# Development and Mining of a Human Biomarker Database

Shaikh Farhad Hossain, Mohammad Bozlul Karim, Takematsu Shotaro, Shigehiko Kanaya, and Md. Altaf-Ul-Amin

Computational Systems Biology Lab, Graduate School of Information Science, Nara Institute of Science and Technology, 8916-5, Takayama, Ikoma, Nara 630-0192, Japan

Email: {hossain.shaikh_farhad.hr7, amin-m}@is.naist.jp

*Abstract*—**Biomarkers are the indicator of a biological state or naturally occurring molecule, gene, or characteristic whose detection indicates the presence of a disease in a living organism. Biomarkers have played an important role in medicine for improving disease diagnosis and prognosis. The biomarker is a key factor in the analysis of diseases and also for analyzing inter disease relations. There is no open source comprehensive online biomarker database. We design and develop a biomarker database based on the preceding research paper, medical report and reliable URL. This biomarker database might be helpful instrumental in identifying a new approach to treat a new disease. This biomarkers database will provide information on different diseases, protein, biochemical, metabolite that cover 5000 above diseases and biomarkers association. Information is gathered from multiple sources which include NBCI, published patents, Clinical Trials, Data from the scientific conference, PMA database, FDA, EMEA, PMDA approved documents, google scholar, PubMed and regulatory-approved documents.**

*Index Terms*—**biomarkers, NBCI, personalized medicine, PubChem ID, KNApSAcK database**

## I. INTRODUCTION

As a part of the Creative and International Competitiveness Project (CICP)-2018 of Nara Institute of Science and Technology (NAIST), Japan, literature search was conducted to find the relations between diseases and biomarkers. The diseases and biomarkers obtained from the literature search were recorded in a MySQL database. References URL, abstract and authors name were also recorded in the database. This database was submitted to CICP-2018 on November 2018, as an open source "biomarkers database".

Clear and credible biomarkers data are collected from the reliable articles so that data can be used in the diagnosis, prognosis or treatment. Collected Biomarker data were analyzed, sorted, and displayed. Source of data is hyperlinked to the references.

Biomarkers and its reference and its source are allowed to primarily store according to the following criteria:
1) We follow biomarker definition by the National Institutes of Health (NIH). NIH definition of a biomarker is: "a characteristic that is objectively measured and evaluated as an indicator of normal biologic processes, pathogenic processes, or pharmacologic responses to a therapeutic intervention."[1]
2) PubMed, Scientific Conferences, regulatory-approved documents and clinical trials are mandatory as the literature source.
3) To pass the initial review step, all collected article must be carefully gone through by the CICP-2018 scholar group.
4) Selected biomarker must have exposure, susceptibility and effect.
5) Adverse health outcome is also considered as criteria for selection (e.g. developmental disability/autism, unintentional injury, respiratory health/asthma).

literature that did not match the above mentioned criteria were discarded.

Under the Windows system, all web applications are implemented. NetBeans, Navicat and PHP languages have been used to execute the web interface. The Biomarker database is implemented using the MySQL server. We installed the XAMPP package where Apache as the HTTP service. An online database user can import biomarker data in MS Excel and PDF formats.

Disease patterns change constantly and identification of accurate and reproducible disease biomarker is also an important challenge. For finding and predicting active medicines, researchers need to read the case study, mining big data. It is tough to find the right drug, for the right patient, within the right time. This credible data will provide more accurate and comprehensive information as a time-saving and powerful tool for further biomarkers research. The biomarker database will provide details information to the users like individual biomarker information, application of biomarkers, specimen, source, authors and publication resources. This database will also provide a good environment for the researchers who are interested in protein, metabolite, disease pattern, diseases similarly, novel medicine discovery, medicine characteristic and design high-quality research. Researchers can submit their novel biomarker to our database through the submission URL page and communicate with us via chat and mail option.

## II. Biomarker Definations

Biomarkers (short for biological markers) is defined as a biochemical, cellular, or molecular alteration that is measurable [2] in biological media, such as blood (Fig 1), body fluids, tissues or cells.
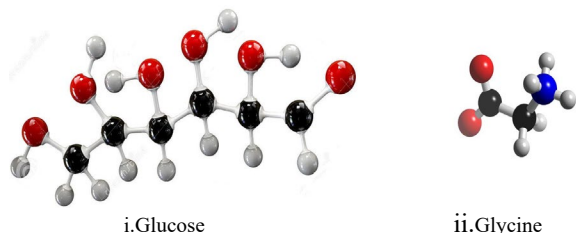


i.Glucose            ii.Glycine

Figure 1. Biomarkers of human body. i. Diabetes biomarker (Glucose - $C_6H_{12}O_6$) ii. Encephalopathy biomarker (Glycine - $C_2H_5NO_2$).

It is a sign of disease or condition which indicate normal or abnormal state of a body. Biomarkers do not include markers of prognosis only, medical or clinical test results (such as from sepsis workups, bone mineral density, x-ray analysis, or EEGs), behavioral/cognitive functioning test results, or growth or other physical measurements or observations (such as birth weight, length, fingerprint ridge count, or head circumference) also included [3]. For example, human health check includes an assessment of cholesterol, heart rate, blood pressure, triglycerides and fasting glucose levels. Body measurements such as body mass index (BMI), weight and waist-to-hip ratio are routinely used for assessing conditions such as metabolic disorders and obesity. When dealing with exposure assessment, biomarkers are generally classified into three groups:

  1) Biomarkers of exposure
  2) Biomarkers of effect
  3) Biomarkers of susceptibility

### A. Biomarkers of Exposure

Biomarkers of exposure include concentrations of the susceptibility characteristics, exogenous parent chemical, its metabolites, or changes in the body fluids or tissues (e.g., blood lead, etc.) [4]. It reflects only the absence or presence of a substance. It provides the quantification of the amount of substance present in the body and routine surveillance. Biomarkers of exposure are used extensively because they can provide information on the pathway, route and source of exposure.

### B. Biomarkers of Effect

Biomarkers of effect are the quantifiable changes (e.g., biomarkers of early loss of pregnancy, antigen production, benzo-pyrene-DNA adducts, gene suppression, tumor secretions [5] etc.), which shows an exposure to a compound and may show a resulting health effect [6]. Biomarkers measured in tumor tissue are excluded because the disease is diagnosed prior to the bio-measure, and the biomarker is used to ascertain prognosis rather than effect.

### C. Biomarkers of Susceptibility

Biomarkers of susceptibility indicate the detection of a polymorphism (such as genetic markers of cancer susceptibility) or particular genotype or a natural characteristics of an organism [7]. It may indicate the existence of or the potential for disease or a potential protection against negative health effects. It can assist to define the sensitivity of susceptible as well as critical times when exposures are more harmful.

For example, the intensity in asthmatics will indicate how susceptible that person would be to the respiratory symptom during exposure to brevotoxin, the toxic aerosolize produced during a red tide [8].

## III. Importance of Biomarker

The importance of biomarkers is radically increasing in all areas of diagnosing, monitoring disease, clinical practice and whether to predict. In every step of patient care, biomarkers are keeping the vital role. For example, blood sugar may be used to diagnose diabetes whilst HbA1c (glycosylated hemoglobin) monitors blood sugar control [9]. Biomarkers assays are becoming more and more important in clinical development. These assays can be used to understand the mechanism of action of a drug as a surrogate marker for monitoring clinical efficacy. Prognostic biomarkers are mostly identified from observational data and used to identify patients more likely to have a certain outcome. To identify a predictive biomarker, there should be a comparison system of the treatment to monitor in patients with and without the biomarker.

Biomarkers are of high importance in personalized medicine (PM) in adverse responses and detecting therapeutic in patient stratification [10]. PM is provided to the individual patient based on the disease pattern when no drug affects every patient in the same way (Fig. 2).
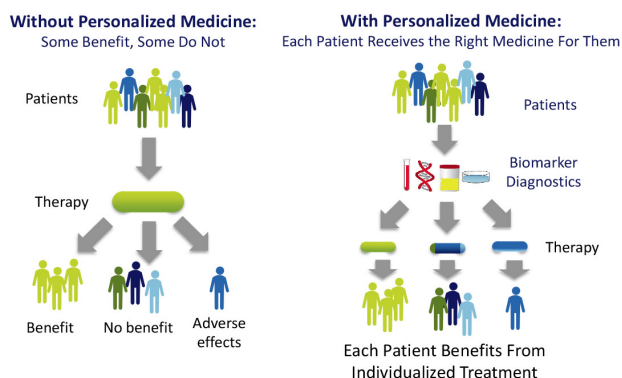


Figure 2. Personalized medicine using biomarker diagnostics [11].

According to the PM Coalition, recent five years 2009-2016, new 132 PM was in the market and 27% of new molecular entities approved by the FDA in 2016. It was a big jump from 2005, when PM accounted for just 5% of NME approvals [12].

## IV. Importance of Biomarker Database

Each day, new biomarkers are discovered all over the world by medical researchers. Biomarker data is constantly increasing and evolving but this data is stored in a scattered way. As a result, researchers need huge

time to collect literature, to read the case study, mining data, recording data, clustering data, finding similarity or dissimilarity and predicting active medicines. So, it is a challenging issue for researcher and companies around the globe. Over the last few years, the pharmaceutical industries have been facing many challenges with increasing time to market, high R & D costs, drug safety, the efficacy of treatment, patient stratification, late attrition and drug resistance [13]. In this situation, biomarkers may be the potential solution for the challenge. In every single step of drug development (Fig. 3), biomarkers are closely interconnected [14].
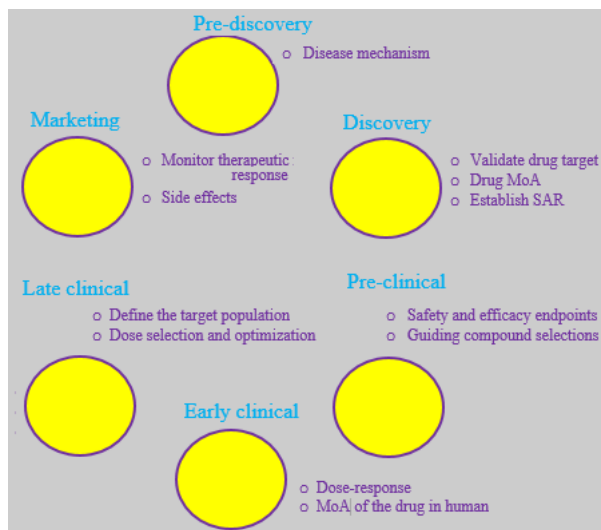


Figure 3.   Drug development chain.

In this case, our created biomarker database will play an important role because it has huge storage data that will assist the researcher to take decision efficiently. We think that our created biomarker database will create a completely new dimension (Fig. 4) to solve this problem.
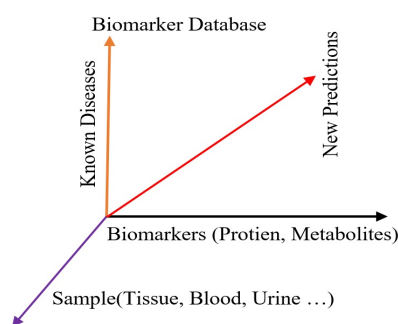


Figure 4.   Prediction of new drug ingredient using biomarker database.

In this biomarker database, data are managed within a single software application. So one can easily retrieve their expected biomarker data by running a single query and can get a large amount of information within a very short time. As data are stored in an organized way in the database so that a researcher can easily analyze the data. Using this database system, biomarker researcher or biomarker research organization can increase their efficiency and reduces overall searching time. By properly recording and updating data as a regular basis, we can address this challenging issue and this biomarker

data can be immense potential resource for researcher and medicinal company.

## V.   BIOMARKERS DATABASE

The clinically evaluated biomarker, exploratory biomarker and pre-clinical biomarker are stored in our database. The database is designed to serve three purposes:
1) to allow users to search the biomarker details information.
2) to allow users to search disease and its respective biomarkers
3) to allow users to show PubChem ID and signature

During the search procedure, user can type a search string in any text field or associated text field of a combo box to minimize the populated data in a combo box. Users that registered the screening criteria are allowed for downloading data in MS Excel and PDF formats. Also, users are allowed to view the contents of the biomarkers database without making any additions or modifications but are allowed to comment in "comment box". A short statistical description of National Center for Biotechnology Information (NCBI) classes of disease [15] and biomarker data that stored in our biomarker database is given below in Table 1.

TABLE I      SUMMARY OF KEYWORD SEARCH STRINGS

| NCBI Disease Name | Keyword Search Strings | Approx. Biomarker |
|---|---|---|
| Blood and Lymph Diseases | White, lead, fecal | 40 |
| Cancers | Asbestos,bladder, breast | 65 |
| The Digestive System | Burn, chronic | 30 |
| Ear, Nose, and Throat | Airway, nosebleed | 35 |
| Diseases of the Eye | Damage, pursuit | 18 |
| Female-Specific Diseases | Genital, neoplasm | 45 |
| Glands and Hormones | Thyroid, steroid | 60 |
| The Heart and Blood Vessels | Congestive , ischemic | 45 |
| Immune System | Autoimmune, Tolerance | 35 |
| Male-Specific Diseases | Reproductive | 18 |
| Muscle and Bone | Actin, skeletal, antigen | 15 |
| Neonatal Diseases | Urinary, gamma | 15 |
| The Nervous System | Central, neuropathy | 40 |
| Nutritional and Metabolic Diseases | Syndrome,bone, Hippocampal | 30 |
| Respiratory Diseases | Viral , Wegener | 20 |
| Skin and Connective Tissue | Melanoma, Basal | 25 |
| The Urinary System | Urothelial, obstruction | 44 |
| Mental and behavioral disorders | Developmental, Retardation, | 40 |

### A.   Salient Features of Our Biomarker Database

1.Quick and easy access

2.Online data view without any registration

3.Instant reporting and downloading for registered user

4.Interface with comprehensive search features

5.String searching

6.Excel and PDF data export options

7.Intelligence analysis

8. Comment, submission and chat option
9. Data sharing with no restriction
10. Notification system for the registered user
11. Biomarker addition into the database
12. Dedicated server and backup server
13. Routinely update with auto messaging
14. Webex session with user on-demand basics
15. Suggested clustering tools and analysis

### B. Main Menu of Biomarker Database

https://biomarkerdatabase.com.jp is the URL name of the biomarker database and user can search their expected biomarker data. User will see the default page or home page of our biomarker database. In the home page (Fig. 5), there is some menus option named Code File, Entry Information, Report. After clicking any individual menu, submenu under every individual menu option will be popped up. Like:

- **Code File**
  - *Biomarker info*
  - *Disease Info*
  - *User Info*
- **Entry Information**
  - *Biomarker Disease*
  - *Reference*
- **Report**
  - *Search Database*
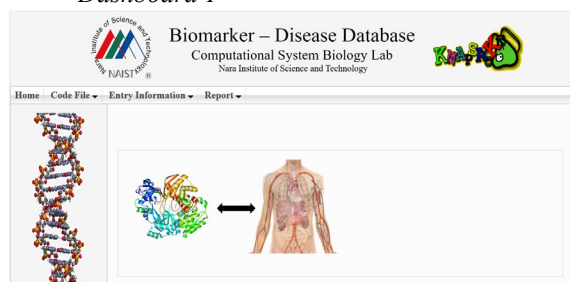  - *Dashboard 1*
  - *Dashboard 1*



Figure 5.   Home page of biomarker database.

### C. Code File

Code File menu consists of three tables called biomarker, disease, and user. In Fig. 6, by clicking the submenu "Biomarker info", a user will see 10 biomarkers per page and clicking "next" button, the user will be able to see more 10 biomarkers and so on. The number of biomarker pages are 307, a user is able to see the whole biomarker of our database. Moreover, by writing full or partial biomarker name or biomarker PubChem ID or both in the specific combo box and clicking the search button a user can specify search terms for any or all of the following fields using advanced search features. By clicking reference URL, a user will get more details information of a specific biomarker.

Similarly, in Fig. 7, by clicking the submenu "Disease Info", the user will see 10 diseases per page and clicking "next" button, the user will be able to see more 10 diseases and so on. The number of disease pages are 137,

a user is able to see the whole disease information of our database. Moreover, by writing full or partial disease name or disease alias name or disease description a user can specify search terms for any or all of the following fields using advanced search features. Using $2^3$-1 or 7 combinations of fields, a user can search at the same time.



Figure 6.   Basic table of biomarker of biomarker database.



Figure 7.   Basic table of disease of biomarker database.



Figure 8.   User request form of biomarker database.

And "User Info" (Fig. 8) is the third submenu of Code File, a user should fill-up the fields of an individual or organization basic information. There are three types of user 1. admin 2. data entry operator 3.web user. Based on data entry operator and web user request, the admin will provide the different level of access privileges. There are many reasons to fill out a registration form of users (e.g. Communications with the user, feedback from the user, support, type of the user). Users are requested to send valuable suggestions and comments as well as any

queries about our database and web services. To access web application, users may use any browser but release after 2015 (recommended) with JavaScript enabled.

### D. Entry Information

Entry Information menu consists of two tables called Biomarker Disease and Reference. In Fig. 9, by clicking the submenu "Biomarker Disease", a user will see 10 biomarkers in the first page and its respective diseases. By clicking "next" button, the user will able to see biomarkers and its respective diseases of next page and so on. The total number of Biomarker Disease page is 464. If a user writes full or partial (minimum 2 characters) disease name or biomarker name or both in the respective left text box, then right combo box will automatically show some options to choose by filtering data using the text box. After selecting the option, clicking the "search" button will take a little bit time to retrieve result from database. Thus the search engine will show biomarker name with its respective diseases and references.
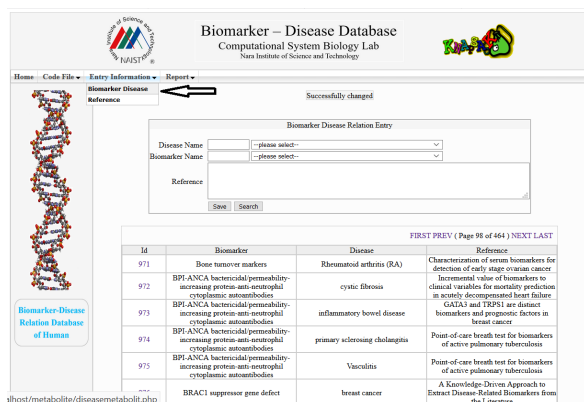


Figure 9.   Biomarker disease searching from biomarker database.

### E. Report

Report is the last and more important menu containing three submenu options named search database, dashboard 1 and dashboard 2 (Fig. 10). Submenu dashboard 1 and dashboard 2 are not yet completed. By writing full or partial biomarker name or PubChem ID or disease name or alias name or reference or any combinations ($2^5$-1 or 31) of fields can be searched at the same time.
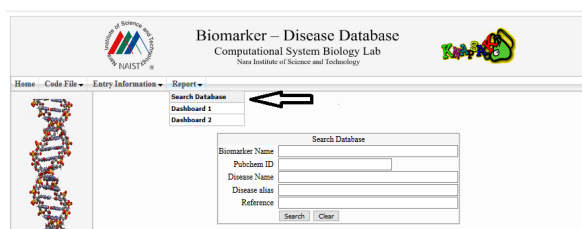


Figure 10.  Searching biomarker or disease by typing string.

For example, if a user can remind a partial name of a biomarker i.e. Cy (full name Cyclin) he can easily search on entire data by using the partial mnemonic. If the user writes in Biomarker Name field only "cy" and in Disease Name field "cancer", it is enough for the user to search his expected outcome. Our database will provide specific

and details information that string is containing "cy" in biomarker name field and related to cancer disease.

In Fig. 11, selecting the designed query result, the online user can import biomarker data in MS Excel and PDF formats. For downloading or importing data, the user needs to register in our website. In short, register users are allowed to download or import data.
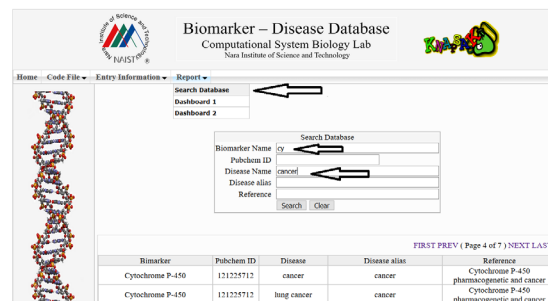


Figure 11. Searching result and import data.

In our biomarker webpage, there is a hyperlink of KNApSAcK database [16] at the top right KNApSAcK icon (Fig. 12). KNApSAcK database is a comprehensive species-metabolite relationship database developed by Computational Systems Biology Lab, NAIST, Japan. Our developed biomarker database is also a part of this lab. By clicking the icon, the user can access
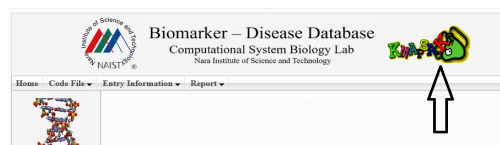


Figure 12. Species-metabolite database hyperlink icon.

The species-metabolite database (Fig. 13). Right now, this database is only accessible for NAIST authenticate student. Later, it will be opened for all.



Figure 13. Species-metabolite KNApSAcK database.

## VI.  Limitations

Dashboard1 and dashboard 2 are not yet completed. We are trying to the linkage of biomarkers statistical tools for analysis.

## VII.  Conclusion

In our database, more than 5000 biomarker-disease associations information is included. Specific, accurate and sensitive data published in various journals and scientific conferences are collected by CICP scholar group. Recorded biomarker data are reliable for multipurpose usage like monitoring disease progression,

diagnosis and prognosis. All of the biomarker information sources are linked to the valid references. A user friendly web interface is designed for data extraction by writing a string or query. Using different search combination, user is able to extract meaningful data. Entire data is stored in a single platform using a relational database. User can import biomarker data with different formats like Excel or PDF.

## CONFLICTS OF INTEREST

All other authors declare no conflicts of interest.

## AUTHOR CONTRIBUTIONS

Shaikh Farhad Hossain and Takematsu Shotaro collect the biomarkers data, Mohammad Bozlul Karim develops the website, Shigehiko Kanaya and Md. Altaf-Ul-Amin verifies the data. All authors had approved the final version.

## ACKNOWLEDGMENT

## REFERENCES

[1] K. Strimbu and A. T. Jorge, "What are biomarkers?" *Current opinion in HIV and AIDS*, vol. 5, no. 6, pp. 463-466, 2010.

[2] L. E. Walker, *et al.*, "WONOEP appraisal: Molecular and cellular biomarkers for epilepsy," vol. 57, no. 9, pp. 1354-1362, 2016.

[3] L. J. Orchinik, H. G. Taylor, and K. A. Espy, "Cognitive outcomes for extremely preterm/extremely low birth weight children in kindergarten," *Int Neuropsychol Soc.*, vol. 17, no. 6, pp. 1067-1079, 2011.

[4] W. P. Watson and A. Mutti, "Role of biomarkers in monitoring exposures to chemicals: Present position, future prospects," *Biomarkers*, vol. 9, no. 3, pp. 211-242, 2004.

[5] H. Hamouchene, V. M. Arlt, I. Giddings, and D. H. Phillips, "Influence of cell cycle on responses of MCF-7 cells to benzo[a]pyrene," *BMC Genomics*, vol. 12, p. 333, 2011.

[6] G. F. Nordberg, *et al.*, "Prevalence of kidney dysfunction in humans-relationship to cadmium dose, metallothionein, immunological and metabolic factors," *Biochimie*, vol. 91, no. 10, pp. 1282-1285, 2009.

[7] I. Iavicoli, V. Leso, and P. A. Schulte, "Biomarkers of susceptibility: State of the art and implications for occupational exposure to engineered nanomaterials," *Toxicol Appl Pharmacol*, vol. 299, pp. 112-124, 2015.

[8] L. E. Fleming, B. Kirkpatrick, and L. C. Backer, "Aerosolized red-tide toxins (brevetoxins) and asthma," *Chest*, vol. 131, no. 1, pp. 187-194, 2007.

[9] M. Leow, "Glycated hemoglobin (HbA1c): Clinical applications of a mathematical concept," *Acta Inform Med*, vol. 24, no. 4, pp. 233-238, 2016.

[10] R. La Russa, *et al.*, "Personalized medicine and adverse drug reactions: The experience of an Italian teaching hospital," *Current Pharmaceutical Biotechnology*, vol. 18, no. 3, pp. 274-281, 2017.

[11] *Personal Medicines in Development Chartpack / A New Treatment Paradigm / A New Treatment Paradigm.* (2018). [Online]. Available: https://chartpack.phrma.org/personal-medicines-in-development-chartpack/a-new-treatment-paradigm/a-new-treatment-paradigm

[12] *Personalized Medicine an Infographic*. (2018). [Online]. Available: https://invivo.pharmaintelligence.informa.com/IV005059/Personalized-Medicine-An-Infographic

[13] N. J. D. Gower, *et al.*, "Drug discovery in ophthalmology: Past success, present challenges, and future opportunities," *BMC Ophthalmol*, vol. 16, no. 11, 2016.

[14] A. B. Halim, "Biomarkers in drug development: A useful tool but discrepant results may have a major impact," *Intechopen*, 2011.

[15] NCBI. (2018). [Online]. Available: https://www.ncbi.nlm.nih.gov/

[16] Y. Nakamura, *et al*. (2018). *KNApSAcK: A Comprehensive Species-Metabolite Relationship Database*. [Online]. Available: http://kanaya.naist.jp/KNApSAcK/

**Shaikh Farhad Hossain** received his M.Sc. degree in Information and Communication Technology from Bangladesh University of Engineering and Technology (BUET), Dhaka, Bangladesh, in 2017. He is currently pursuing his PhD in Biomarkers, Proteins and Gene expression in Nara Institute of Science and Technology (NAIST), Nara, Japan.

**Mohammad Bozlul Karim** received his M.Engg. degree in Information and Communication Technology from Bangladesh University of Engineering and Technology (BUET), Dhaka, Bangladesh, in 2010. He is currently pursuing his PhD in Graph Clustering and its application on different Biological Network analysis in Nara Institute of Science and Technology (NAIST), Nara, Japan.

**Shigehiko Kanaya** received B.Sc. degree in Bio-science from Science University of Tokyo, Japan in 1985, and Ph.D. from Toyohashi University of Technology, Japan in 1990. He served as an Assistant Professor in Information Engineering at Yamagata Univ. in 1990, Guest Associate Professor at National Institute of Genetics in 1996, Associate Professor at Electronic and Information Engineering in 1999, Associate Professor, Applied Bio system Engineering at Yamagata Univ. in 2000, Guest researcher at Bio radical Institute (Yamagata Prefecture), in 2000 Associate Professor, Research and Education Center for Genetic Information at NAIST in 2001, Associate Professor, Graduate School of Information Science at NAIST in 2002, Professor, Graduate School of Information Science at NAIST in 2004 and is currently working as a Professor, Nara Institute of Science and Technology, Japan.

**Md. Altaf-Ul-Amin** received B.Sc. degree in Electrical and Electronic Engineering from Bangladesh University of Engineering and Technology (BUET), Dhaka, M.Sc. degree in Electrical, Electronic and Systems Engineering from University Keban gsaan Malaysia (UKM) and PhD degree from Nara Institute of Science and Technology (NAIST), Japan. He received the best student paper award in the IEEE 10th Asian Test Symposium. Also, he received other best paper awards as a co-author of journal and conference articles. He previously worked in several universities in Bangladesh, Malaysia and Japan. Currently he is working as an associate professor in Computational Systems Biology Lab of NAIST.