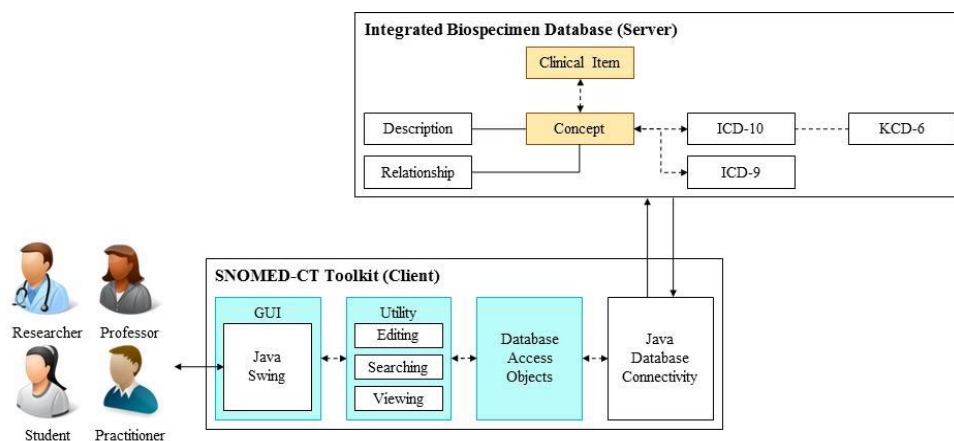    pulala@dhu.ac.kr

A three-stage data mining process (specification, classification, and standardization) was used to construct a database composed of 18 tables containing 7,197,252 raw data items. To clarify the meanings of the items, clinical items were defined using the Systematized Nomenclature of Medicine -- Clinical Terms (SNOMED-CT), and specimen items were defined using the Logical Observation Identifiers Names and Codes (LOINC). In total, 2,504 (70.4%) items were mapped and the remainder withheld because of uncertainties in the collected data.

The most efficient means of maximizing the semantic interoperability of a database would be to map the biobank data that could enhance understanding of the significance of an item to the international terminology directly. This necessitates the development of a tool that is easily maintained. However, the existing SNOMED-CT toolkit [4-6] provides only the search concept and term, outputs to the web, and includes an installation program; it is impossible to research and edit the mapped identifier or automatically map the searched concept identifier in conjunction with a specific database. Therefore, this study developed a customized toolkit that enables all users who wish to take advantage of the biological data stored in the biobank to perform clinical item searches and SNOMED-CT mapping, editing, and viewing upon accessing the database. Specifically, experts designed the function and user interface of the toolkit after conducting a requirement analysis and implemented an output function that links the terminology based on previously developed toolkit [7].

## 2. Methods

### 2.1. Expert Requirement Analysis

The user requirements were analyzed by three SNOMED-CT experts such that they could develop a toolkit enabling users to utilize the biological data. The toolkit included various functions, such as searching clinical items, mapping clinical items with concept identifiers, viewing concept information, and editing mapped information in connection with the database (Figure 1).



**Figure 1. Architecture of SNOMED-CT Toolkit**

The detailed function requirements were separated into five domains (Table 1). The search domain included the SNOMED-CT concept, reference terminology, and a clinical item search function, and the output domain included detailed information on the selected SNOMED-CT concept and a clinical item metadata output function. Specifically, we implemented an edit domain function that automatically manages and inputs the mapping results and a reference terminology domain function that utilizes the benefits of SNOMED-CT to differentiate it from existing toolkits.

**Table 2. Result of Requirement Analysis**

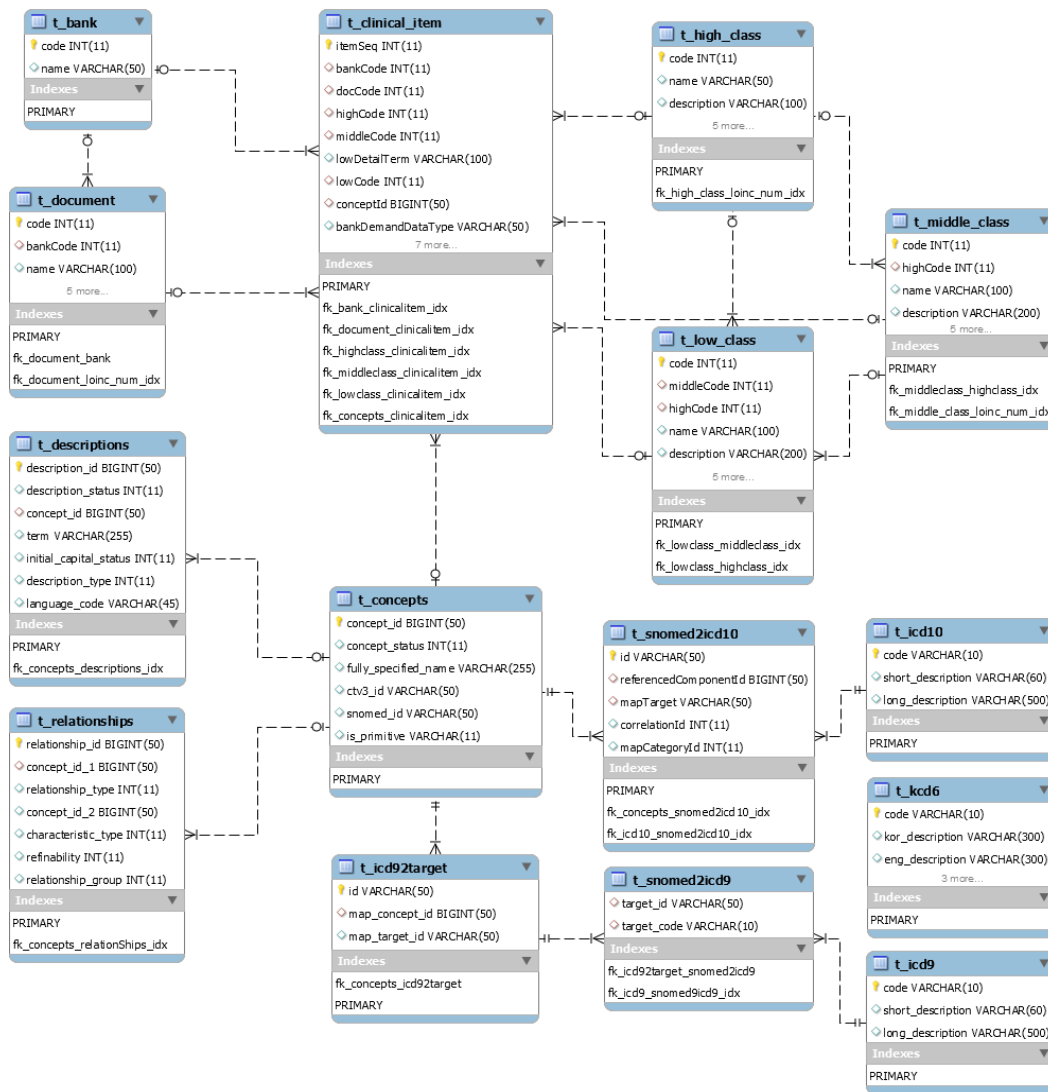| Domain | Requirement |
|---|---|
| Search | search the SNOMED-CT concept, descriptions using the identifier |
| | search the SNOMED-CT concept in accordance with string conditions |
| | search the clinical item in accordance with metadata conditions |
| | search the reference terminology concept |
| Output | output the searched SNOMED-CT concept in table form |
| | output the hierarchy and top-level information of selected SNOMED-CT concept |
| | output the detailed information of selected SNOMED-CT concept |
| | output the post-coordinated mapping attributes of selected SNOMED-CT concept |
| | output the post-coordinated expression |
| | output the clinical item metadata in table form |
| Edit | pre-coordinated mapping and automatically input |
| | post-coordinate mapping and automatically input |
| | edit the mapping information of selected clinical item |
| Reference Terminology | refer the ICD-10 code about selected SNOMED-CT concept |
| | refer the KCD-6 code about selected SNOMED-CT concept |
| | refer the ICD-9 code about selected SNOMED-CT concept |

SNOMED-CT contains 400,000 medical concepts; this vast number includes multiple terms for a single concept [8-9]. It can be mapped to other terminologies, because it consists of a hierarchy of various depths based on 19 top-level concepts. By mapping the International Classification of Diseases (ICD)-10, 9 and the Korean Classification of Diseases (KCD-6) codes corresponding to the identifier of a specific concept, it is possible to attain compatibility with existing terminology. Therefore, the user is not only provided with the detailed SNOMED-CT information for the searched term but also reference information for other terminologies.

## 2.2. Database Design

To provide clinical item information and database connection, we added nine new tables to the existing entity relationship diagram (ERD) (Figure 2). The ERD of the clinical domain consisted of seven tables, including biobank and document information tables (t_bank, t_document), clinical item classification tables (t_high_class, t_middle_class, and t_low_class), and the t_clinical_item table, which contains metadata for 1,796 clinical items.

The primary key (PK) to the classification tables is t_clinical_item, which includes data on the clinical item, values, descriptions, and SNOMED-CT concept identifier fields.

The three SNOMED-CT tables are t_concepts, t_descriptions, and t_relationships, which contain 397,787, 1,182,734, and 1,454,681 records, respectively. The unique SNOMED-CT concepts are stored in t_concepts, t_descriptions stores multiple terms for a single concept, and t_relationships stores the relationships between the source and target concepts. Therefore, t_descriptions and t_relationships refer to the PK of t_concepts. To meet these requirements, a relationship is created between the concept-id field in t_clinical_item and the PK in t_concepts. Other terminology tables (t_icd9, t_icd10, and t_kcd6) were added, and mapping tables (t_snomed2icd10, t_icd92target, and t_snomed2icd9) were placed between the t_concepts.
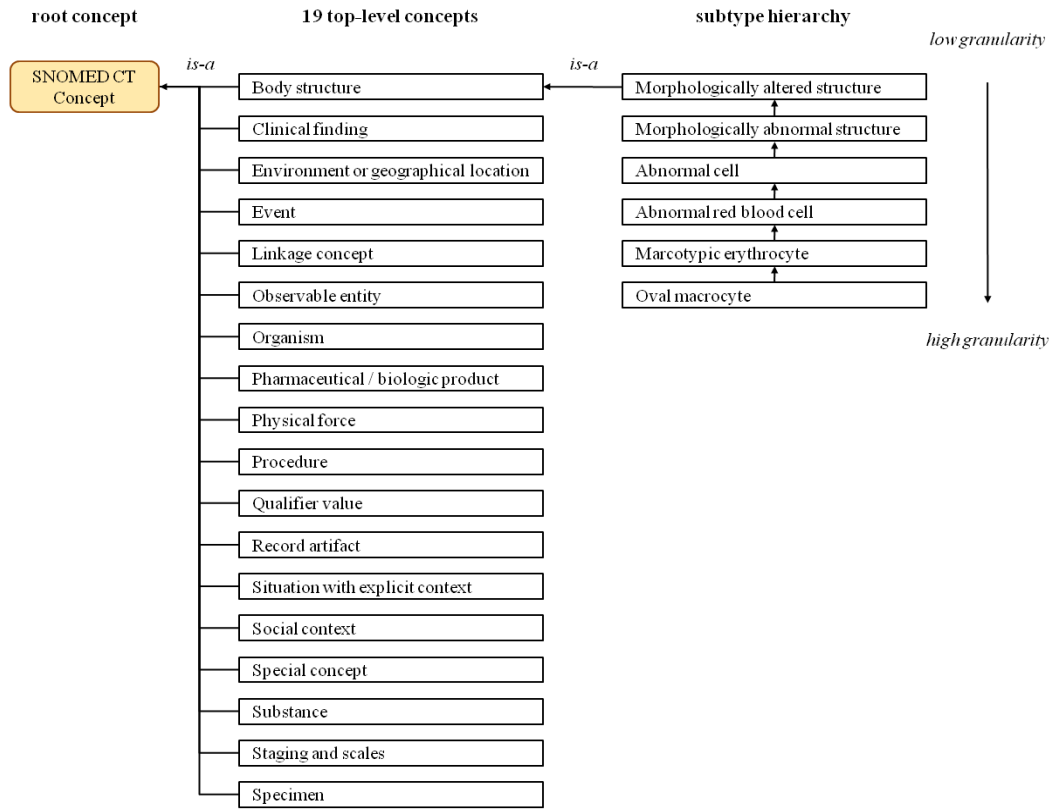
**Figure 2. Entity Relationship Diagram of SNOMED-CT Toolkit**

The mapping tables used were derived from an official source in the Unified Medical Language System (UMLS), which is updated regularly [9]. Each table was mapped to other terminologies to appropriately determine the depth of the SNOMED-CT hierarchy using predetermined rules. For example, ICD-10 was mapped to 54,262 of the 69,902 codes, which belong to three top-level concepts (51,900 clinical findings, 1,008 events, and 1,354 situations with explicit context). KCD-6 might be compatible with SNOMED-CT without the use of a mapping table by using the ICD-10 mapping table, because KCD-6 is a set of Korean diseases within the ICD-10.

## 2.3. Concept Hierarchy Algorithm

The SNOMED-CT concept consists of a hierarchy that source-target concepts using an "is-a" relationship from the "SNOMED-CT Concept" root concept. The 19 top-level concepts connected directly to the root concept with "is-a" relationships are clearly divided into those involving a semantic relationship, while the remainder are connected hierarchically with "is-a" relationships through a descendant of one or more of the top-level concepts (Figure 3). For example, "lung abscess" is the top level concept of "clinical finding"; parent concepts are "lung disorder", "inflammatory disorder of the respiratory

tract", and "thoracic abscess", and the child concepts include "abscess of lung and mediastinum", "single lung abscess", *etc*.



**Figure 3. Hierarchy of SNOMED-CT**

The SNOMED-CT concept needs to decipher the meaning of a term through its relative position in the hierarchy and detailed information, because it does not have a clear definition. However, due to the enormous number of concepts and the complex hierarchy of SNOMED-CT, a user cannot remember all of the parents or children of a specific concept. Therefore, it is necessary to display multiple parents and children in a sub-graph for a selected concept to enable rapid understanding by the user. SNOMED-CT stores a variety of relationships and attribute types for the source and target concepts in t_relationships. Any efficient method for generating a hierarchy has to explore the source and target concepts using an "is-a" relationship_type field. The toolkit was designed as a concept hierarchy algorithm that explores the parent and children of the selected concept based on "is-a" concept identifiers (Figure 4).

**Figure 4. Algorithm of SNOMED-CT Hierarchy**

The designed algorithm was divided into search areas for the parent and child concepts using two threads, and a node in the sub-tree search was performed using a recursive call method [11]. The parent concept searches the concept_id_2 field (target concept) after setting the selected concept to the concept_id_1 field (source concept). This process is repeated until the next target concept is the root concept. By contrast, the child concept is repeated until there is no next source concept after setting the selected concept to the target concept. Since the number of parent and child concepts makes the search very inefficient, we added a terminate condition that is a recursive call that continues until there is a maximum of three levels between the current concept and the parent and child concepts. In addition, if the parent and child concepts are selected, the research on the parent and child concepts has optimized the production rate of the hierarchy.

## 3. Results

This study developed a SNOMED-CT toolkit based on the designed ERD that met the cited requirements, using the Java programming language. We used the JDBC-MySQL Connector to connect to a database server developed using MySQL Server 6.0, and implemented the user interface by using SWING. The developed toolkit consist of a clinical item tab and a SNOMED-CT tab, and it is automatically executed when you log in to the previously database server using the account information that is given to each biobank users.

### 3.1. Clinical Item Tab

The clinical item tab consists of a search condition panel and a results table panel (Figure 5). The search condition panel can filter the mapping of clinical items using the combo boxes, check box, and text field. The combo boxes consist of items that document the appropriate biobank and items for the fields (high class, middle class, low class, data types, value type, units, and concept identifier) to search in conjunction with a text field. The user can filter the clinical items in a specified document using a document combo box, and the search can be divided into items not mapped and items mapped to SNOMED-CT using the combo box.
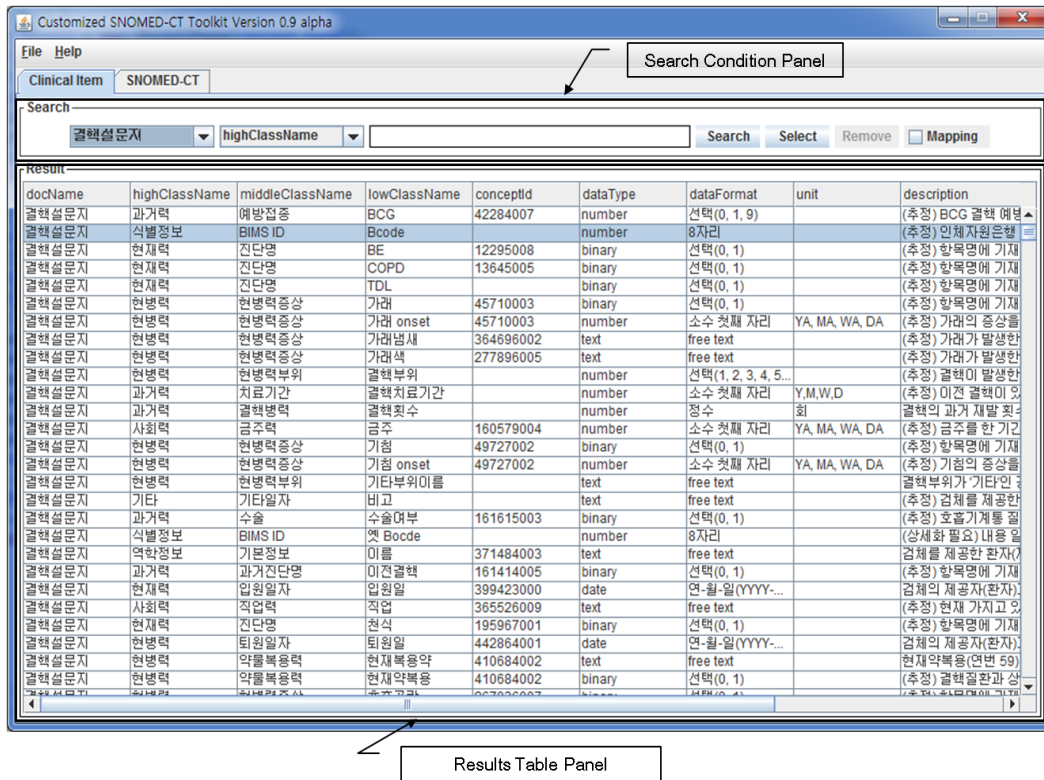
Customized SNOMED-CT Toolkit Version 0.9 alpha

File  Help

Clinical Item    SNOMED-CT

Search Condition Panel

Search

결핵설문지 ▼    highClassName ▼    [            ]    Search    Select    Remove    ☐ Mapping

Result

| docName | highClassName | middleClassName | lowClassName | conceptId | dataType | dataFormat | unit | description |
|---|---|---|---|---|---|---|---|---|
| 결핵설문지 | 과거력 | 예방접종 | BCG | 42284007 | number | 선택(0, 1, 9) | | (추정) BCG 결핵 예방 |
| 결핵설문지 | 식별정보 | BIMS ID | Bcode | | number | 8자리 | | (추정) 인체자원은행 |
| 결핵설문지 | 현재력 | 진단명 | BE | 12295008 | binary | 선택(0, 1) | | (추정)항목명에 기재 |
| 결핵설문지 | 현재력 | 진단명 | COPD | 13645005 | binary | 선택(0, 1) | | (추정)항목명에 기재 |
| 결핵설문지 | 현재력 | 진단명 | TDL | | binary | 선택(0, 1) | | (추정)항목명에 기재 |
| 결핵설문지 | 현병력 | 현병력증상 | 가래 | 45710003 | binary | 선택(0, 1) | | (추정)항목명에 기재 |
| 결핵설문지 | 현병력 | 현병력증상 | 가래 onset | 45710003 | number | 소수 첫째 자리 | YA, MA, WA, DA | (추정) 가래의 증상을 |
| 결핵설문지 | 현병력 | 현병력증상 | 가래냄새 | 364696002 | text | free text | | (추정) 가래가 발생한 |
| 결핵설문지 | 현병력 | 현병력증상 | 가래색 | 277896005 | text | free text | | (추정) 가래가 발생한 |
| 결핵설문지 | 현병력 | 현병력부위 | 결핵부위 | | number | 선택(1, 2, 3, 4, 5... | | (추정) 결핵이 발생한 |
| 결핵설문지 | 과거력 | 치료기간 | 결핵치료기간 | | number | 소수 첫째 자리 | Y,M,W,D | (추정) 이전 결핵이 있 |
| 결핵설문지 | 과거력 | 결핵병력 | 결핵횟수 | | number | 정수 | 회 | 결핵의 과거 재발 횟 |
| 결핵설문지 | 사회력 | 금주력 | 금주 | 160579004 | number | 소수 첫째 자리 | YA, MA, WA, DA | (추정) 금주를 한 기간 |
| 결핵설문지 | 현병력 | 현병력증상 | 기침 | 49727002 | binary | 선택(0, 1) | | (추정)항목명에 기재 |
| 결핵설문지 | 현병력 | 현병력증상 | 기침 onset | 49727002 | number | 소수 첫째 자리 | YA, MA, WA, DA | (추정) 기침의 증상을 |
| 결핵설문지 | 현병력 | 현병력부위 | 기타부위이름 | | text | free text | | 결핵부위가 '기타'인 |
| 결핵설문지 | 기타 | 기타일자 | 비고 | | text | free text | | (추정) 검체를 제공한 |
| 결핵설문지 | 과거력 | 수술 | 수술여부 | 161615003 | binary | 선택(0, 1) | | (추정) 호흡기계통 질 |
| 결핵설문지 | 식별정보 | BIMS ID | 옛 Bcode | | number | 8자리 | | (상세화 필요)내용 일 |
| 결핵설문지 | 역학정보 | 기본정보 | 이름 | 371484003 | text | free text | | 검체를 제공한 환자(가 |
| 결핵설문지 | 과거력 | 과거진단명 | 이전결핵 | 161414005 | binary | 선택(0, 1) | | (추정)항목명에 기재 |
| 결핵설문지 | 현재력 | 입원일자 | 입원일 | 399423000 | date | 연-월-일(YYYY-... | | 검체의 제공자(환자) |
| 결핵설문지 | 사회력 | 직업력 | 직업 | 365526009 | text | free text | | (추정) 현재 가지고 있 |
| 결핵설문지 | 현재력 | 진단명 | 천식 | 195967001 | binary | 선택(0, 1) | | (추정)항목명에 기재 |
| 결핵설문지 | 현병력 | 퇴원일자 | 퇴원일 | 442864001 | date | 연-월-일(YYYY-... | | 검체의 제공자(환자) |
| 결핵설문지 | 현병력 | 약물복용력 | 현재복용약 | 410684002 | text | free text | | 현재약복용(연번 59) |
| 결핵설문지 | 현병력 | 약물복용력 | 현재약복용 | 410684002 | binary | 선택(0, 1) | | (추정) 결핵질환과 상 |

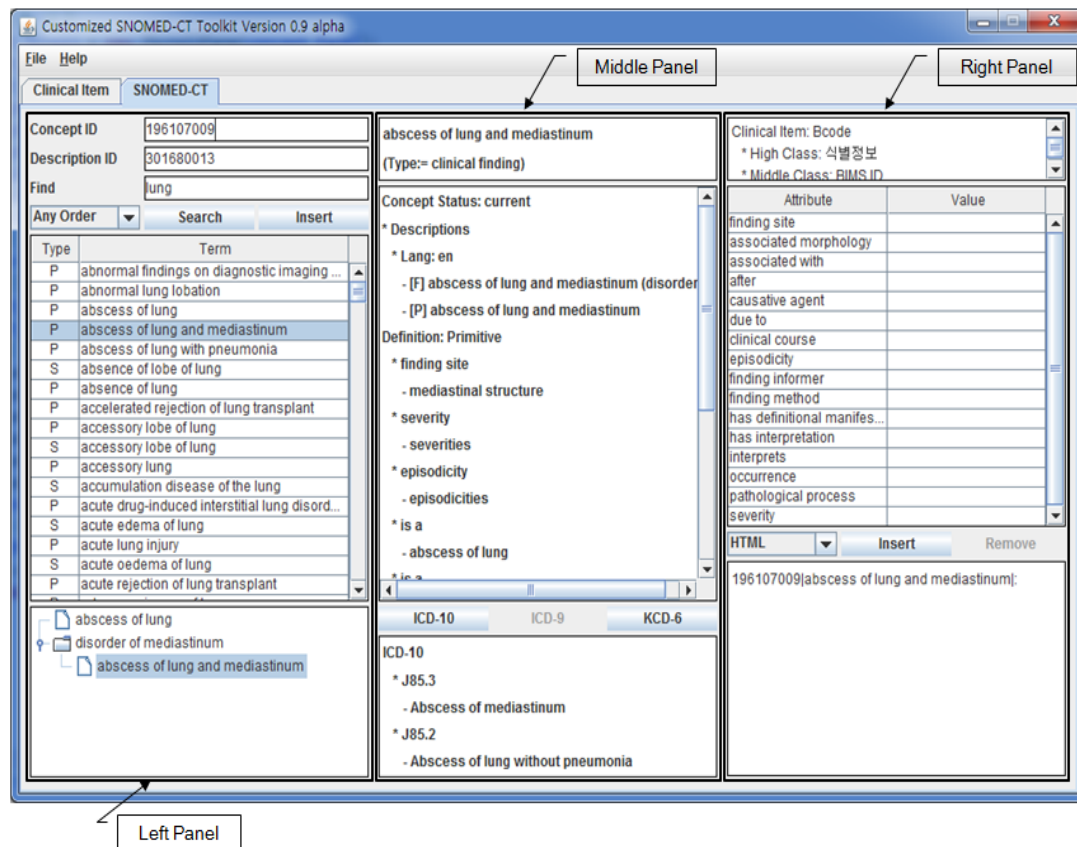Results Table Panel

**Figure 5. Screen of Clinical Item Tab**

By entering a string in the text field, it is possible to search only the corresponding clinical items. It other words, the component selected by the user is added to the conditional query statement and the database server returns the results. Specifically, pressing the deletion button removes the mapping information of the selected clinical items, and pressing the selection button transmits the current output results to the SNOMED-CT tab to be mapped automatically by the recording unit.

The results table panel outputs clinical item metadata regarding the corresponding search condition in table form. The fields in the table consist of document, high class, middle class, low class values and concept identifier, data types, value type, units, and descriptions. These can be filtered to output the result through combo boxes and a check box placed on top by using the relationship between t_clinical_item and the adjacent tables.

### 3.2. SNOMED-CT Tab

The SNOMED-CT tab consists of three detail panels, which provide the functions of analyzed requirement (Figure 6). The left panel provides three functions: *i.e.*, identifier of concept and description input, term input, search results table, and the hierarchy output function of the selected concept. An term added to the query statement with search condition in the combo box returns the results from t_descriptions, and provides a search results table.

**Figure 6. Screen of SNOMED-CT Tab**

The search function using the description identifier outputs detailed information regarding the corresponding concept by simply comparing the value of the description_id field in t_descriptions. However, the concept identifier can have a number of descriptions for a single concept, and thus the middle panel must select detailed information to output. Therefore, the toolkit was designed such that the information on the unique fully specified name (FSN) outputs descriptions for the search for the concept identifier. The combo box contains four condition items (any order, starts with, ends with, and identical), and a wildcard can be added to the term field based on selected item.

The search results table outputs information on the searched term using the term and description_type fields of t_descriptions. The types of terms are divided into preferred term and synonym, and these are output to S and P, respectively. If an outputted term exceeds the size of the table, it is outputted to the tooltip box. Specifically, the left panel manages the information output by another panel, and therefore the concept identifier of the term selected by the user should be shared by all instances. Therefore, the value of concept_id field was stored in class variables.

The hierarchy function outputs parent-child concepts with an "is-a" relationship using the t_relationships fields based on the selected concept. The hierarchy outputs the concept of multiple parents and children in the form of a JTree using the designed algorithm. The parent and child nodes that are outputted to JTree are updated using the information from the other panel using the concept identifier based on the event selected by the user.

The center panel provides three functions: the top level of the selected concept, detailed information about the concept, and information on other terminology. The top level outputs information about the corresponding concept of the 19 top-level concepts using a recursive call method as used for exploring the parent concept in the designed algorithm. Top-level information is output along with the FSN of the selected concept, which uses the fully_specified_name field of the t_concept. Detailed concept information

outputs the preferred term, synonym, FSN, parent concept, and relationship type of the selected concept by using all fields in the SNOMED-CT key tables. Other terminology is output to the code and description of the term through the activated button when the identifier of the selected concept is present in the mapping tables. Also, the user can search for other than the mapped terms by using the dialog output through the ICD-10, KCD-6, ICD-9 item in the help menu at the top (Figure 7).
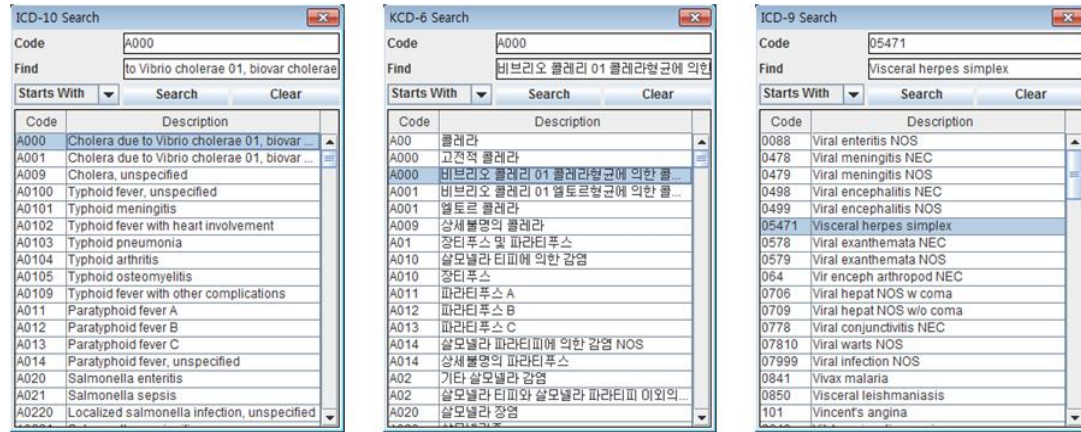


**Figure 7. Screen of ICD-10, KCD-6, ICD-9 Dialog**

The right panel provides two functions for SNOMED-CT mapping: clinical item metadata output and post-coordinated mapping. Clinical item metadata output is a function for automatically outputting information relating to the items necessary for SNOMED-CT mapping by the recording unit, which utilizes the results of the item selected in the clinical item tab. Thus, the user select one of the concept from the outputted results table through the appropriate search, the concept identifier is automatically stored in the conceptId field of t_clinical_item when the insert button pressed. Post-coordinated mapping is a method for information not represented in pre-coordinated mapping or expressing a concept in detail. It is expressed by repetition of one or more attribute concepts and the corresponding value concepts related to the selected concept.

The structure of the post-coordinated expression is of the form "concept id|concept FSN|:attribute concept id|attribute concept FSN|=value concept id|value concept FSN|".

The attribute concept can be used for only 9 of the 19 top levels (clinical finding, procedure, evaluation procedure, specimen, body structure, pharmaceutical biological product, situation with explicit, event, and physical object); "clinical finding" includes finding site, due to, severity, *etc*. These concepts output the array of attributes stored in accordance with the top-level concepts of the selected concept in table form.

## 4. Conclusions

The aim of this research was to maximize the semantic interoperability of a database developed previously that presented a customized SNOMED-CT toolkit for efficient use of biobank resources by users. The requirement analysis by SNOMED-CT experts ensured that the developed toolkit provides optimal functionality to users. Users can easily search clinical item metadata using the developed toolkit, and SNOMED-CT concept searching, pre-coordinated mapping, and post-coordinated mapping can be used for non-mapped items. Unlike existing toolkits, it is possible to confirm the reference terminology information of the selected concept by using mapping tables (ICD-10, KCD-6, ICD-9), and to search the reference terminology, by selecting the appropriate item from a menu. Therefore, users can

efficiently search and map international terminology, while maximizing their utilization of biobank data with the developed toolkit, which meets all of the predetermined requirements.

In this work, a customized toolkit was developed for SNOMED-CT as a solution to the semantic interoperability issues that must be addressed in the utilization and sharing of biobank data. The developed toolkit was designed to use only information stored in existing databases to enable optimal user functionality. Since SNOMED-CT uses international terminology that can be applied to various domains and in all medical institutions with a biobank, a plug-and-play function must be added to the developed toolkit to improve its versatility regardless of the type and structure of the database. In the future, we plan to extend the toolkit such that it can be linked directly to a user's database via a web application.
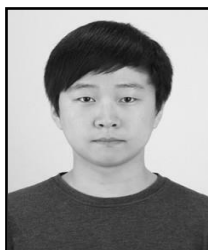
## Acknowledgments

## References

[1]   I. Hirtzlin, C. Dubreuil, N. Preaubert, J. Duchier, B. Jansen, J. Simon, P. Lobato De Faria, A. Perez-Lezaun, B. Visser, G. D. Williams and A. Thomsen, "An empirical survey on biobanking of human genetic material and data in six EU countries", European Journal of Human Genetics, vol. 11, no. 6, (2003), pp. 475-488.

[2]   S. Oster, S. Langella, S. Hastings, D. Ervin, R. Madduri, J. Phillips, T. Kurc, F. Siebenlist, P. Covitz, K. Shanbhag, I. Foster and J. Saltz, "caGrid 1.0: An Enterprise Grid Infrastructure for Biomedical Research", Journal of the American Medical Informatics Association, vol. 15, no. 2, (2008), pp. 138-149.

[3]   C. Viertler and K. Zatloukal, "Biobanking and Biomolecular Resource Research Infrastructure (BBMRI) implications for pathology", Pathologe, vol. 29, (2008), pp. 138-149.

[4]   IHTSDO SNOMED CT Browser, http://browser.ihtsdotools.org/.

[5]   NPE SNOMED CT Browser, http://www.snomedbrowser.com/.

[6]   CliniClue Xplore, http://www.cliniclue.com/.

[7]   H. S. Park, H. S. Kim and H. Cho, "Development of a Customized SNOMED-CT Toolkit for Efficient Clinical Term Searches and Mapping Biobank Resources in Korea", Advanced Science and Technology Letters Healthcare and Nursing, vol. 88, (2015), pp. 212-217.

[8]   I. Alecu, C. Bousquet and M. C. Jaulent, "A case report: using SNOMED CT for grouping Adverse Drug Reactions Terms", BMC Medical Informatics and Decision Making, (2008).

[9]   S. Lusignan, T. Chan and S. Jones, "Large complex terminologies: more coding choice, but harder to find data reflections on introduction of SNOMED CT as an NHS standard", Informatics in primary care, vol. 19, no. 3, (2011), pp. 3-5.

[10]  Unified Medical Language System, http://www.nlm.nih.gov/research/umls/.

[11]  K. Xu, C. Yuan and X. Wei, "The Hierarchical Structure and Bridging Member of k-Clique Community", International Journal of Database Theory and Application, vol. 7, no. 3, (2014), pp. 201-218.

## Authors

**Hyun Sang Park**, B. S. is a researcher and M.S. student at Department of Medical Informatics, Kyungpook National University, Daegu, South Korea. His research interests are medical informatics standard, mobile healthcare, and ISO/IEEE 11073.

**Hune Cho**, Ph. D. is a professor of the Department of Medical Informatics at Kyungpook National University School of Medicine. He has the Ph.D. degree in Medical Informatics from Utah State University. His research interests are ubiquitous healthcare, mobile healthcare, and HL7.

**Sung Hee Lee**, Ph. D. is a professor of the College of Nursing at Kyungpook National University. She is a RN and has the Ph.D. degree in Nursing Science from Kyungpook National University. Hers research interests are low fertility, fetal attachment behaviors and sexual assault.

**Hwa Sun Kim**, Ph.D. is a professor of the Department of Medical Information Technology at Daegu Haany University. She is a RN and has the Ph.D. degree in Medical Informatics from Kyungpook National University. Hers research interests are mobile healthcare, hospital information system, and standard terminology.