

Down Syndrome Prediction/Screening Model Based on Deep Learning and Illumina Genotyping Array

Bing Feng

School of Computer Science and Technology
Tianjin University
Tianjin, China
Department of Computer Science and Engineering
University of South Carolina
Columbia, SC, USA
e-mail: bingf@email.sc.edu

William Hoskins

Department of Computer Science and Engineering
University of South Carolina
Columbia, SC, USA
e-mail: hoskinsw@email.sc.edu

Yan Zhang

Department of Computer Science and Engineering
University of South Carolina
Columbia, SC, USA
e-mail: zhang348@email.sc.edu

Zibo Meng

Department of Computer Science and Engineering
University of South Carolina
Columbia, SC, USA
e-mail: mengz@email.sc.edu

David C. Samuels

Vanderbilt University School of Medicine
Vanderbilt University
Nashville, TN, USA
e-mail: David.c.samuels@vanderbilt.edu

Yan Guo*

Department of Internal Medicine
The University of New Mexico
Albuquerque, NM USA
e-mail: YaGuo@salud.unm.edu

Jijun Tang*

School of Computer Science and Technology
Tianjin University
Tianjin, China
Department of Computer Science and Engineering
University of South Carolina
Columbia, SC, USA
e-mail: jtang@cse.sc.edu

Abstract—Down syndrome (DS) is a genetic disorder with genome dosage imbalances and micro-duplications of human chromosome 21. It is usually associated with a group of serious diseases, including intellectual disabilities, cardiac diseases, physical abnormalities, and other abnormalities. Currently, since there is no cure for human DS, screening and early detection have become the most efficient way for DS prevention. In this study, we used deep learning techniques to build accurate DS prediction/screening models based on the analysis of newly introduced Illumina genotyping array. Specifically, we built chromosome SNP maps based on clinical genotyping data collected by Vanderbilt University Medical Center. Then we proposed a convolutional neural network (CNN) architecture with ten layers and two merged CNN models, which took two input chromosome SNP maps in combination. Our CNN DS prediction/screening model achieved over 99.3% average accuracy, as well as very low false positive and false negative rate, which are critical to disease prediction and screening in medical practice. It also had better performances in terms of all evaluating metrics when compared with three conventional machine-learning algorithms. Finally, we visualized the feature maps and the trained filter weights from intermediate layers of our trained

CNN model. We further discussed the advantages of our method and the underlying reasons for its robust performance.

Keywords- Deep Learning, Convolutional Neural Network, Down Syndrome, Illumina Genotyping

I. INTRODUCTION

Down syndrome (DS) is a genetic disorder caused by genome dosage imbalances and micro-duplications of human chromosome 21 (HSA21) [1]. It is usually associated with intellectual disabilities, congenital heart defects, childhood Leukaemia, Alzheimer's disease, early ageing, physical abnormalities, and other abnormalities [1–3]. Even though DS occurs in a high rate worldwide (1 per 1,000 live births) [4], and it has been well studied, researchers haven't found any effective cure method yet [5]. No environmental factor or parents' behavioral factor has been discovered to cause the Human DS either [6]. Currently, Human DS therapy studies are mainly concentrating on early intervention, educational therapy [7,8], physical therapy [9,10], as well as emotional and behavioral therapies [8,11]. These therapies only have limited effects and they never cure DS

fundamentally [5]. Therefore, screening and early detection have become the most efficient way for human DS prevention.

DS screening has been studied since the 1960s. A few DS biomarkers have been discovered, such as alpha-fetoprotein levels, human chorion gonadotropin, and unconjugated estriol [12, 13]. Currently, DS screening studies include ultrasound measurement of fetal nuchal translucency [14], blood test [15], sequencing test [16], and combined genetic test [17]. However, 1 out of every 16 DS screening test positive women are still suffering from further high-risk invasive diagnostic procedures, such as amniocentesis and chorionic villus sampling, which have a 1/200 chance to result in fetal loss [16,18].

Human chromosome 21(Hsa21) encodes more than 500 genes [19,20], including protein modifiers, transcription factors, RNA splicing factors/modifiers, cell surface receptors and adhesion molecules, and various biochemical pathways components [20,21]. However, only 165 genes are annotated as protein-coding genes or microRNAs. More than 350 genes have unassigned functions [22]. Recent GWAS studies have discovered that the Single-nucleotide polymorphisms (SNPs) variations, copy number variations, and unidentified genetic variations are highly associated with the genetic disorders of human DS genomes [23–26]. Illumina has introduced a new exome genotyping array technique that targets rare SNPs of human genome. Vanderbilt University Medical Center and Vanderbilt Epidemiology Center have developed chip-processing protocols for processing Illumina genotyping array and collected the clinical data for various diseases [24].

Machine-learning techniques have been widely used in disease prediction, disease diagnosis, and bio-marker identification [27–29]. Commonly used machine-learning algorithms include conventional artificial neural network, deep learning neural networks, support vector machine (SVM) [27], random forest [29], decision tree [30], and Bayesian classifiers [31]. Recently, deep convolutional neural networks (CNN) and recurrent neural networks (RNN) have been successfully demonstrated in a variety of disease screening, diagnosis, and predictions problems [27,32–34]. To the best of our knowledge, only a few conventional machine-learning techniques have been applied to DS studies [20]. Most of them were conducted on mice DS models [20,24,35]. Clara et al. analyzed expression levels of 77 proteins and designed a self-organizing map to identify biological differences in DS mice model Ts65Dn [20]. Cao et al. used a naive bayes classifier to predict the level of locomotor activity under the treatments of the N-methyl-D-aspartate receptor in mice models Ts65Dn and Ts1Cje [35]. Zhao et al. proposed a hierarchical constrained local model with independent component analysis to detect DS of young pediatric patients [36]

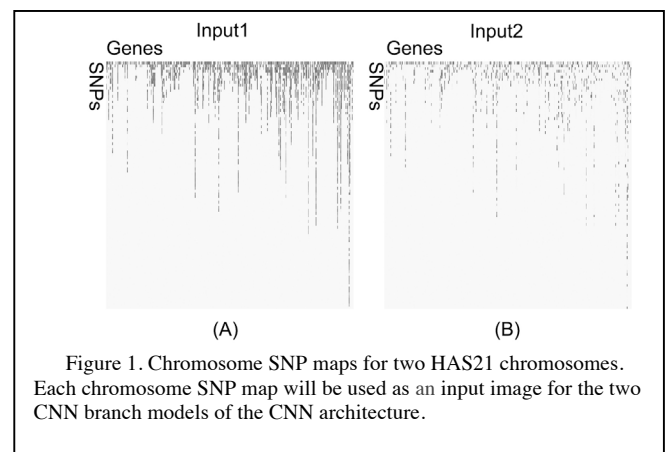
In this study, we used deep learning techniques to construct an accurate human DS prediction/screening model from newly introduced Illumina exome genotyping array data. First, we built the chromosome SNP maps based on the clinical genotyping array data collected by Vanderbilt University Medical Center. Next, we proposed convolutional

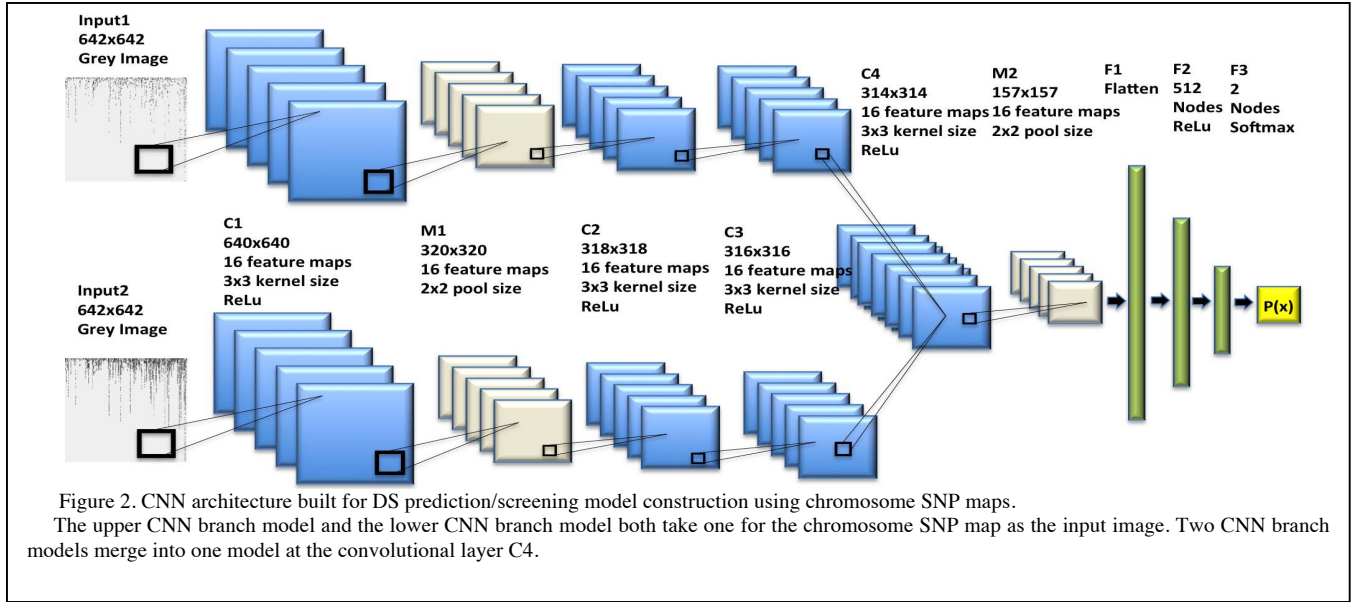
neural network (CNN) architecture with nine layers and two merged branch CNN models, which took two chromosome SNP maps as inputs. We also used three conventional supervised learning algorithms SVM, random forest, and decision tree to construct DS prediction/screening models with the same genotyping data. We further evaluated and compared the performances of different models, and concluded that the CNN model achieved the best performances in all evaluation metrics over all models. Finally, we visualized the feature maps and the trained filter weights from intermediate layers of our trained CNN model. We further discussed the advantages of our method and the underlying reasons for its powerful performance.

II. METHODS

A. Data

We analyzed an Illumina exome genotyping array dataset that targeted at the rare single-nucleotide polymorphisms (SNPs). This genotyping dataset was collected by Vanderbilt University Medical Center containing 378 samples, 315 non-DS patients and 63 DS patients. Each sample contained intensity information of total 5458 SNPs from 321 HSA21 coding genes. For the CNN model construction, we built two chromosomes SNP maps and used them as input images for each sample. As shown in Figure 1, each chromosome SNP map represents intensities of all SNP site on one single chromosome of HSA21. Each pixel of chromosome SNP map represents the intensity of one SNP site from one gene. Figure 2 gives two input chromosome SNP maps for two HSA21 chromosomes. Columns of the chromosome SNP map represent all adjacent genes on the chromosome. Rows represent adjacent SNP sites within the same gene. For conventional machine-learning model constructions, we used the original Illumina genotyping array dataset with 5458 SNPs features for model training and testing. All SNPs data were normalized into the interval [0,1]. To construct robust and reliable models, we generated ten parallel training and testing datasets by randomly select first 75% samples for training and the rest for testing. Then we built ten parallel DS prediction/screening models and calculated average performance metrics to provide a systematic evaluation of all models used in this study.





B. CNN Architecture Construction

The CNN architecture constructed in this study was merged from two branch CNN models. Each branch model contained five layers, including one input layer, three convolutional layers, and one max pooling layer. Each branch model took one chromosome SNP map as the input image. We further merged the two branches CNN model into one model in layer 6, which was another convolutional layer. Figure 2 showed the structures of our CNN architecture. The model was compiled with binary cross-entropy loss function and stochastic gradient descent optimizer (SGD) with learning rate of 0.01, 1e-6 decay, and 0.9 nesterov momentum. CNN model construction and training were implemented with Keras and Tensorflow as the backend, using a NVIDIA GeForce GTX TITAN Pascal 12GB GPU on Ubuntu 14.04.5 LTS.

C. Conventional Machine-learning Algorithms.

The random forest, SVM, and decision tree algorithms were implemented by Python and Scikit Learn package [37]. Random forest algorithm used Gini impurity and “entropy” to measure the quality of split for a feature. There was no limit on the maximum depth of each sub-tree until all its leaves were pure or had less than two samples left. The decision tree was implemented by Classification and Regression Trees algorithm, which was similar to the C4.5 algorithm and constructed binary trees using the feature and threshold with the largest information gain for each node. It also used Gini impurity and “entropy” to measure the quality of split for a feature. The maximum number of features was set to the total number of features in the dataset. In addition, there was no depth limit for the constructed models. The SVM algorithm was implemented by C-Support Vector Classification algorithm.

III. RESULT

A. Convolutional Neural Network (CNN) Architecture

Our CNN architecture extracted genomic features from the chromosome SNP maps built from Illumina exome genotyping arrays. For each patient, we built two chromosome SNP maps of different chromosomes of HSA21. Each pixel of the chromosome SNP map represented the intensity of one SNP site. Figure 2 showed the CNN architecture that used for human DS prediction/screening. Our CNN architecture took two chromosome SNP maps as input images and fed them into two branch CNN models in combination. Each branch model had five layers, including one input layer, three convolutional layers, and one max-pooling layer. Then we merged two branch models into one model in the fourth convolutional layer (Figure 2, C4), which was followed by one max-pooling layer and three fully connected layers. All hidden layers were followed by dropouts to reduce over-fitting.

B. DS Prediction/Screening Model of Convolutional Neural Network

In this section, we used the CNN architecture (shown in Figure 2) to construct the human DS prediction/screening model with Illumina genotyping array data. The genotyping array dataset contained 378 samples, comprised of 63 DS patients and 315 non-DS patients. First, we normalized intensities of all SNPs into the interval [0,1]. Next, for each sample, we built two chromosome SNP maps based on the SNPs of two HSA1 chromosomes. As shown in Figure 1, columns of the chromosome SNP map represented adjacent genes on the chromosome. Rows represented adjacent SNP sites within the same genes. Each pixel represented the intensity of one SNP site of one gene. Then, we fed two input images into our CNN architectures in combination for

Algorithms	Evaluation Metrics of Different Models					
	Accuracy	Precision	Recall	F-score	False Positive rate	False Negative rate
CNN	99.3(± 0.4)%	99.2%	98.4%	99.3%	0.6%	1.1%
SVM	96.7(± 0.9)%	92.7%	95.9%	94.2%	2.9%	5.3%
Random Forest	97.1(± 0.7)%	94.4%	94.9%	94.7%	1.9%	8.1%
Decision Tree	96.9(± 1.0)%	94.1%	95.4%	94.6%	2.2%	8.0%

Table 1. Evaluation metrics of different DS prediction/screening models.

model training and testing. To build a robust and reliable model, we generated ten parallel training and testing datasets by randomly selecting the first 75% samples for training and using the rest for testing. Then we built ten parallel DS prediction/screening models and calculated average performance metrics to provide a systematic evaluation. As shown in table 1, the DS prediction/classification model built from our CNN architectures achieved an average accuracy of 99.3%. It also achieved very high scores in precision, recall, and F-score, which were 99.2%, 98.4%, and 99.3% respectively. It is worth to notice that the average false positive rate and false negative rates were only 0.6% and 1.1%. In all ten parallel experiments, only five samples with non-DS and two samples with DS in total were mis-predicted to the wrong classes. The above results demonstrated that our CNN architectures constructed robust and accurate CNN model for human DS prediction/screening, which could be used for DS screening using SNP chromosome maps.

C. DS Prediction/Screening Models of Conventional Machine-learning Algorithms

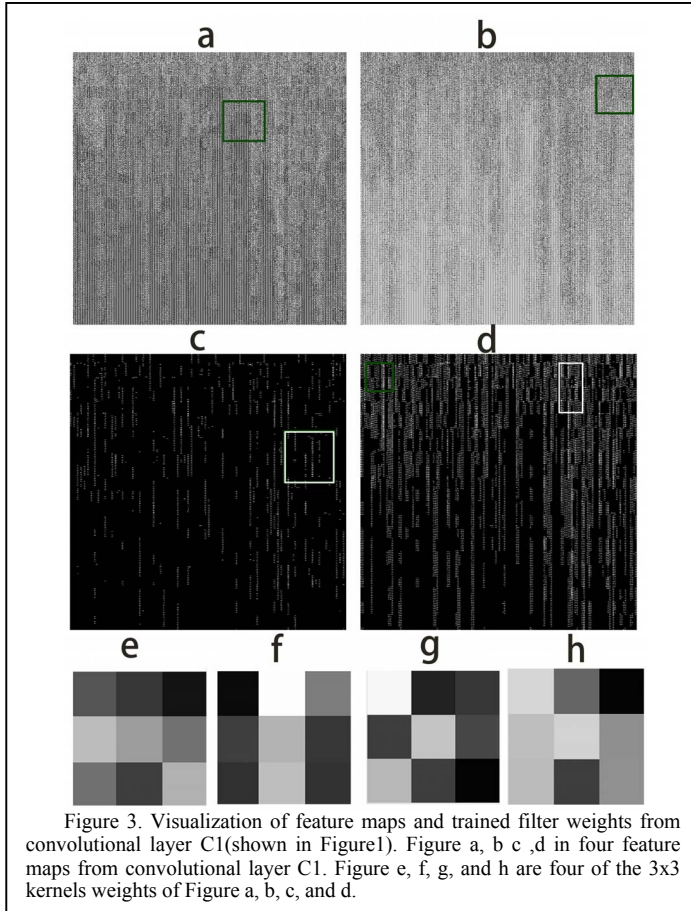
We evaluated and compared the performances of our CNN model and three conventional machine-learning algorithms. We generated ten parallel training and testing datasets by randomly selecting the first 75% data for training and the rest for testing. Next, we used the original Illumina genotyping array data with total 5458 SNP features to construct DS prediction/screening models using support vector machine (SVM), random forest, and decision tree algorithms. Learning algorithms, model constructions and configurations were implemented by Scikit Learn package [37], which were available in the Methods section. Finally, we evaluated and compared the performances of these three conventional supervised learning algorithms and our CNN model.

As Table 1 shown, models built from SVM, random forest, and decision tree could achieve over 96% average accuracies. The random forest model achieved the best performance over all three conventional supervised learning models, which had an average accuracy of 97.1%. The SVM model had the worst performance with average accuracy of 96.7%. The performance of decision tree was between other two models with an average accuracy of 96.9%. However, Table 1 also showed that our CNN models achieved higher average accuracy (99.3%) than all three conventional learning algorithms. It also produced higher average

precisions, recalls, and F-scores. Even though the three conventional models achieved good accuracy scores, their false negative rates were very high, which were 5.3%, 8.1%, and 8.0% respectively. Their high false negative rates made them hard to be put into medical practice. In comparison, our CNN DS prediction/screening model only had an average false negative rate of 1.1% and false positive rate of 0.6%. These findings demonstrated that our CNN model had notably better performances in all evaluating metrics when compared to the other three conventional machine-learning models.

D. Visualization of Feature Maps and Trained Filter Weights

We further visualized the feature maps and their corresponding trained filters from the intermediate layers of our trained CNN model. We discussed the advantages of our CNN methods and the underlying reasons for its powerful performance. First, the CNN model took two chromosome SNP maps as inputs and extracted genomic features from two chromosomes in combination. Each pixel of the chromosome SNP map represented the intensity of one SNP site. We selected the size 3x3 kernel and 16 filters, which would not neglect any detailed information or pattern of adjacent pixels. Figure 3 showed the output feature maps and their corresponding filter weights from the convolutional hidden layer C1 (showed in Figure 2). As shown in the feature maps, firstly, our trained model could also sharpen the input images, and highlight the informative and dense regions (showed in green rectangles). Secondly, our model could detect the lines that represented the intensities of continuous genes and SNP sites (showed in white rectangles). Furthermore, our CNN model could highlight the regional visual patterns and distinctive local motifs from small regions on the chromosomes SNP maps. These regional patterns represented the correlated genetic variations within local genomic regions, as well as adjacent genes and SNP sites. Therefore, our CNN model could keep the most informative and critical SNP sites, and remove less important sites from input images. It could also detect the comprehensive characterization, and extract patterns from the neighboring genomic regions and their simultaneous or causal SNPs variations. However, conventional supervised learning algorithms treat all SNP sites as independent features. Their discriminative power was limited due to the high computational costs in identifying definitive features for



subset characterization and optimization. Also, they tended to construct models with a global view from all available features, and might overlook the regional patterns and correlations among neighboring genes and SNPs.

IV. DISCUSSION

Current DS screening studies include ultrasound measurement of the fetal nuchal translucency [14], blood test [15], sequencing test [16], and combined genetic tests [17]. One of the major benefits of DS screening is offering a non-invasive test to determine the risk of DS. Patient with low-risk could then decide to avoid further invasive diagnostic procedures that might result in fetal loss. Currently, SNPs array analyses of fetal genomes can be performed on fetal trophoblast cells by non-invasive procedures as early as five weeks of gestation [25,26]. In this study, our method made each SNP site as one pixel on the chromosome SNP map, which significantly reduced the image size and training complexity. The two branch CNN models of CNN architecture could learn the local genomic pattern features and correlations of the adjacent genes and SNPs from two input chromosome SNP maps in combination. However, conventional machine-learning algorithms treated all SNPs as independent features. In addition, the discriminative power of these algorithms was limited due to computational costs in identifying definitive features for subset characterization and optimization. Therefore, it was hard for

the conventional machine-learning algorithms to detect the regional patterns and correlations among features. In future work, CNN and other deep learning techniques based on Illumina genotyping array could be further used for diseases predictions and genomic patterns detections.

Previous studies illustrated that the expressions and variations of the local group of genes and SNPs were highly correlated within local chromosome regions in human disease genomes [38]. GWAS studies also proved that SNPs variation, copy number variations and unidentified genetic variation were highly associated with human DS [24–26]. In this study, our CNN DS prediction/screening model took advantage of the comprehensive characterization and regional genetic features extractions. As shown in Figure 3, convolutional layers and filters could detect the genetic patterns and distinctive local motifs from the chromosome SNP maps. Our model could highlight the correlated genetic variations within local genomic regions, and simultaneous or causal variations among adjacent genes and SNP sites. Our results demonstrated that the CNN model not only enabled accurate DS prediction from human genotyping arrays, but also provided an effective analysis of their local genomic characteristics. It could locate the most critical and informative genes and corresponding SNPs associated with human DS disease from the outputs of the intermediate convolutional layers, which provided insights into the mechanisms of DS pathological genomics. Currently, treatments of DS were mainly concentrated on educational therapy [7,8], physical therapy [9,10], emotional and behavioral therapies [8,11], which had limited effects and couldn't cure DS fundamentally. Traditional drugs haven't shown any clear effect or benefit of human DS treatments either [39]. The genetic patterns, correlated genes, and SNPs variation identified by our CNN model provided opportunities to study the genomic markers and pathway components associated with human DS, which further facilitated DS gene therapy studies and genetic medicine developments.

REFERENCES

- [1] Antonarakis, S. E. Down syndrome and the complexity of genome dosage imbalance. *Nat. Rev. Genet.* (2016).
- [2] Patterson, D. Molecular genetic analysis of down syndrome. *Hum. genetics* 126, 195–214 (2009).
- [3] Wiseman, F. K., Alford, K. A., Tybulewicz, V. L. & Fisher, E. M. Down syndrome—recent progress and future prospects. *Hum. molecular genetics* 18, R75–R83 (2009).
- [4] Weijerman, M. E. & De Winter, J. P. Clinical practice. *Eur. journal pediatrics* 169, 1445–1452 (2010).
- [5] Down syndrome. US Department of Health and Human Services | National Institutes of Health (2017)
- [6] Parker, S. E. et al. Updated national birth prevalence estimates for selected birth defects in the united states, 2004–2006. *Birth Defects Res. Part A: Clin. Mol. Teratol.* 88, 1008–1016 (2010).
- [7] Guralnick, M. J. Early intervention approaches to enhance the peer-related social competence of young children with developmental delays: A historical perspective. *Infants young children* 23, 73 (2010).
- [8] Wuang, Y.-P., Chiang, C.-S., Su, C.-Y. & Wang, C.-C. Effectiveness of virtual reality using wii gaming technology in children with down syndrome. *Res. developmental disabilities* 32, 312–321 (2011).

- [9] Pitetti K, Baynard T, Agiovlasitis S. Children and adolescents with Down syndrome, physical fitness and physical activity[J]. *Journal of Sport and Health Science*, 2013, 2(1): 47-57.
- [10] Wentz, E. E. Importance of initiating a “tummy time” intervention early in infants with down syndrome. *Pediatr. Phys. Ther.* 29, 68–75 (2017).
- [11] Greenspan, S. I., Wieder, S. & Simons, R. The child with special needs: Encouraging intellectual and emotional growth.(Addison-Wesley/Addison Wesley Longman, 1998).
- [12] Brock, D. J. & Sutcliffe, R. G. Alpha-fetoprotein in the antenatal diagnosis of anencephaly and spina bifida. *The Lancet* 300, 197–199 (1972).
- [13] Wald, N. J. et al. Maternal serum screening for down’s syndrome in early pregnancy. *Bmj* 297, 883–887 (1988).
- [14] Spencer, K., Souter, V., Tul, N., Snijders, R. & Nicolaides, K. A screening program for trisomy 21 at 10–14 weeks using fetal nuchal translucency, maternal serum free β -human chorionic gonadotropin and pregnancy-associated plasma protein-a. *Ultrasound Obstet. Gynecol.* 13, 231–237 (1999).
- [15] Ehrlich, M. et al. Noninvasive detection of fetal trisomy 21 by sequencing of dna in maternal blood: a study in a clinical setting. *Am. journal obstetrics gynecology* 204, 205–e1 (2011).
- [16] Palomaki, G. E. et al. Dna sequencing of maternal plasma to detect down syndrome: an international clinical validation study. *Genet. medicine* 13, 913–920 (2011).
- [17] Driscoll, D. A. & Gross, S. Prenatal screening for aneuploidy. *New Engl. J. Medicine* 360, 2556–2562 (2009).
- [18] American College of Obstetricians and Gynecologists. ACOG Practice Bulletin No. 88, December 2007. Invasive prenatal testing for aneuploidy[J]. *Obstetrics and gynecology*, 2007, 110(6): 1459.
- [19] Sturgeon, X. & Gardiner, K. J. Transcript catalogs of human chromosome 21 and orthologous chimpanzee and mouse regions. *Mammalian Genome* 22, 261–271 (2011).
- [20] Higuera, C., Gardiner, K. J. & Cios, K. J. Self-organizing feature maps identify proteins critical to learning in a mouse model of down syndrome. *PloS one* 10, e0129126 (2015).
- [21] Dierssen, M. & de la Torre, R. Pathways to cognitive deficits in down syndrome. *Down Syndr. From Underst. Neurobiol. to Ther.* 197, 73 (2012).
- [22] Gardiner, K. et al. Down syndrome: from understanding the neurobiology to therapy. *J. Neurosci.* 30, 14943–14945 (2010).
- [23] Ramachandran, D. et al. Genome-wide association study of down syndrome-associated atrioventricular septal defects. *G3:Genes, Genomes, Genet.* 5, 1961–1971 (2015).
- [24] Sailani, M. R. et al. The complex snp and cnv genetic architecture of the increased risk of congenital heart defects in down syndrome. *Genome research* 23, 1410–1421 (2013).
- [25] Jain, C. V. et al. Fetal genome profiling at 5 weeks of gestation after noninvasive isolation of trophoblast cells from the endocervical canal. *Sci. translational medicine* 8, 363re4–363re4 (2016).
- [26] Petry, C. et al. Associations between a fetal imprinted gene allele score and late pregnancy maternal glucose concentrations. *Diabetes & Metab.* (2017).
- [27] Roth, H. R. et al. Improving computer-aided detection using convolutional neural networks and random view aggregation. *IEEE transactions on medical imaging* 35, 1170–1181 (2016).
- [28] Park, Y. & Kellis, M. Deep learning for regulatory genomics. *Nat. biotechnology* 33, 825–826 (2015).
- [29] Gray, K. R. et al. Random forest-based similarity measures for multi-modal classification of alzheimer’s disease. *NeuroImage* 65, 167–175 (2013).
- [30] Anbarasi, M., Anupriya, E. & Iyengar, N. Enhanced prediction of heart disease with feature subset selection using genetic algorithm. *Int. J. Eng. Sci. Technol.* 2, 5370–5376 (2010).
- [31] Schwarz, J. M., Ro’delsperger, C., Schuelke, M. & Seelow, D. Mutationtaster evaluates disease-causing potential of sequence alterations. *Nat. methods* 7, 575–576 (2010).
- [32] Esteva, A. et al. Dermatologist-level classification of skin cancer with deep neural networks. *Nat.* 542, 115–118 (2017).
- [33] Sun, W., Tseng, T.-L. B., Zhang, J. & Qian, W. Enhancing deep convolutional neural network scheme for breast cancer diagnosis with unlabeled data. *Comput. Med. Imaging Graph.* 57, 4–9 (2017).
- [34] Faust, O. et al. Computer aided diagnosis of coronary artery disease, myocardial infarction and carotid atherosclerosis using ultrasound images: A review. *Phys. Medica* (2016).
- [35] Nguyen, C. D., Costa, A. C., Cios, K. J. & Gardiner, K. J. Machine-learning methods predict locomotor response to mk-801 in mouse models of down syndrome. *J. neurogenetics* 25, 40–51 (2011).
- [36] Zhao, Q. et al. Hierarchical constrained local model using ica and its application to down syndrome detection. *International Conference on Medical Image Computing and Computer-Assisted Intervention*, 222–229 (Springer, 2013).
- [37] Pedregosa, F. et al. Scikit-learn: Machine-learning in python. *J. Mach. Learn. Res.* 12, 2825–2830 (2011).
- [38] Farh, K. K.-H. et al. Genetic and epigenetic fine-mapping of causal autoimmune disease variants. *Nat.* 518, 337 (2015).
- [39] Gardiner, K. J. Pharmacological approaches to improving cognitive function in down syndrome: current status and considerations. *Drug Des Devel Ther* 9, 103–125(2015)