

# Evaluation of Hyperbolic Attention in Histopathology Images

Renyu Zhang  
*Department of Computer Science*  
*University of Chicago*  
 Chicago, U.S.  
 zhangr@uchicago.edu

Aly A. Khan  
*Department of Pathology*  
*University of Chicago*  
 Chicago, U.S.  
 aakhan@uchicago.edu

Robert L. Grossman  
*Department of Medicine & Computer Science*  
*University of Chicago*  
 Chicago, U.S.  
 robert.grossman@uchicago.edu

**Abstract**—We bring together into a common framework three key ideas — multi-scale medical image analysis, the attention mechanism, and hyperbolic embeddings. The formulation and evaluation of hyperbolic-attention models for multi-scale medical image analysis have not been previously explored. In this paper, we evaluate a hyperbolic-attention model on two classification tasks using histopathology image datasets. The experiments show improvement compared to other commonly used models. Our method directly captures the multi-scale structure of histopathology images, and we speculate that the hyperbolic attention mechanism naturally singles out one or more structures at one or more scales that are most discriminatory.

**Index Terms**—hyperbolic attention; histopathology images

## I. INTRODUCTION

Over the past decade, deep learning based methods have achieved great success in computer vision, including state-of-the-art performance in image classification and image retrieval. In the biomedical community, deep neural networks have been able to learn fine-grained features directly from X-ray, CT, MRI and histopathology images and have been applied to complex tasks such as disease and outcome prediction [1], [2]. Such deep neural networks typically perform a series of convolutional transformations that represent images as an embedding in the Euclidean space. However, there is growing evidence that certain types of data, such as data exhibiting a hierarchical or multi-scale structure, are not efficiently represented in Euclidean space [3]–[6].

At the same time, the size of medical images, such as histopathology images are often too large for most modern GPU and off-the-shelf deep learning models. As a result, images are either rescaled to a lower-resolution or segmented into smaller tiles, which can fit into GPU memory. On one hand, rescaling images to lower-resolution can result in distortion and loss of salient image details. On the other hand, the segmentation of images into tiles, though high-resolution, can result in loss of contextual and spatial information. Here, the middle ground provides a tantalizing question: Can we combine differently scaled versions of an image and both harness lower resolution contextual information alongside higher resolution detailed views?

Histopathology images can be viewed under different magnifications or scales, inducing a hierarchical or multi-scale structure. However, robust methods for harnessing multi-scaled

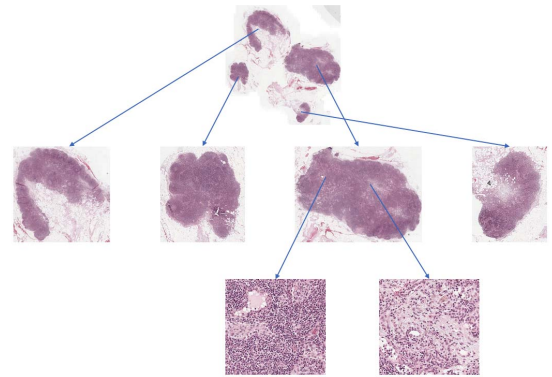


Fig. 1. Example digital hematoxylin and eosin (H&E) stained histopathology slide image with differently scaled views and the relative hierarchy.

views from histopathology images have not been significantly explored. In this work, we propose and evaluate a simple hyperbolic modification to existing deep learning architectures to enable sharing of information from multi-scaled views (Figure 1).

Hyperbolic space is diffeomorphic to standard Euclidean space but has constant negative sectional curvature. A useful model for hyperbolic space is the Poincaré ball model in which the distance from the origin to the boundary increases exponentially. As such, Poincaré embeddings can be thought of as a continuous analog to trees [3]. Intuitively, this makes it suitable for learning embeddings that reflect a natural hierarchy (e.g., image sections at different magnifications) [4]. Our approach is inspired by recent work using Poincaré embeddings to model hierarchical structures [5]–[8].

In this work, we derive a hyperbolic version of standard attention-based models often used in medical imaging. We show through examples that using hyperbolic spaces may improve performance in tissue classification tasks in H&E images. The key contributions are as follows: (1) we extend prior work in computer vision [8] by implementing a novel formulation of hyperbolic attention; (2) we present evidence that multi-scale views of H&E images may be better modeled using hyperbolic spaces. Finally, we conclude the paper with discussion on current challenges and future opportunities.

## II. RELATED WORK

Deep learning methods have typically segmented histopathology images into a collection of tiles due to computing and memory limitations. For example, in [9], a deep convolutional neural network (CNN) was trained on various tiles extracted from whole-slide images, and their predictions were aggregated in order to classify LUAD, LUSC or normal lung samples. We refer to this general approach as a baseline method. Recently, new methods have used the attention mechanism to weight the different tiles, similar to a multiple instance learning (MIL) problem. For example, Katharopoulos and Fleuret [10] define the *Deep MIL* model by pooling tiles based on attention weights. We refer to this general approach as a standard attention-based model.

Hyperbolic spaces have been examined in different deep learning settings due to the increased capacity and structural biases of the embedding space [3], such as link prediction in networks and modeling lexical entailment. Recently, Ganea et. al. [4] and Khurlov et. al. [5] proposed hyperbolic neural networks and hyperbolic attention networks, respectively. Moreover, Lempitsky et. al. [8] demonstrated hyperbolic embeddings can provide a better alternative in many practical scenarios, e.g., few-shot learning and person re-identification. However, the formulation and evaluation of hyperbolic-attention models on medical images has not been significantly explored.

## III. METHOD

We introduce a hyperbolic generalization to traditional CNN architectures in order to extract visual features from different scales and perform operations in hyperbolic space. Briefly, a traditional CNN, such as ResNet18 [11], can be used to generate feature representations of tiles extracted from whole-slide H&E images. Next, we define a bijective mapping of the generated features to hyperbolic space. We then define linear, multinomial regression and attention layers operating in hyperbolic space. We note the attention layer operates in hyperbolic space and computes a slide-level representation. Finally, we define classification using a multinomial regression layer. We derive and present our formal definitions in the following subsections and adopt notations from [4] and [8]. We denote the resulting model as a hyperbolic-attention model.

### A. Poincaré ball model

The Poincaré model  $(\mathbb{D}^n, g^{\mathbb{D}})$  is defined by the manifold  $\mathbb{D}^n = \{x \in \mathbb{R}^n : \|x\| < 1\}$  equipped with Riemannian metric  $g_x^{\mathbb{D}} = \lambda_x^2 g^E$  where  $\lambda_x := \frac{2}{1-\|x\|^2}$ .  $g^E = \mathbf{I}_n$  is the Euclidean metric tensor. To make use of the Poincaré ball of radius  $c \geq 0$ , we denote  $\mathbb{D}_c^n := \{x \in \mathbb{R}^n | c\|x\|^2 < 1\}$ . If  $c = 0$ ,  $\mathbb{D}_c^n = \mathbb{R}^n$ ; If  $c > 0$ , it is an open ball with radius  $1/\sqrt{c}$ .

### B. Möbius addition

For a pair of  $\mathbf{x}, \mathbf{y} \in \mathbb{D}_c^n$ , the Möbius addition is defined as follows

$$\mathbf{x} \oplus_c \mathbf{y} := \frac{(1 + 2c\langle \mathbf{x}, \mathbf{y} \rangle + c\|\mathbf{y}\|^2) \mathbf{x} + (1 - c\|\mathbf{x}\|^2) \mathbf{y}}{1 + 2c\langle \mathbf{x}, \mathbf{y} \rangle + c^2\|\mathbf{x}\|^2\|\mathbf{y}\|^2} \quad (1)$$

With  $c \rightarrow 0$  we can obtain the Euclidean distance of two vectors in  $\mathbb{R}^n$ .

### C. Exponential and logarithmic maps

In order to do operations in hyperbolic space, bijective maps are defined to map from  $\mathbb{R}^n$  to  $\mathbb{D}_c^n$ . The exponential map  $\exp_x^c$  is a function from  $\mathbb{R}^n$  to  $\mathbb{D}_c^n$

$$\exp_x^c(\mathbf{v}) := \mathbf{x} \oplus_c \left( \tanh \left( \sqrt{c} \frac{\lambda_x^c \|\mathbf{v}\|}{2} \right) \frac{\mathbf{v}}{\sqrt{c\|\mathbf{v}\|}} \right) \quad (2)$$

The inverse map is defined as

$$\log_x^c(\mathbf{y}) := \frac{2}{\sqrt{c\lambda_x^c}} \operatorname{arctanh}(\sqrt{c}\|\mathbf{x} \oplus_c \mathbf{y}\|) \frac{-\mathbf{x} \oplus_c \mathbf{y}}{\|\mathbf{x} \oplus_c \mathbf{y}\|} \quad (3)$$

In practice, we use the maps  $\exp_0^c$  and  $\log_0^c$  for transition between the Euclidean and Poincaré ball representations of a vector.

### D. Hyperbolic linear layer

Similar to [8], we define a hyperbolic linear layer a map from  $\mathbb{D}_c^{n_1}$  to  $\mathbb{D}_c^{n_2}$ . For input  $\mathbf{x} \in \mathbb{D}_c^{n_1}$  to this layer and a trainable matrix  $\mathbf{M}$  of size  $n_2 \times n_1$ , if  $\mathbf{M}\mathbf{x} \neq 0$ , the output of this layer is

$$\mathbf{M}^c(\mathbf{x}) := \frac{1}{\sqrt{c}} \tanh \left( \frac{\|\mathbf{M}\mathbf{x}\|}{\|\mathbf{x}\|} \operatorname{arctanh}(\sqrt{c}\|\mathbf{x}\|) \right) \frac{\mathbf{M}\mathbf{x}}{\|\mathbf{M}\mathbf{x}\|} \quad (4)$$

otherwise  $\mathbf{M}^c(\mathbf{x}) := 0$ . For a bias vector  $\mathbf{b} \in \mathbb{D}_c^{n_2}$ , the corresponding linear layer is  $\mathbf{M}^c(\mathbf{x}) \oplus_c \mathbf{b}$ .

### E. Klein model

In order to define hyperbolic attention model, we will make use of Klein model and hyperbolic averaging. Similar to Poincaré model, it is defined in  $\mathbb{K}^n = \{\mathbf{x} \in \mathbb{R}^n : \|\mathbf{x}\| < 1\}$ . Let  $\mathbf{x}^{\mathbb{D}}$  and  $\mathbf{x}^{\mathbb{K}}$  denote the coordinates of the same point in the Poincaré and Klein models. We use the following formulas to map from each other.

$$\mathbf{x}^{\mathbb{D}} = \frac{\mathbf{x}^{\mathbb{K}}}{1 + \sqrt{1 - c\|\mathbf{x}^{\mathbb{K}}\|^2}} \quad (5)$$

$$\mathbf{x}^{\mathbb{K}} = \frac{2\mathbf{x}^{\mathbb{D}}}{1 + c\|\mathbf{x}^{\mathbb{D}}\|^2} \quad (6)$$

### F. Hyperbolic attention

Given a query, key and value  $\mathbf{q}_i^{\mathbb{D}}, \mathbf{k}_j^{\mathbb{D}}, \mathbf{v}_{ij}^{\mathbb{D}} \in \mathbb{D}^n$ , the corresponding coordinates in Klein model are  $\mathbf{q}_i^{\mathbb{K}}, \mathbf{k}_j^{\mathbb{K}}, \mathbf{v}_{ij}^{\mathbb{K}}$ , we define the attention weights  $\alpha_{ij}$  as follows.

$$\alpha_{ij} = f(\mathbf{q}_i^{\mathbb{D}}, \mathbf{k}_j^{\mathbb{D}}) \quad (7)$$

The function  $f(\cdot)$  is a hyperbolic neural network followed by softmax or sigmoid. The outputs  $m_i^{\mathbb{K}}$  for the hyperbolic attention module are as follows.

$$m_i^{\mathbb{K}} \left( \{\alpha_{ij}\}_j, \{\mathbf{v}_{ij}^{\mathbb{K}}\}_j \right) = \sum_j \left[ \frac{\alpha_{ij} \gamma(\mathbf{v}_{ij}^{\mathbb{K}})}{\sum_{\ell} \alpha_{i\ell} \gamma(\mathbf{v}_{i\ell}^{\mathbb{K}})} \right] \mathbf{v}_{ij}^{\mathbb{K}} \quad (8)$$

where the  $\gamma(\mathbf{v}_{ij}^{\mathbb{K}})$  are the Lorentz factors,

$$\gamma(\mathbf{v}_{ij}^{\mathbb{K}}) = \frac{1}{\sqrt{1 - c \|\mathbf{v}_{ij}^{\mathbb{K}}\|^2}} \quad (9)$$

Similar to [12], we modify the attention mechanism to get slide-level representations. The attention weight  $\alpha_{ij}$  is only computed based on  $\mathbf{v}_{ij}^{\mathbb{D}}$

$$\alpha_{ij} = f(\mathbf{v}_{ij}^{\mathbb{D}}) \quad (10)$$

After we get the hyperbolic attention output  $m_i^{\mathbb{D}}$  of Klein model, we can map it to Poincaré model.

### G. Multiclass logistic regression

The resulting formula for hyperbolic multiclass logistic regression for  $K$  classes is written below; here  $p_k \in \mathbb{D}_c^n$  and  $a_k \in T_{\mathbf{p}_k} \mathbb{D}_c^n \setminus \{\mathbf{0}\}$  are learnable parameters.

$$p(y = k | \mathbf{x}) \propto \exp \left( \frac{\lambda_{\mathbf{p}_k}^c \|\mathbf{a}_k\|}{\sqrt{c}} \operatorname{arcsinh} \left( \frac{2\sqrt{c} \langle -\mathbf{p}_k \oplus_c \mathbf{x}, \mathbf{a}_k \rangle}{(1 - c \|\mathbf{p}_k \oplus_c \mathbf{x}\|^2) \|\mathbf{a}_k\|} \right) \right) \quad (11)$$

## IV. RESULTS

We compare performance of our hyperbolic-attention model with a baseline model and the *Deep MIL* attention-based model on simple tissue classification tasks using two well-known public datasets. We provide a thorough description of each experiment in following subsections.

### A. Camelyon16

The Camelyon16 challenge [13] was organized by the IEEE International Symposium on Biomedical Imaging. It evaluated various machine learning models to detect cancer metastasis. There are 159 slides with normal tissue class labels and 111 slides with tumor tissue class labels in the training set. The test data set contains 80 slides with normal tissue and 50 slides with tumor tissue. In this work, we evaluate our model on classifications based on slide-level annotations. We perform minimal data pre-processing. Tiles sizes of 500x500 pixels (500px) and 1000x1000 pixels (1000px) are extracted without overlap, and 2000x2000 pixels (2000px) tiles are extracted with step size 1000. In order to filter out non-tissue containing or background tiles, we only keep those tiles with average intensity less than 0.85 and greater than 0.2.

We implemented a baseline approach similar to [9]. We used ResNet18 [11] pretrained on ImageNet [14]. We fine-tuned the model by updating all layers to classify tiles extracted from H&E slides. Tiles of different sizes were re-scaled to the default ResNet18 input layer size, and generated embeddings of length 10. We labeled all tiles with the same label as the slide from which they are extracted. During validation or testing, the model aggregated all the tile predictions of a slide by taking the average and used this average as a slide-level prediction.

We implemented the *Deep MIL* attention-model by randomly sampling 5 tiles as input for each slide. We used ResNet18 pretrained on ImageNet as a feature extractor. We again fine-tuned the model by updating all layers to classify tiles extracted from H&E slides. Similarly, tiles of different sizes were re-scaled to the default ResNet18 input layer size, and generated embeddings of length 10. We performed test-time augmentation for each slide in the test data set, where the model generated a prediction for each test slide 10 times and randomly sampled 5 tiles each time. We calculated the average of all the 10 predictions as the final output.

We sought to compare performance between the baseline model and the *Deep MIL* attention-based model using tiles from a single fixed scale view. We independently evaluated performance on 3 different scales (Table I). The models were tested in 5-fold cross-validation manner. We choose 4 folds as training and 1 fold as a validation data set for each assignment. We trained the models 4 times for each assignment. Both models were trained with Adam optimization method with learning rate=1e-4,  $\beta_1 = 0.9$ ,  $\beta_2 = 0.999$  and  $\epsilon=1e-8$ . To ensure numerical stability, clipping by norm, similar to [8], is performed. For all the 20 checkpoints of each model, we tested the performance on the test dataset and the mean AUC and confidence interval (CI) are reported in Table I. We can see that for the baseline model the best tile size is 1000. *Deep MIL* performs similarly for the 3 scales and demonstrates it is always better than the baseline model.

TABLE I  
PERFORMANCE OF SINGLE SCALE MODELS

Scale	Baseline model		Deep MIL	
	Mean AUC	CI(0.95)	Mean AUC	CI(0.95)
500px	0.569	[0.474,0.663]	0.619	[0.524,0.714]
1000px	0.597	[0.467,0.726]	0.602	[0.505,0.700]
2000px	0.583	[0.460,0.706]	0.613	[0.540,0.687]

Next, we sought to compare the performance between the *Deep MIL* attention-model and our hyperbolic attention model in a multi-scale setting. We combined and used 15 tiles, 5 tiles from each of the previous 3 scales, to train and evaluate slide-level classification accuracy. We also examined the effect of different embedding lengths, by setting lengths to 5, 10 and 100. The results are reported in Table II. By combining tiles from different scales, the *Deep MIL* model does not outperform single scale. However, the hyperbolic attention model is better than *Deep MIL* when using multi-scale views. This suggests that our hyperbolic attention model can learn better multi-scale embeddings.

TABLE II  
PERFORMANCE OF MULTIPLE SCALE MODELS

Embed size	Deep MIL		Hyperbolic attention	
	Mean AUC	CI(0.95)	Mean AUC	CI(0.95)
5	0.602	[0.513,0.690]	0.623	[0.536,0.710]
10	0.606	[0.519,0.693]	0.615	[0.537,0.692]
100	0.592	[0.491,0.694]	0.637	[0.574,0.700]

## B. TCGA

We consider a second general tissue classification task involving normal and lung cancer subtypes (LUAD and LUSC) as presented in [9]. We downloaded H&E lung slides from the TCGA Genomic Data Commons [15]. There were 811 LUAD slides, 745 LUSC slides, and 585 adjacent normal slides. We processed all tiles in the same way as we did for the Camelyon16 data set.

We evaluated our model with the baseline model and the *Deep MIL* model on the lung classification task. All slides were again split into 5-fold and we tested all models on the lung data set in 5-fold cross-validation manner. The models were trained 4 times with the same settings as with the Camelyon16 dataset. The mean Macro-average AUC and CI(0.95) of 20 checkpoints are reported in Table III and Table IV.

Table III shows the performance of baseline models and *Deep MIL* models trained on 3 different scales. We find that the *Deep MIL* models produce slightly better results.

TABLE III  
PERFORMANCE OF SINGLE SCALE MODELS

Scale	Baseline model		Deep MIL	
	Mean AUC	CI(0.95)	Mean AUC	CI(0.95)
500px	0.969	[0.954,0.985]	0.973	[0.957,0.989]
1000px	0.971	[0.955,0.986]	0.972	[0.954,0.990]
2000px	0.966	[0.945,0.987]	0.967	[0.949,0.985]

Table IV shows the performance of *Deep MIL* and our hyperbolic attention model when we combine tiles from 3 scales. Note that in this example, single-scale and multiple-scale *Deep MIL* produce comparable results. Our model's performance is better when the embedding size are 5 or 10. When the embedding size is 100, the performances of the two models are comparable. Overall, this again suggests that our hyperbolic attention model can learn better multi-scale embeddings.

TABLE IV  
PERFORMANCE OF MULTIPLE SCALE MODELS

Embed size	Deep MIL		Hyperbolic attention	
	Mean AUC	CI(0.95)	Mean AUC	CI(0.95)
5	0.967	[0.952,0.983]	0.971	[0.958,0.985]
10	0.970	[0.948,0.992]	0.974	[0.961,0.988]
100	0.973	[0.958,0.988]	0.972	[0.958,0.986]

## V. DISCUSSION AND CONCLUSION

In this paper, we develop a hyperbolic-attention model using a Poincaré ball and Klein model to classify histopathology slide images. Our results suggest that our hyperbolic attention model can efficiently learn multi-scale embeddings. However, we note some important limitations in our work and plan to explore these in future work. First, while our hyperbolic-attention model achieves better performance over the standard baseline and the *Deep MIL* model, there is still room for improvement, especially when the slide number is limited.

Second, we did not explore an integrative interpretation of hyperbolic space geometry and the use of attention to identify salient structures and scales that are associated with improved model performance. Third, while we evaluated our model on two independent datasets, we think future work should also examine additional datasets and data types, such as X-ray images.

## REFERENCES

- [1] G. Litjens, T. Kooi, B. E. Bejnordi, A. A. A. Setio, F. Ciompi, M. Ghahfouari, J. A. Van Der Laak, B. Van Ginneken, and C. I. Sánchez, "A survey on deep learning in medical image analysis," *Medical image analysis*, vol. 42, pp. 60–88, 2017.
- [2] D. Shen, G. Wu, and H.-I. Suk, "Deep learning in medical image analysis," *Annual review of biomedical engineering*, vol. 19, pp. 221–248, 2017.
- [3] M. Nickel and D. Kiela, "Poincaré embeddings for learning hierarchical representations," in *Advances in Neural Information Processing Systems 30*, I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, Eds. Curran Associates, Inc., 2017, pp. 6338–6347. [Online]. Available: <http://papers.nips.cc/paper/7213-poincare-embeddings-for-learning-hierarchical-representations.pdf>
- [4] O. Ganea, G. Becigneul, and T. Hofmann, "Hyperbolic neural networks," in *Advances in Neural Information Processing Systems 31*, S. Bengio, H. Wallach, H. Larochelle, K. Grauman, N. Cesa-Bianchi, and R. Garnett, Eds. Curran Associates, Inc., 2018, pp. 5345–5355. [Online]. Available: <http://papers.nips.cc/paper/7780-hyperbolic-neural-networks.pdf>
- [5] Ç. Gülçehre, M. Denil, M. Malinowski, A. Razavi, R. Pascanu, K. M. Hermann, P. W. Battaglia, V. Bapst, D. Raposo, A. Santoro, and N. de Freitas, "Hyperbolic attention networks," *CoRR*, vol. abs/1805.09786, 2018. [Online]. Available: <http://arxiv.org/abs/1805.09786>
- [6] I. Chami, R. Ying, C. Ré, and J. Leskovec, "Hyperbolic graph convolutional neural networks," 2019.
- [7] C. Gulcehre, M. Denil, M. Malinowski, A. Razavi, R. Pascanu, K. M. Hermann, P. Battaglia, V. Bapst, D. Raposo, A. Santoro *et al.*, "Hyperbolic attention networks," *arXiv preprint arXiv:1805.09786*, 2018.
- [8] V. Khruikov, L. Mirvakhabova, E. Ustinova, I. V. Oseledets, and V. S. Lempitsky, "Hyperbolic image embeddings," *CoRR*, vol. abs/1904.02239, 2019. [Online]. Available: <http://arxiv.org/abs/1904.02239>
- [9] N. Coudray, P. S. Ocampo, T. Sakellaropoulos, N. Narula, M. Snuderl, D. Fenyö, A. L. Moreira, N. Razavian, and A. Tsirigos, "Classification and mutation prediction from non-small cell lung cancer histopathology images using deep learning," *Nature medicine*, vol. 24, no. 10, pp. 1559–1567, 2018.
- [10] A. Katharopoulos and F. Fleuret, "Processing megapixel images with deep attention-sampling models," *CoRR*, vol. abs/1905.03711, 2019. [Online]. Available: <http://arxiv.org/abs/1905.03711>
- [11] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," *CoRR*, vol. abs/1512.03385, 2015. [Online]. Available: <http://arxiv.org/abs/1512.03385>
- [12] M. Ilse, J. M. Tomczak, and M. Welling, "Attention-based deep multiple instance learning," *CoRR*, vol. abs/1802.04712, 2018. [Online]. Available: <http://arxiv.org/abs/1802.04712>
- [13] G. Litjens, P. Bandi, B. Ehteshami Bejnordi, O. Geessink, M. Balkenhol, P. Bult, A. Halilovic, M. Hermsen, R. van de Loo, R. Vogels, Q. F. Manson, N. Stathonikos, A. Baidoshvili, P. van Diest, C. Wauters, M. van Dijk, and J. van der Laak, "1399 H&E-stained sentinel lymph node sections of breast cancer patients: the CAMELYON dataset," *GigaScience*, vol. 7, no. 6, 05 2018, giy065. [Online]. Available: <https://doi.org/10.1093/gigascience/giy065>
- [14] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei, "ImageNet Large Scale Visual Recognition Challenge," *International Journal of Computer Vision (IJCV)*, vol. 115, no. 3, pp. 211–252, 2015.
- [15] R. L. Grossman, A. P. Heath, V. Ferretti, H. E. Varmus, D. R. Lowy, W. A. Kibbe, and L. M. Staudt, "Toward a shared vision for cancer genomic data," *New England Journal of Medicine*, vol. 375, no. 12, pp. 1109–1112, 2016.