# A trial of web-based method to teach impact of genome variation on protein structures at high school

*Saki Katagiri[1], Mayu Shibata[1], Misae Kudo[1] and Kei Yura[1,2]

The advancement of genome sequencing technology has unveiled whole DNA sequences of individuals. Comparison of the genome sequences has uncovered differences among human population which may relate to the differences in the susceptibility to diseases of each person. The wisdom acquired by the state-of-the-art technology should be shared in the society, especially among the next generation to advance their quality of life. We developed an easy-to-use web interface to promote the idea of genome sequence variation, effect on protein structures and the relationship between the variations and diseases in humans. We used the web tool in the Super Science High School Program class and assessed the impact of using the tools for the promotion of understanding among high school pupils.

## 1. Introduction

Education on cancer was proposed to be included in a curriculum for compulsory school in Japan in 2012 and the report for the methodology was announced in 2016 [1]. The report emphasized the importance of the accurate knowledge of cancer, its prevention and early discovery with cure. The knowledge should be disseminated among younger generation to pre-empt or at least reduce the future risk of diseases, thereby reducing the medical care cost of the society.

The report proposed that the cancer education be introduced in health and physical education (H.P.E.) [1]. It is natural to include cancer education in H.P.E., because a lifelong health care has already been taught in H.P.E. However, cancer is a complicated phenomenon and its understanding requires advanced biology including genetics, immunology, and so forth, hence cancer education is strongly related to biology classes, too. In the context of biology, cancer can be understood as changes in gene expression levels, in gene expression patterns, in protein structure and function and in other related phenomena, which are caused by the changes in genome sequences. The changes in genome sequences may occur by environmental or hereditary factors. The thorough understanding of these processes, in effect, facilitates the civil understanding of cancer.

The realization of the ideal education in high school is a tall order. Long hour working of high school teachers is considered a social problem in Japan. In 2019, more than half of high school teachers were reported to work overtime. The working hours added up to more than 45 hours per month [2]. Decreasing the workload of high school teachers is a priority to keep a sound education, or introducing a new task in the current curriculum is virtually impossible.

---

* SK, MS and MK contributed equally to this study.

We, therefore, propose a web-based ready-to-use educational material to let the pupils understand the relationship between change in genome sequences and changes in protein structures by themselves. The recent worldwide pandemic of COVID-19 emphasized the need of online learning materials in high school, too. Although in the advanced bioinformatics study, there already are many web tools regarding this topic such as VaProS [3], most of the web tools are too complicated for high school pupils. Here we report a new learning tool called Genome Literacy in Education (GLE), available at http://cib.cf.ocha.ac.jp/gle/. The web tool offers an opportunity to learn the impact of genetic changes in the context of molecular biology and can help the pupils to understand that the critical role of various genetic conditions in genetic diseases. We used GLE in a seminar of the Super Science High School Program in the summer of 2021 and assessed the usability of the tool.

## 2. Material and Methods

*Acquisition of human gene and protein sequence data*

In order to focus on proteins familiar to high school pupils, we selected ten human proteins introduced in a high school biology textbook. Protein coding sequences (CDSs) in DNA which code for the amino acid sequences of proteins were retrieved from Ensembl database (ver. 104) [4]. Positions of variation on DNA sequence of each protein were extracted from ClinVar, a database for variations of human genome [5]. Among all the obtained variations on the selected CDSs, we arbitrarily chose ten variations on each CDS for simplicity. The three-dimensional (3D) structures of the proteins were obtained from PDBj [6]. The 3D structure data in PDBj contain multiple protein chains and ligands. The minimum set of protein chains and ligands which was sufficient to interpret the function of the protein were used in the web tool.

*Matching data from different data sources*

Correspondence of DNA sequence data obtained from Ensembl, Variation data from ClinVar, and protein 3D structure from PDBj was made on amino acid sequences. A CDS in Ensembl was conceptually translated to amino acid sequence. A variation position from ClinVar was mapped onto the CDS and the nucleotides was altered. The altered CDS was translated to an amino acid sequence and a variant sequence was obtained. Each variation on CDS was classified into one of the following four categories based on the changes in amino acid sequences; synonymous variation (no change in amino acid sequence), nonsense variation (the codon changes to stop codon), missense variation (changes in amino acid residue) and frameshift variation (insertion or deletion in CDS). All of the frameshift variation resulted in forming a stop codon close to the variation site and hence the coding sequence was shortened. The obtained amino acid sequences were aligned to the sequences extracted from PDBj. The 3D structure data of the protein was not necessarily complete, but some of the parts were often missing, and hence the correspondence should be carefully checked.

*Prediction of protein 3D structures*

Based on the correspondence of the amino acid residues between the sequence and the 3D structure, the 3D structures of original and variation proteins were built. For the 3D structure without variation, if the original sequence and the sequence in PDBj were identical, then the 3D structure in PDBj was used. If there were one amino acid residue discrepancy, the region around the discrepancy site was modeled by

MODELLER (ver. 10.1) [7]. For modeling, the discrepancy site ±3 residues at least were modeled and positions of other residues were fixed. If there were more than one amino acid discrepancy, the whole structures were modeled using ModWeb [7].

The protein structure modeling software can deduce the structure of the protein, yet the stability of the protein should be assessed independently. The structural stability of the variant proteins was assessed by PremPS [8] before modeling the structure. A variant sequence with $\Delta\Delta G \geq 1.0$ kcal/mol, where $\Delta\Delta G$ is a difference of unfolding Gibbs free energy between the original and the variant proteins, was considered not to be stable enough to fold. Out of missense variations, sequences that passed the criterion above were modeled by MODELLER [7] on the original 3D structure. The quality of deduced 3D structure was further evaluated based on the structural change from the template structure. If the backbone root mean square deviation was less than 1.0Å, the variation was considered to cause no big structural change. For nonsense and frameshift variations, the quality of the structure was assessed by ProSA [9]. ProSA checks the reliability of the structure by checking the compatibility of the modeled structure with the experimentally solved protein 3D structures in the public database.

*Data integration and presentation in Genome Literacy in Education*

The data prepared above were integrated in a single web tool named Genome Literacy in Education (GLE) as shown in Figure 1. Each protein was presented with three types of biological information: the nucleotide sequence that codes for the protein, the amino acid sequence, and the 3D structure in a single window. A variation can be selected by a pull-down menu on the nucleotide sequence viewer. When a user chooses a variation, an alignment between the variant and the original amino acid sequences is shown inside the amino acid sequence box. The variant and original structures are displayed in the structure box. The 3D structure of the protein is drawn by NGL viewer (ver. 0.10.4-1) [10]. The protein is shown in ball-and-stick model colored in standard atom color, and a backbone trace is emphasized in ribbon model in cobalt green. The backbone of the altered amino acid residues is colored in pink. The atoms of the altered amino acid residues are shown in ball-and-stick model with greater radius. Molecules that coexisted in the protein structure database are copied as is and are displayed in color different from green and pink.
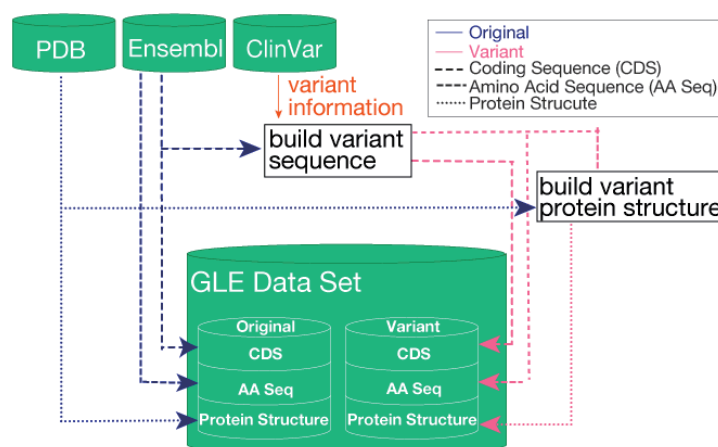


**Figure 1.** The architecture of Genome Literacy in Education (GLE), a web tool to visualize the effect

of variation on genome in protein structure.

## 3. Results and Discussion

*Proteins in GLE*

The proteins selected for GLE are shown in Table 1. The crystal structure of collagen α-3(IV) was a chimeric protein of three different proteins including the transcript of COLA43. This protein chain was divided into three independent chains. A structure of keratin, fibrinogen β chain and thyroxine-binding globulin included one amino acid residue site that was occupied by different amino acid type from the one in translated sequence, and hence the original structure was modeled. For actin and rhodopsin proteins, since there were no 3D structures of the proteins derived from *Homo sapiens*, the structures were modeled based on the structure derived from different organisms. The 3D structure data in PDBj rarely had the coordinates of hydrogen atoms and therefore, atoms shown in GLE were limited to non-hydrogen atoms.

**Table 1.** Genetic and structural data used in GLE

| Protein | Gene | Interacting Molecule | Ref. |
|---|---|---|---|
| Actin (N.A.) | ACTA1 (1578.1) | - | 11 |
| Breast cancer type 1 susceptibility protein (1JM7) | BRCA1 (11453.1) | BRCA1-associated RING domain protein 1, $Zn^{2+}$ | 12 |
| Catalase (1F4J) | CAT (7891.1) | Protoporphyrin IX | 13 |
| Collagen α-3(IV) (6WKU) | COL4A3 (42829.1) | Collagen α-4(IV), Collagen α-5(IV) | 14 |
| Fibrinogen β chain (3E1I) | FGB (3786.1) | Fibrinogen α chain, Fibrinogen γ chain, $Ca^{2+}$ | 15 |
| Haemoglobin α subunit (4HHB) | HBA1 (10399.1) | Haemoglobin β subunit, Protoporphyrin IX | 16 |
| Insulin (1GUJ) | INS (7729.1) | - | 17 |
| Keratin, type I cytoskeletal 10 (6UUI) | KRT10 (11377.1) | Keratin type II cytoskeletal 1 | 18 |
| Rhodopsin (N.A.) | RHO (3063.1) | - | 19 |
| Thyroxine-binding globulin (2XN6) | SERPINA7 (14518.1) | 3,5,3',5'-tetraiolo-L-thyronine | 20 |

Parenthesis in Protein refers to ID of PDBj and that in Gene to consensus coding sequence code in NCBI. N.A. indicates that there is no corresponding ID.

*GLE interface*

GLE was developed for the Japanese high school pupils to learn the relationship between the genome variation and the change in protein structure/function (Figure 2). A user can select one of the ten proteins that are familiar to Japanese high school pupils on the top menu, and observe the nucleic acid sequence, amino acid sequence and 3D structure of those proteins. A user can change nucleotides in the given

sequence at the top of the interface, and observe the changes in amino acid sequence and protein structure at the bottom of the interface.

The number of variations displayed in GLE was limited to ten in each protein. In ClinVar, there were a huge number of variations, even in the coding region of the genome in each protein. For example, in human actin proteins, there were more than 800 variations including the ones on homologous proteins [21]. When all these variations were compiled, the display became too complicated and the web tool was unlikely to meet the goal. Therefore, we randomly selected the variations and limited the number to ten to assure an easy use of the interface to high school pupils.

*Evaluating GLE*

The usefulness of GLE was tested in an online seminar for high school pupils in 2021. We used GLE for the explanation of the relationship between genome variation and protein structure changes. GLE apparently helped pupils understand the relationship through the central dogma of molecular biology. Throughout the seminar, 14 out of 15 pupils mentioned that using GLE made their learning much easier compared to the situation without interactive online learning material. This result suggests the usefulness of GLE to promote understanding of the central dogma and the impact of changes in genome sequences which may result in hereditary diseases.



**Figure 2.** A snapshot of GLE. Actin protein is examined on the tool.

*The use of GLE for cancer education*

GLE can be one of the supporting materials for cancer education in high school. GLE best depictures the connection between the changes in nucleotides and amino acid residues. BRCA1 in GLE is one of the genes strongly associated with hereditary breast and ovarian cancer. Variations of nucleotides on BRCA1 gene may cause cancer [22]. Most of the variations stored in GLE for BRCA1 were annotated as pathogenic variations in ClinVar [5] with only one exception which is a variation found in healthy population. Selection of the first variation closest to the 5' terminus of the DNA sequence immediately tells that the amino acid sequence is truncated and half of the 3D structure of the protein is lost. This process of analysis intuitively makes the pupil recognize that something unusual may happen by the nucleotide variation. Since the variation is strongly correlated with breast-ovarian cancer patients, the pupil can imagine that the unusual phenomena caused by the variation results in the disease. The causality between the variation and the disease cannot be taught only with GLE, but the hypothesis can be nurtured in pupils which will be a starting point to learn the detail of molecular biology of cancer.

## 4. Conclusion

We built GLE, a web tool that visualizes the genetic variations and their impacts on protein structures for high school pupils. GLE included ten different proteins which are familiar to them. For each protein, a set of biological information, namely, coding sequence, amino acid sequence and protein 3D structure was shown with at most ten variations. The usefulness of GLE was confirmed in the online seminar. We showed that the web tool can facilitate teacher-friendly cancer education class focusing on its molecular basics. We will have a roadmap to increase the number of proteins in GLE, once we have positive feedback from high school pupils and teachers.

## 5. Acknowledgments

## References

1. Task force for considering the cancer education at Ministry of Education, Culture, Sports, Science

and Technology. Report on the methods for cancer education at School (2016). https://www.gankyouiku.mext.go.jp/download/cancer_education_report.pdf2.

2. Ministry of Education, Culture, Sports, Science and Technology. Results of Survey on Efforts to Reform School Work Styles by the Board of Education in 2020 (2020). https://www.mext.go.jp/content/ 20201224-mxt_zaimu-000011455_1.pdf

3. Gojobori, T. Ikeo, K., Katayama, Y., Kawabata, T., Kinjo, AR., Kinoshita, K., Kwon, Y., Migita, O., Mizutani, H., Muraoka, M., Nagata, K., Omori, S., Sugawara, H., Yamada, D., Yura, K. VaProS: A Database-Integration Approach for Protein/Genome Information Retrieval. *Journal of Structural and Functional Genomics,* **17**, 69-81 (2016). https://doi.org/10.1007/s10969-016-9211-3

4. Howe, K.L., Achuthan, P., Allen, J., Allen, J., Alvarez-Jarreta, Amode, J.M.R., Armean, I.M., Azov, A.G., Bennett, R., Bhai, J., Billis, K., Boddu, S., Charkhchi, M., Cummins, C., Fioretto, L.D.R., Davidson, C., Dodiya. K., Houdaigui, B.E., Fatima. R., Gall, A., Giron, C.G., Grego, T., Guijarro-Clarke, C., Haggerty, L., Hemrom, A., Hourlier, T., Izuogu, O.G., Juettemann, T., Kaikala, V., Kay, M., Lavidas, I., Le, T., Lemos, D., Martinez, J.G., Marugán, J.C., McMahon, T.M.A.C., Mohanan, S., Moore, B., Muffato, M, Oheh, D.N., Paraschas, D., Parker, A., Parton, A., Prosovetskaia, I., Sakthivel, M.P., Salam, A.I.A., Schmitt, B.M., Schuilenburg, H., Sheppard, D., Steed, E., Szpak, M., Szuba, M., Taylor, K., Thormann, A., Threadgold, G., Walts, B., Winterbottom, A., Chakiachvili, M., Chaubal, A., Silva, N.D., Flint, B., Frankish, A., Hunt, S.E., IIsley. G.R., Langridge, N., Loveland, J.E., Martin, F.J., Mudge, J.M., Morales, J., Perry, E., Ruffier, M., Tate, J., Thybert, D., Trevanion, S.J., Cunningham, F., Yates, A.D., Zerbino, D.R., Flicek, P. Ensembl 2021. *Nucleic Acids Research*, **49**, 884–891 (2021). https://doi.org/10.1093/nar/gkaa942

5. Landrum, M. J., Chitipiralla, S., Brown, G. R., Chen, C., Gu, B., Hart, J., Hoffman, D., Jang, W., Kaur, K., Liu, C., Lyoshin, V., Maddipatla, Z., Maiti, R., Mitchell, J., O'Leary, N., Riley, G. R., Shi, W., Zhou, G., Schneider, V., Maglott, D., Holmes, J.B., Kattman, B. L. ClinVar: improvements to accessing data. *Nucleic Acids Research,* **48**, D835-D844 (2020). https://doi.org/10.1093/nar/gkz972

6. Bekker, G-J., Yokochi, M., Suzuki, H., Ikegawa, Y., Iwata, T., Kudo, T., Yura, K., Fujiwara, T., Kawabata, T., Kurisu, G. Protein Data Bank Japan: Celebrating our 20th anniversary during a global pandemic as the Asian hub of 3D macromolecular structural data. *Protein Science*, **31**, 173-186 (2022). https://doi.org/10.1002/pro.4211

7. Sali, A., Blundell, T.L. Comparative protein modelling by satisfaction of spatial restraints. *Journal of Molecular Biology*, **234**, 779-815 (1993). https://doi.org/10.1006/jmbi.1993.1626, https://modbase.compbio.ucsf.edu/modweb/

8. Chen, Y., Lu, H., Zhang, N., Chen, Y., Zhu, Z., Wang, S., Li, M. PremPS: Predicting the impact of missense mutations on protein stability. *PLoS Computational Biology,* **16**, e1008543 (2020). https://doi.org/ 10.1371/journal.pcbi.1008543

9. Wiederstein, M., Sippl, M.J. (2007) ProSA-web: interactive web service for the recognition of errors in three-dimensional structures of proteins. *Nucleic Acids Research* **35**, W407-W410 (2007). https://doi.org/10.1093/nar/gkm290, https://prosa.services.came.sbg.ac.at/prosa.php

10. Rose, A.S., Bradley, A.R., Valasatava, Y., Duarte, J.M., Prlic, A., Rose, P.W. NGL viewer: web-based molecular graphics for large complexes. *Bioinformatics*. **34**, 3755-3758 (2018). https://doi.org/10.1093/bioinformatics/bty419

11. Dominguez, R., Madasu, Y., Rao, J.N. Actin cytoskeleton. Mechanism of actin filament pointed-end capping by tropomodulin. *Science*, **345**, 463-467 (2014). https://doi.org/10.1126/science.1256159

12. Brzovic, P.S., Hoyt, D.W., King, M.C., Klevit, R.E., Rajagopal, P. Structure of a BRCA1-BARD1 heterodimeric RING-RING complex. *Nature. Structural & Molecular Biology*, **8**, 833-837 (2001). https://doi.org/10.1038/nsb1001-833

13. Abraham, D.J., Ko, T.P., Musayev, F.N., Safo, M.K., Wu, S.H. Structure of tetragonal crystals of human erythrocyte catalase. *Acta Crystallography Section D*, **57**, 1-7 (2001). https://doi.org/10.1107/s0907444900013767

14. Bauer, R., Boudko, S.P., Chetyrkin, S.V., Hudson, B.G., Ivanov, S., Smith, J., Voziyan, P.A. Collagen IV$^{\alpha345}$ dysfunction in glomerular basement membrane diseases. II. Crystal structure of the α 345 hexamer. *Journal of Biological Chemistry*, **296**, 100591-100591 (2021). https://doi.org/10.1016/j.jbc.2021.100591

15. Bowley, S.R., Lord, S.T. Fibrinogen variant BbetaD432A has normal polymerization but does not bind knob "B". *Blood*, **113**, 4425-4430 (2009). https://doi.org/10.1182/blood-2008-09-178178

16. Fermi, G., Fourme, R., Perutz, M.F., Shaanan, B. The crystal structure of human deoxyhaemoglobin at 1.74 A resolution. *Journal of Molecular Biology*, **175**, 159-174, (1984). https://doi.org/ 10.1016/0022-2836(84)90472-8

17. Brange, J., Chance, K., Dodson, G.G., Finch, J., Scott, D.J., Whittingham, J.L., Wilson, A. Insulin at pH 2: Structural analysis of the conditions promoting insulin fibre formation. *Journal of Molecular Biology*, **318**, 479-490, (2002). https://doi.org/10.1016/S0022-2836(02)00021-9

18. Bunick, C.G., Milstone, L.M. The X-Ray crystal structure of the keratin 1-keratin 10 helix 2B heterodimer reveals molecular surface properties and biochemical insights into human skin disease. *Journal of Investigative Dermatology*, **137**, 142-150 (2017). https://doi.org/10.1016/j.jid.2016.08.018

19. Adaixo, R., Dawson, R.J., Deupi, X., Flock, T., Maeda, S., Marino, J., Matile, H., Mohammed, I., Muehle, J., Pamula, F., Schertler, G., Stahlberg, H., Taylor, N.M., Tsai, C.J. Cryo-EM structure of the rhodopsin-Gαi-βγ complex reveals binding of the rhodopsin C-terminal tail to the Gβ subunit. *elife*, 8, e46041 (2019). https://doi.org/10.7554/eLife.46041

20. Carrell, R.W., Chan, W.L., Ley, S.V., Loiseau, F., Milroy, L.G., Myers, R.M., Qi, X., Read, R.J.,

Wei, Z., Yan, Y.,Zhou, A. Allosteric modulation of hormone release from thyroxine and corticosteroid binding-globulins. *Journal of Biological Chemistry*, **286**, 16163-16173 (2011). https://doi.org/10.1074/jbc.M110.171082

21. Duong, H.T.T., Suzuki, H., Katagiri, S., Shibata, M., Arai, M., Yura, K. Computational study of the impact of nucleotide variations on highly conserved proteins: In the case of actin. *Biophysics and Physicobiology,* **19**, 3190025 (2022). https://doi.org/10.2142/biophysico.bppb-v19.0025

22. Momozawa, Y., Sasai, R., Usui, Y., Shiraishi, K., Iwasaki, Y., Taniyama, Y. Parsons, M.T., Mizukami, K., Sekine, Y., Hirata, M., Kamatani, Y., Endo, M., Inai, C., Takata, S., Ito, H., Kohno, T., Matsuda, K., Nakamura, S., Sugano, K., Yoshida, T., Nakagawa, H., Matsuo. K., Murakami, Y., Spurdle, A.B., Kubo, M. Expansion of Cancer Risk Profile for BRCA1 and BRCA2 Pathogenic Variants. *JAMA Oncology*, **8**, 871‑878 (2022). https://10.1001/jamaoncol.2022.0476

Saki Katagiri

Address: Graduate School of Humanities and Sciences, Ochanomizu University, 2-1-1 Otsuka, Bunkyo,
Tokyo 112-8610, Japan

E-mail: g2170501@edu.cc.ocha.ac.jp


Mayu Shibata

Address: Graduate School of Humanities and Sciences, Ochanomizu University, 2-1-1 Otsuka, Bunkyo,
Tokyo 112-8610, Japan

E-mail: g2170503@edu.cc.ocha.ac.jp


Misae Kudo

Address: Graduate School of Humanities and Sciences, Ochanomizu University, 2-1-1 Otsuka, Bunkyo,
Tokyo 112-8610, Japan

E-mail: g2140547@edu.cc.ocha.ac.jp


Kei Yura

Address: Graduate School of Humanities and Sciences, Ochanomizu University, 2-1-1 Otsuka, Bunkyo,
Tokyo 112-8610, Japan. Center for Interdisciplinary AI and Data Science, Ochanomizu
University, 2-1-1 Otsuka, Bunkyo, Tokyo 112-8610, Japan.

E-mail: yura.kei@ocha.ac.jp