# Unsupervised Learning for Cell-level Visual Representation in Histopathology Images with Generative Adversarial Networks

Bo Hu♯, Ye Tang♯, Eric I-Chao Chang, Yubo Fan, Maode Lai and Yan Xu*

*Abstract*—The visual attributes of cells, such as the nuclear morphology and chromatin openness, are critical for histopathology image analysis. By learning cell-level visual representation, we can obtain a rich mix of features that are highly reusable for various tasks, such as cell-level classification, nuclei segmentation, and cell counting. In this paper, we propose a unified generative adversarial networks architecture with a new formulation of loss to perform robust cell-level visual representation learning in an unsupervised setting. Our model is not only label-free and easily trained but also capable of cell-level unsupervised classification with interpretable visualization, which achieves promising results in the unsupervised classification of bone marrow cellular components. Based on the proposed cell-level visual representation learning, we further develop a pipeline that exploits the varieties of cellular elements to perform histopathology image classification, the advantages of which are demonstrated on bone marrow datasets.

*Keywords*—*unsupervised learning, representation learning, generative adversarial networks, classification, cell.*

## I. Introduction

HISTOPATHOLOGY images are considered to be the gold standard in the diagnosis of many diseases [1]. In many situations, the cellular components are an important determinant. For example, in the biopsy sections of bone marrow, the abnormal cellular constitution indicates the presence of blood disease [2]. Bone marrow is the key component of both the hematopoietic system and the lymphatic system

Bo Hu, Ye Tang, Yubo Fan and Yan Xu are with the State Key Laboratory of Software Development Environment and the Key Laboratory of Biomechanics and Mechanobiology of Ministry of Education and Research Institute of Beihang University in Shenzhen and Beijing Advanced Innovation Centre for Biomedical Engineering, Beihang University, Beijing 100191, China (email: bohu1996@gmail.com; yetang1995@gmail.com; yubofan@buaa.edu.cn; xuyan04@gmail.com).

Maode Lai is with the Department of Pathology, School of Medicine, Zhejiang University (email: lmd@zju.edu.cn).

Eric I-Chao Chang, and Yan Xu are with Microsoft Research, Beijing 100080, China (email: echang@microsoft.com; v-yanx@microsoft.com).
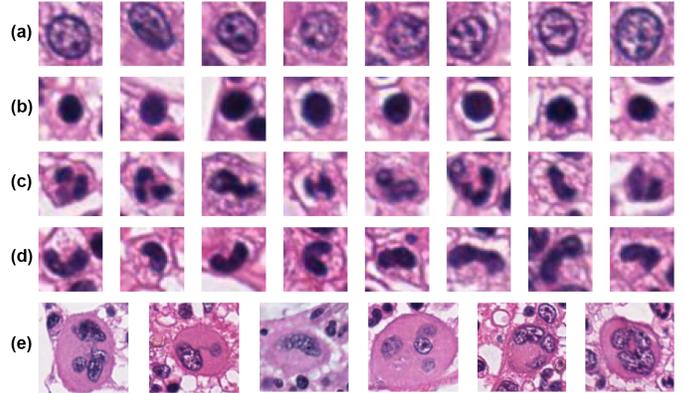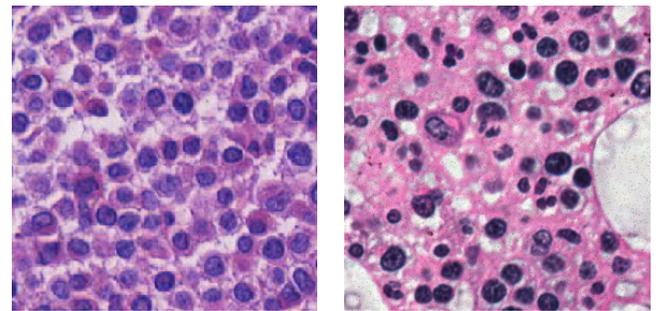


Fig. 1. Examples of five types of cellular elements in bone marrow: (*a*) granulocytes precursors such as myeloblasts, (*b*) cells with dark, dense, and close phased nuclei, the candidates of which are most likely lymphocytes and normoblasts, (*c*) granulocytes such as neutrophils, (*d*) monocytes, and (*e*) megakaryocytes. Five types of cells can be distinguished by the chromatin openness, the density of nuclei, and if nuclei show the appearance of being segmented. Megakaryocytes appear the least often, as well are the most distinguished due to their massive size.



(a) abnormal       (b) normal

Fig. 2. Examples of bone marrow images sliced from Whole Slide Images (WSI). Too many myeloblasts in (a) indicate the presence of blood disease.

by producing large amounts of blood cells. The cell lines undergoing maturation in the marrow mostly include myeloid cells (granulocytes, monocytes, megakaryocytes, and their precursors), erythroid cells (normoblasts), and lymphoid cells (lymphocytes and their precursors). Figure 1 are examples of five main cellular components in bone marrow. These components are significant to both the systemic circulation and the immune system. Several kinds of cancer are characterized

by the cellular constitution in bone marrow [2]. For instance, too many granulocytes precursors such as myeloblasts indicate the presence of chronic myeloid leukemia. Having large, abnormal lymphocytes heralds the presence of lymphoma. Figure 2 shows the difference between normal and abnormal bone marrow histopathology images from the perspective of cells.

As described above, cell-level information is irreplaceable for histopathology image analysis. Cell-level visual attributes such as the morphological features of nuclei and the openness of chromatin are helpful for various tasks such as cell-level classification and nuclei segmentation. We define cell-level images as the output from nuclei segmentation. Each cell-level image contains only one cell. We opt to perform representation learning on these cell-level images, in which the visual attributes such as the nuclei morphology and chromatin openness are distinguished. The learned features are further utilized to assist tasks such as cell counting to highlight the quantification of certain types of cells.

To achieve this, the main obstacle is the labeling of cells. There are massive amounts of cells in each histopathology image, which makes manual labeling ambiguous and laborious. Therefore, an unsupervised cell-level visual representation learning method based on unlabeled data is believed to be more reasonable than fully supervised methods. Unsupervised cell-level visual representation learning is known to be difficult. First, geometrical and morphological appearances of cells from the same category can have a distinct diversity due to factors such as cell cycles. Furthermore, the staining conditions of histopathology images can be pretty diverse, resulting in inconsistent color characteristics of nuclei and cytoplasm.

Recently, deep learning has been proven to be powerful in histopathology image analysis such as classification [3], [4], segmentation [5], [6], and detection [7], [8]. Generative Adversarial Networks (GANs) [9] are a class of generative models that use unlabeled data to perform representation learning. GAN is capable of transforming noise variables into visually appealing image samples by learning a model distribution that imitates the real data distribution. Several GAN architectures such as Deep Convolutional Generative Adversarial Nets (DCGAN) [10] have proven their advantages in various natural images datasets. Recently, Wasserstein-GAN (WGAN) [11] and WGAN with gradient penalty (WGAN-GP) [12] have greatly improved the stability of training GAN. More complex network structures such as residual networks [13] can now be fused into GAN models.

Meanwhile, Information Maximizing Generative Adversarial Networks (InfoGAN) [14] makes a modification that encourages GAN to learn interpretable and meaningful representations. InfoGAN maximizes the mutual information between the chosen random variables and the observations to make variables represent interpretable semantic features. The problem is that InfoGAN utilizes a DCGAN architecture, which requires meticulous attention towards hyperparameters. For our problem, it suffers a severe convergence problem.

Inspired by WGAN-GP and InfoGAN, we present an unsupervised representation learning method for cell-level images using a unified GAN architecture with a new formulation of loss, which inherits the superiority from both WGAN-GP and InfoGAN. We observe great improvements followed by the setting of WGAN-GP. Introducing mutual information into our formulation, we are capable of learning interpretable and disentangled cell-level visual representations, as well as allocate cells into different categories according to their most significant semantic features. Our method achieves promising results in the unsupervised classification of bone marrow cellular components.

Based on the cell-level visual representations, the quantification of each cellular component can be obtained by the trained model. Followed by this, cell proportions for each histopathology image can then be calculated to assist image-level classification. We further develop a pipeline combining cell-level unsupervised classification and nuclei segmentation to conduct image-level classification of histopathology images, which shows its advantages via experimentations on bone marrow datasets.

The contributions of this work include the following: (1) We present an unsupervised framework to perform cell-level visual representation learning using generative adversarial networks. (2) A unified GAN architecture with a new formulation of loss is proposed to generate representations that are both high-quality and interpretable, which also endows our model the capability of cell-level unsupervised classification. (3) A pipeline is developed that exploits the varieties of cell-level elements to perform image-level classification of histopathology images.

## II. RELATED WORKS

### A. Directly Related Works

*1) Generative Adversarial Networks:* Goodfellow et al. [9] propose GANs, a class of unsupervised generative models consisting of a generator neural network and an adversarial discriminator neural network. While the generator is encouraged to produce synthetic samples, the discriminator learns to discriminate between generated and real samples. This process is described as a minimax game. Radford et al. [10] propose one of the most frequently used GAN architectures DCGAN.

Arjovsky et al. [11] propose WGAN, which modifies the objective function, securing the training process to be more stable. For regular GANs, the training process optimizes a lower bound of the Jensen-Shannon (JS) divergence between the generator distribution and the real data distribution. WGAN modifies this by optimizing an approximation of the Earth-Mover (EM) distance. The only challenge is how to enforce the Lipschitz constraint on the discriminator. While Arjovsky et al. [11] use weight-clipping, Gulrajani et al. [12] propose WGAN-GP, which adds a gradient penalty on the discriminator. For our bone marrow datasets, even if we have tried multiple hyperparameters, DCGAN still suffers from a severe convergence difficulty. While DCGAN leads to the failure for our datasets, WGAN-GP greatly eases this problem.

Chen et al. [14] introduce mutual information into GAN architecture. Mutual information describes the dependencies between two separate variables. Maximizing mutual information between the chosen random variables and the generated

samples, InfoGAN produces representations that are meaningful and interpretable. To exploit the varieties of cellular components, the superior ability of InfoGAN in learning disentangled and discrete representations is what a regular GAN lacks.

Therefore, we propose a unified GAN architecture with a new formulation of loss, which inherits the superiority of both WGAN-GP and InfoGAN. The outstanding stability of WGAN-GP eases the difficulty in tuning the complicated hyperparameters of InfoGAN. Introducing mutual information into our model, we are capable of learning interpretable cell-level visual representations, as well as allocate cells into different categories according to their most significant semantic features.

*2) Classification of Blood Disease:* Nazlibilek et al. [15] propose a system to help automatically diagnose acute lymphocytic leukemia. This system consists of several stages: nuclei segmentation, feature extraction, cell-level classification, and cell counting. In their future work, they claim that the result of cell counting can be used for further diagnosis of acute lymphocytic leukemia.

In our work, we design a similar workflow which consists of nuclei segmentation, cell-level classification, and image-level classification. Our advantages lie in the novelty of an unsupervised setting and the convincing performance of image-level classification based on the calculated cell proportions.

### B. Cell-level Representation

The representation of individual cells can be used for a variety of tasks such as cell classification. Traditional cell-level visual representation for classification tasks can be categorized into four categories [16]: morphological [17], texture [18], [19], intensity [20], and cytology features [21]. These traditional methods have been employed in the representation of white blood cells [22], [23], [24]. However, the features used above need to be manually designed by experienced experts according to the characteristics of different types of cells. While images suffer from a distinct variance, discovering, characterizing and selecting good handcraft features can be extremely difficult.

To remedy the limitations of manual features in cell classification, Convolutional Neural Network (CNN) learns higher-level latent features, whose convolution layer can act as a feature extractor [25]. Xie et al. [26] propose Deep Embedding Clustering (DEC) that simultaneously learns feature representations and cluster assignments using deep neural networks.

Variational Autoencoder (VAE) [27] serves as a convincing unsupervised strategy in cell-level visual representation learning [28], [29], [30]. However, how to use VAE to learn categorical and discrete latent variables is still under investigation. Dilokthanakul et al. [31] and Jiang et al. [32] design models combining VAE with Gaussian Mixture Model (GMM). But they demonstrate their experiment on one-dimensional datasets such as MNIST. To perform clustering and embedding on a higher-dimensional dataset, their methods still need a feature extractor.

GANs such as Categorical GAN [33] can merge categorical variables into the model with little effort, which makes learned

representations disentangled and interpretable. This ability is critical in medical image analysis where accountability is especially needed.

### C. Cell-level Histopathology Image Analysis

*1) Classification:* Cell classification has been performed in diverse histopathology related works such as breast cancer [34], acute lymphocytes leukemia [35], [36], and colon cancer [37].

Based on the result of cell classification, some approaches have been proposed to determine the presence or location of cancer [21], [38]. In prostate cancer, Nguyen et al. [21] innovatively employ cell classification for automatic cancer detection and grading. They distinguish the cancer nuclei and normal nuclei, which are combined with textural features to classify the image as normal or cancerous and then detect and grade the cancer regions. In the diagnosis of Glioma, Hou et al. [38] apply CNN to the classification of morphological attributes of nuclei. They also claim that the nuclei classification result provides clinical information for diagnosing and classifying glioma into subtypes and grades. Zhang et al. [39], [40], [41] and Shi et al. [42] use either supervised or semi-supervised hashing models for cell-level analysis.

All of these works require a large amount of accurately annotated data. Obtaining such annotated data is time-consuming and labor-intensive while GAN can optimally leverage the wealth of unlabeled data.

*2) Segmentation:* Nuclei segmentation is of great importance for cell-level classification. Nuclei segmentation methods can be roughly categorized as follows: intensity thresholding [43], [44], morphology operation [45], [46], deformable models [47], watershed transform [48], clustering [49], [50], and graph-based methods [51], [52]. The methods above have been broadly applied to the segmentation of white blood cells.
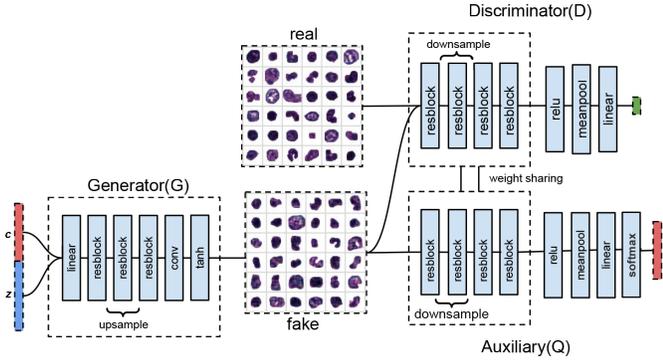
### D. Generative Adversarial Networks in Medical Images

Recently, several works involving GAN have gathered great attention in medical image analysis.
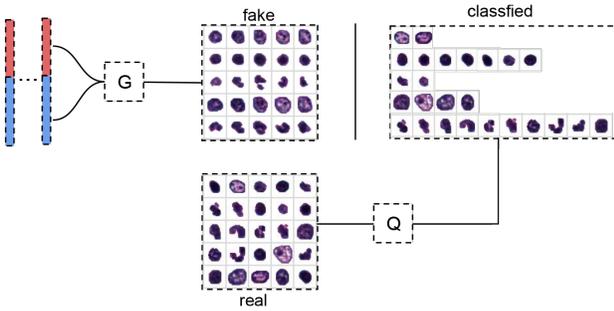
In medical image synthesizing, Nie et al. [53] estimate the CT image from its corresponding MR image with context-aware GAN. In medical image reconstruction, Li et al. [54] use GAN to reconstruct medical images with the thinner sliced thickness from regular thick-slice images. Mahapatra et al. [55] propose a super resolution method that takes a low-resolution input fundus image to generate a high-resolution super-resolved image. Wolterink et al. [56] employ GAN to reduce the noise in low-dose CT images. All these recent works demonstrate the great potential of GAN in solving complicated medical problems.
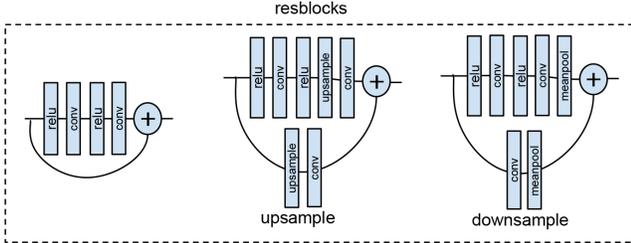
### III. METHODS

In this section, we first introduce an unsupervised method for cell-level visual representation learning using GAN. Then we present the details of how image-level classification is performed on histopathology images based on cell-level representation.

(a) Training process. Random variables are composed of Gaussian variables $z$ and the discrete variable $c$. Besides playing the minimax game between the generator ($G$) and the discriminator ($D$) through the EM distance, we also minimize the negative Log-likelihood between $c$ and the output of the auxiliary network ($Q(c|G(c, z))$ to maximize mutual information.



(b) Test process. Real samples are classified into five categories by the auxiliary network $Q$. At the same time, fake samples are generated by giving noises with the chosen $c$ for each class. In the example of generated samples (fake), one row contains five samples from the same category in $c$, and a column shows the generated images for 5 possible categories in $c$ with $z$ fixed.



(c) Illustration of residual blocks (resblocks) in the architecture. There are three different types of residual blocks considering whether they include nearest-neighbor upsampling or mean pooling for downsampling. Batch normalization layers are used in our generator to help stabilize training.

Fig. 3. Network architecture of our cell-level visual representation learning: (a) Training process. (b) Test process. (c) The architecture of residual blocks (written as resblock in (a) and (b)).

## A. Cell-level Visual Representation Learning

Given cell-level images that come from nuclei segmentation as the real data, we define a generator network $G$, a discriminator network $D$, and an auxiliary network $Q$. The architecture of these networks are shown in Figure 3. In the training process, we learn a generator distribution that matches the real data distribution by playing a minimax game between

$G$ and $D$ by optimizing an approximation of the Earth-Mover (EM) distance. Meanwhile, we maximize mutual information between the chosen random variables and the generated samples using an auxiliary network $Q$. In the test process, the generator generates the representations for each category of cells according to different values of the chosen random variables. Cell images can be allocated to the corresponding categories by the auxiliary network $Q$.

*1) Training Process:* Given cell-level images sampled from the real data distribution $x \sim \mathbb{P}_r$, the first goal is to learn a generator distribution $\mathbb{P}_g$ that matches the real data distribution $\mathbb{P}_r$.

We first define a random noise variable $z$. The input noise $z$ is transformed by the generator into a sample $\tilde{x} = G(z), z \sim p(z)$. $\tilde{x}$ can be viewed as following the generator distribution $\mathbb{P}_g$. Inspired by WGAN [11], we optimize networks through the WGAN objective $W(\mathbb{P}_r, \mathbb{P}_g)$:

$$W(\mathbb{P}_r, \mathbb{P}_g) = \sup_{\|f\|_{L \leq 1}} \mathbb{E}_{x \sim \mathbb{P}_r}[f(x)] - \mathbb{E}_{\tilde{x} \sim \mathbb{P}_g}[f(\tilde{x})]. \quad (1)$$

$W(\mathbb{P}_r, \mathbb{P}_g)$ is an efficient approximation of the EM distance, which is constructed using the Kantorovich-Rubinstein duality [11]. The EM distance measures how close the generator distribution and the data distribution are. To distinguish two distributions $\mathbb{P}_g$ and $\mathbb{P}_r$, the adversarial discriminator network $D$ is trained to learn the function $f$ that maximizes $W(\mathbb{P}_r, \mathbb{P}_g)$. To make $\mathbb{P}_g$ approach $\mathbb{P}_r$, the generator instead is trained to minimize $W(\mathbb{P}_r, \mathbb{P}_g)$. The value function $V(D, G)$ is written as follows:

$$V(D, G) = \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] - \mathbb{E}_{z \sim p(z)}[D(G(z))]. \quad (2)$$

This minimax game between the generator and the discriminator is written as:

$$\min_G \max_{D \in \mathcal{D}} V(D, G). \quad (3)$$

Followed by the work of WGAN-GP [12], a gradient penalty is added on the discriminator to enforce the Lipschitz constraint to make sure that the discriminator lies within the space of 1-Lipschitz functions $D \in \mathcal{D}$. The loss of the discriminator with a hyperparameter $\lambda_1$ is written as:

$$L_D = \mathbb{E}_{z \sim p(z)}[D(G(z))] - \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] + \lambda_1 \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}}[\|\nabla_{\hat{x}} D(\hat{x})\|_p - 1]^2, \quad (4)$$

where $\mathbb{P}_{\hat{x}}$ is defined sampling uniformly along straight lines between pairs of points sampled from the data distribution $\mathbb{P}_r$ and the generator distribution $\mathbb{P}_g$.

In this way, our model is capable of generating visually appealing cell-level images. But still, it fails to exploit information of categories of cells since the noise variable $z$ doesn't correspond to any interpretable feature. Motivated by this, our second goal is to make the chosen variables represent meaningful and interpretable semantic features of cells. Inspired by InfoGAN [14], we introduce mutual information into our model:

$$I(X; Y) = H(X) - H(X|Y) = H(Y) - H(Y|X). \quad (5)$$

$I(X; Y)$ describes the dependencies between two separate variables $X$ and $Y$. It measures the different aspects of
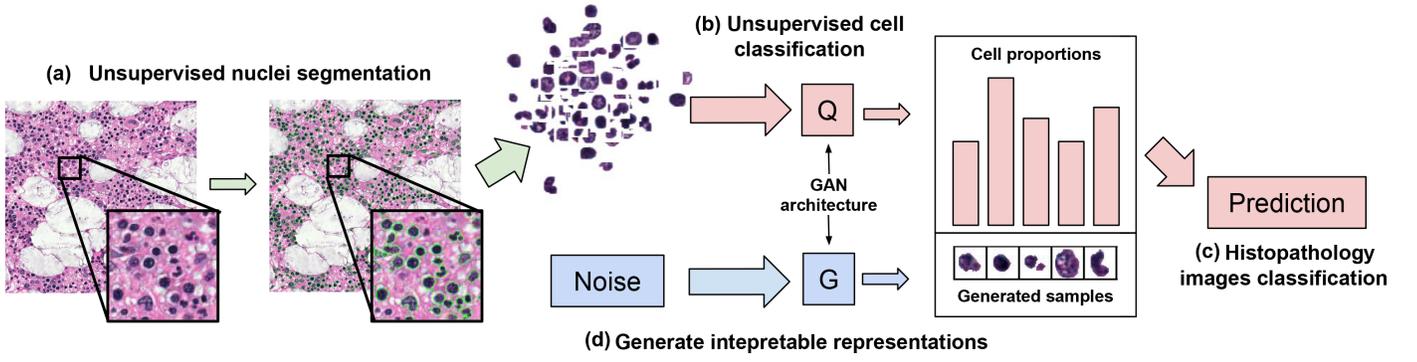
Fig. 4. Overview of our pipeline as follows: (a) Nuclei segmentation is performed on histopathology images. (b) Using the trained GAN architecture, Cell-level clustering is performed using the learned auxiliary network $Q$. Cell proportions are then calculated for each histopathology image. (c) Image-level prediction is given based on cell proportions. (d) For visualization, the generator $G$ can generate the interpretable representation for each category of cells by changing the noises.

the association between two random variables. If the chosen random variables correspond to certain semantic features, it's reasonable to assume that mutual information between generated samples and random variables should be high.

We define a latent variable $c$ sampled from a fixed noise distribution $p(c)$. The concatenation of the random noise variable $z$ and the latent variable $c$ is then transformed by the generator G into a sample $G(z, c)$. Since we encourage the latent variable to correspond with meaningful semantic features, there should be high mutual information between $c$ and $G(z, c)$. Therefore, the next step is to maximize mutual information $I(c; G(z, c))$, which can be written as:

$$I(c; G(z, c)) = H(c) - H(c|G(z, c)). \qquad (6)$$

Followed by this, a lower bound $L_I$ is given by:

$$L_I(G, Q) = \mathbb{E}_{z \sim p(z), c \sim p(c)}[\log Q(c|G(z, c))] + H(c), \quad (7)$$

where $H(c)$ is the entropy of the variable sampled from a fixed noise distribution. Maximizing this lower bound, we maximize mutual information $I(c; G(z, c))$. The proof can be found in InfoGAN [14].

Since we introduce the latent variable $c$ into the model, the value function $V(D, G)$ is replaced by:

$$V(D, G) \leftarrow \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] - \mathbb{E}_{z \sim p(z), c \sim p(c)}[D(G(z, c))]. \ (8)$$

As we combine the adversarial process with the process of maximizing mutual information, this information-regularized minimax game with a hyperparameter $\lambda_2$ can be written as follows:

$$\min_{G, Q} \max_{D \in \mathcal{D}} V(D, G) - \lambda_2 L_I(G, Q). \qquad (9)$$

The loss of $D$ can be replaced by:

$$L_D \leftarrow \mathbb{E}_{z \sim p(z), c \sim p(c)}[D(G(z, c))] - \mathbb{E}_{x \sim \mathbb{P}_r}[D(x)] + \lambda_1 \mathbb{E}_{\hat{x} \sim \mathbb{P}_{\hat{x}}}[||\nabla_{\hat{x}} D(\hat{x})||_p - 1]^2,$$
$$(10)$$

Since $H(c)$ can be viewed as a constant, the loss of the auxiliary network $Q$ can be written as the negative log-likelihood between $Q(c|G(c, z))$ and the discrete variable $c$. The losses of $G$ and $Q$ can be interpreted as below:

$$L_G = -\mathbb{E}_{z \sim p(z), c \sim p(c)}[D(G(z, c))], \qquad (11)$$

$$L_Q = -\lambda_2 \mathbb{E}_{z \sim p(z), c \sim p(c)}[\log Q(c|G(z, c))]. \qquad (12)$$

Figure 5 shows how noises are transformed into interpretable samples during the training process.
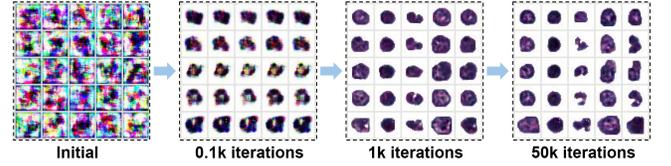


Fig. 5. Example of how a set of noise vectors are transformed into interpretable image samples over generator iterations. We use a 5-dimensional categorical variable $c$ and 32 Gaussian noise variables $z$ as input. Different rows correspond to different values of $z$. Different columns correspond to different values of $c$. The value of $c$ largely corresponds to cell types.

*2) Test Process:* In the training process, a generator distribution is learned to imitate the real data distribution. An auxiliary distribution is learned to maximize the lower bound. Especially if $c$ is sampled from a categorical distribution, a softmax function is applied as the final layer of $Q$. Under this circumstance, $Q$ can act as a classifier in the test process, since the posterior $Q(c|x)$ is discrete. Assuming that each category in $c$ corresponds to a type of cells, the auxiliary network $Q$ can divide cell-level images into different categories while the generator $G$ can generate the interpretable representation for each category of cells.

### B. Image-level Classification

Based on the cell-level visual representation learning, we propose a pipeline combining nuclei segmentation and cell-level visual representation to highlight the varieties of cellular elements. Image-level classification is performed using the calculated cell proportions. The illustration of this pipeline is shown in Figure 4.

*1) Nuclei Segmentation:* An unsupervised nuclei segmentation approach is ultilized consisting of four stages: normalization, unsupervised color deconvolution, intensity thresholding
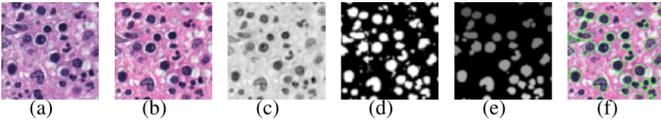
Fig. 6. Overview of segementation process: ($a$) the cropped image, ($b$) the normalized image, ($c$) the separated hematoxylin stain image using color deconvolution, ($d$) the binary image generated by intensity thresholding, ($e$) the labeled image after postprocessing where different grayscale values stand for different segmented instances, and ($f$) the final segmentation image.

and postprocessing to segment nuclei from the background. Figure 6 is an overview of our segmentation pipeline.

**Color Normalization:** We employ Reinhard color normalization [57] to convert the color characteristics of all images into the desired standard by computing the mean and standard deviations of a target image in LAB space.

**Color Deconvolution:** Using the PCA-based 'Macenko' method [58], unsupervised color deconvolution is performed to separate the normalized image into two stains. We project pixels onto a best-fit plane, wherein it selects the stain vectors as percentiles in the 'angle distribution' of the corresponding plane. With the correct stain matrix for color deconvolution, the normalized image can be separated into hematoxylin stain and eosin stain.

**Intensity Thresholding:** To sufficiently segment cells, we apply intensity thresholding in the hematoxylin stain image where the intensity distribution of cells is consistently distinct from the background. By converting the hematoxylin stain image into a binary image with a constant global threshold, the cells are roughly segmented.

**Postprocessing:** In image postprocessing, objects with fewer pixels than the minimum area threshold will be removed from the binary image. Then we employ the method in [44] to remove thin protrusions from cells. Furthermore, we use opening operation to separate a few touched cells.

*2) Classification:* We utilize the model distribution trained in our unsupervised representation learning as the cell-level classifier. Assuming that we use a $k$-dimensional categorical variable as the chosen variable in the training process, the real data (cell-level images) distribution is allocated into $k$ dimensions. In the test process, cell-level images are unsupervised classified into $k$ corresponding categories.

For each histopathology image, we count the numbers of cell-level instances in each category as the representation of its cellular constitution, denoted as $\{X_1, X_2, X_3, \ldots, X_k\}$. For cellular element $i$, the ratio of the number of this cellular element to the total number of the cellular constitution in this image is calculated by $P_i = \frac{X_i}{\sum_{i=1}^{k} X_i}$. We define $P_i$ as the cell proportion of cellular element $i$.

Given cell proportions $\{P_1, P_2, P_3, \ldots, P_k\}$ as the feature vector of histopathology images, we utilize either k-means or SVM to give image-level predictions.

## IV. EXPERIMENTS AND RESULTS

### A. Dataset

All our experiments are conducted on bone marrow histopathology images stained with hematoxylin and eosin. As

described before, the cellular constitution in bone marrow is a determinant in diagnoses of blood disease.

**Dataset A:** Publicly available dataset [59] which consists of eleven images of healthy bone marrow with a resolution of $1200 \times 1200$ pixels. Each image contains around 200 cells. The whole dataset includes 1995 cell-level images in total. We label all cell-level images into four categories: 34 neutrophils, 751 myeloblasts, 495 monocytes, and 715 lymphocytes. Images are carefully labeled by two pathologists. When the two pathologists disagree on a particular image, a senior pathologist makes a decision over the discord.

**Dataset B:** Dataset provided by the First Affiliated Hospital of Zhejiang University which contains whole slides of bone marrow from 24 patients with blood diseases. Each patient matchs with one whole slide. We randomly crop 29 images with a resolution of $1500 \times 800$ pixels from all whole slides. Dataset B contains around 12000 cells in total. For this dataset, we label 600 cell-level images into three categories for evaluation: 200 myeloblasts, 200 monocytes, and 200 lymphocytes. The labeling process is conducted in the same manner as Dataset A.

**Dataset C:** Combination of Datasets A and B, which results in 29 abnormal and 11 normal histopathology images.

**Dataset D:** Dataset includes whole slides from 28 patients with bone marrow hematopoietic tissue hyperplasia (negative) and 56 patients with leukemia (positive). Each patient matchs with one whole slide. We randomly crop images with a resolution of $1500 \times 800$ pixels from all whole slides. This results in 72 negative and 132 positive images. After segmentation, Dataset D contains around 80000 cells in total.

### B. Implementation

**Network Parameters:** Our generator $G$, discriminator $D$ and auxiliary network $Q$ all have the structures of residual networks. In the training process, all three networks are updated by Adam optimizer ($\alpha = 0.0001$, $\beta_1 = 0.5$, $\beta_2 = 0.9$, $lr = 2 \times 10^{-4}$) [61] with a batch size of 64. All our experiments use hyperparameters $\lambda_1 = 10$ and $\lambda_2 = 1$. For each training iteration, we update $D$, $G$ and $Q$ in turn. One training iteration consists of five discriminator iterations, one generator iteration, and one auxiliary network iteration. For each training process, we augment the training set by rotating images with angles $90°$, $180°$, $270°$. We train ten epochs for our model in each experiment.

**Noise Sources:** The noise fed into the network is the combination of a 5-dimensional categorical variable and 32 Gaussian noise variables for the training of Dataset A or Dataset B. We use the combination of a 5-dimensional categorical variable and 64 Gaussian noise variables for Dataset C.

**Segmentation Parameters:** The mean value of the standard image in three channels is $[8.98 \pm 0.64, 0.08 \pm 0.11, 0.02 \pm 0.03]$ for color normalization. Vectors for color deconvolution are picked from 1% to 99% angle distribution while the magnitude below 16 is excluded from the computation. We use the threshold value of 120 for intensity thresholding. In the post-process, objects with pixels smaller than 200 will be removed. An opening operation with $7 \times 7$ kernel size is performed to

separate touched cells. When the edge of the bounding box of a cell-level image is larger than 32 pixels, we rescale the image to make the larger edge match to 32. Each cell is centered in a $32 \times 32$ pixel image where blank is filled with $[255, 255, 255]$.

**Bounding Box:** To prevent the color and texture contrast from troubling the feature extraction process, we use instances without segmentation for baseline methods. If we depose the nuclei in the center with the loose bounding box in the same manner as our previous experiments, cells will suffer from severe overlapping. Thus, we crop the minimum bounding box region along each segmented instance, and then resize it into $32 \times 32$ pixels as our dataset.

**Software:** We implement our experiments on framework Pytorch for deep learning models and framework HistomicsTK for nuclei segmentation. Our model is compared with multiple sources of baselines. Three main types of baselines are claimed to be relevant as follows: (1) feature extractors including manual features, HOG and DNN extractor; (2) supervised classifiers including SVM and DNN; (3) clustering algorithms including DEC and K-means. The rich mix of different sources of baselines, including deep learning algorithms, provides a stronger demonstration to our experiments. We utilize k-means++ [60] to choose the initial values when using k-means to perform clustering. The feature code[1] is Python implementation in all these algorithms.

**Hardware:** For hardware, we use one pair of Tesla K80 GPU for parallel training and testing of neural network models. Other baseline experiments are conducted on Intel(R) Xeon(R) CPU E5-2690 v3 @ 2.60GHz. For our model, with a batch size of 64, using one pair of K80 GPU for parallel computation, each generator iteration costs 3.2 seconds in the training process when each batch costs 0.18 seconds in the test process.

### C. Cell-level Classification Using Various Features

To demonstrate the quality of our representation learning, we apply the trained model as a feature extractor. The experiment is conducted on Dataset A. In this experiment, 1596 cell-level images are used for training; 399 cell-level images are used for testing.

**Comparison:** (1) MF: 188-dimensional manual feature combined of SIFT [62], LBP [63], and $L \times a \times b$ color histogram. (2) DNN: DNN+k-means: DNN features extracted by ResNet-50 trained on Imagenet-1K, on top of which k-means is performed. (3) Our Method: We downsample the features after each residual block of the discriminator into a $4 \times 4$ spatial grid using max pooling. These features are flattened and concatenated to form an 8192-dimensional vector. On top of the feature vectors, an L2-SVM is trained to perform classification.

Different processing strategies are used as follows: (1) w/ Seg: using the output generated by nuclei segmentation; (2) w/o Seg: using the minimum bounding box along each cell-level instance.

**Evaluation:** For each class, we denote the number of true positives $TP$, the number of false positives $FP$ and the

[1]Implementation details can be found at https://github.com/bohu615/nu_gan

number of false negatives $FN$. The precision, recall and F-score ($F_1$) for each class are defined as follows:

$$precision = \frac{TP}{TP + FP},$$
$$recall = \frac{TP}{TP + FN}, \qquad (13)$$
$$F_1 = \frac{2 \cdot precision \cdot recall}{precision + recall}.$$

The average precision, recall and F-socre are calculated weighted by support (the number of true instances of each class).

**Results:** We randomly choose correctly classified and mis-classified samples displayed in Figure 7. The comparison of results is shown as Table I, which proves the advantages of our representation learning method. The manual feature extractor can generate a better result based on the bounding box regions, but its performance is still lower than ours. The color of the background can provide useful information for the color histogram channel in manual features but is viewed as noise for the DNN based extractor. Though the dimensions of the feature vectors of our method are higher, the clustering ability of our model ensures further unsupervised applications. Furthermore, we apply mean pooling on top of feature maps to prove that using less dimensional features can also generate a comparable result. In this manner, we achieve 0.850 F-score using 2048 dimensional features and 0.840 F-score using 512 dimensional features.



Fig. 7. Visualization of cell-level classification performed on Dataset A: ($up$) correctly classified samples and ($down$) misclassified samples. misclassified samples can be illegible for pathologists either.

TABLE I.     Performance of cell-level classification using various features.

| Methods | Precision | | Recall | | F-score | |
|---|---|---|---|---|---|---|
| | w/ Seg | w/o Seg | w/ Seg | w/o Seg | w/ Seg | w/o Seg |
| MF | 0.821 | 0.837 | 0.803 | 0.847 | 0.811 | 0.842 |
| DNN | 0.838 | 0.760 | 0.817 | 0.769 | 0.827 | 0.764 |
| Our Method | **0.865** | / | **0.848** | / | **0.857** | / |

### D. Cell-level Clustering

As the priority of image-level classification of histopathology images, cell-level clustering is performed using the trained auxiliary network $Q$. We conduct experiments on the three datasets described in Section IV-A.

**Comparison:** (1) MF+k-means: Manual features with k-means. (2) DNN+k-means: DNN features extracted by ResNet-50 trained on Imagenet-1K, on top of which k-means is performed. (3) HOG+DEC: Deep Embedded Clustering (DEC) [26] on 2048-dimensional HOG features. (4) Our Method: Cell images are unsupervised allocated to five clusters by the auxiliary network $Q$. We also test models such as Categorical

GAN (CatGAN) [33], InfoGAN (under DCGAN architecture), and Gaussian Mixture VAE (GMVAE) [31] on our datasets under different hyperparameters, but find them fail to converge.

The following processing strategies are also used: (1) w/ Seg: using the output generated by nuclei segmentation; (2) w/o Seg: using the minimum bounding box along each cell-level instance.

**Evaluation:** We evaluate the performance of clustering using the average F-score, purity, and entropy. For the set of clusters $\{\omega_1, \omega_2, \ldots, \omega_K\}$ and the set of classes $\{c_1, c_2, \ldots, c_J\}$, we assume that each cluster $\omega_k$ is assigned to only one class $\arg\max_j(|\omega_k \cap c_j|)$. The F-score for class $c_j$ is then given by Equation 13. The average F-score is given calculated by the number of true instances in each class.

Purity and Entropy are also used as evaluation metrics, which are written as follows:

$$purity = \frac{1}{N} \sum_k \max_j |\omega_k \cap c_j|,$$
$$entropy = -\frac{1}{N} \sum_k |\omega_k| \log \frac{|\omega_k|}{N}. \quad (14)$$

Larger purity and smaller entropy indicate better clustering results.

For nuclei segmentation, we use Intersection over Union (IoU) and the F-score as evaluation metrics. A segmented instance (I) is matched with the ground truth (G) only if they intersect at least 50% (i.e., $|I \cap G| > 0.5G$). For each matched instance and its ground truth, the overlapping pixels are counted as true positive ($TP$). The pixels of instance remain unmatched are counted as false positive ($FP$) while the pixels of ground truth remaining unmatched are counted as false negative ($FN$). The F-score is then calculated using Equation 13.

For k-means based methods, the average F-score is approximately the same ($\pm 0.02$) using either four, five, or six clusters.

**Annotations:** To evaluate the capability of nuclei segmentation, We randomly choose 20 patches from Dataset C with a resolution of $200 \times 200$ pixels. The ground truth is carefully labeled by two pathologists. When the two pathologists disagree on a particular image, a senior pathologist makes a decision over the discord.

**Results:** For nuclei segmentation, our method achieves 0.56 mean IoU and 0.70 F-score.

For cell-level clustering, the comparison shown as Table II shows the superiority of our method. To explicitly reveal the semantic features our model has captured, we randomly choose 60 samples from each of the five clusters displayed in Figure 8, which shows a distinct consistency within each cluster. Reasonable interpretations can be given. Cells are clustered according to the semantic features such as the chromatin openness, the darkness and density of nuclei, and if nuclei show the appearance of being segmented.

When it comes to unsupervised classification, none of the baseline methods can benefit from the bounding box. We observe that the color context of the background can be disturbing when the classification is under the fully unsupervised manner.

TABLE II. Performance of cell-level clustering.

| Dataset | Methods | Purity | | Entropy | | F-score | |
|---|---|---|---|---|---|---|---|
| | | w/ Seg | w/o Seg | w/ Seg | w/o Seg | w/ Seg | w/o Seg |
| A | MF+k-means | 0.579 | 0.442 | 1.376 | 1.598 | 0.603 | 0.510 |
| | DNN+k-means | 0.667 | 0.470 | 1.256 | 1.552 | 0.677 | 0.501 |
| | HOG+DEC | 0.729 | 0.637 | 1.086 | 1.167 | 0.737 | 0.664 |
| | Our Method | **0.855** | / | **0.750** | / | **0.863** | / |
| B | MF+k-means | 0.392 | 0.421 | 1.561 | 1.545 | 0.409 | 0.454 |
| | DNN+k-means | 0.719 | 0.406 | 0.844 | 1.557 | 0.760 | 0.435 |
| | HOG+DEC | 0.771 | 0.681 | 0.697 | 1.161 | 0.812 | 0.693 |
| | Our Method | **0.874** | / | **0.431** | / | **0.841** | / |
| C | MF+k-means | 0.459 | 0.446 | 1.533 | 1.597 | 0.484 | 0.514 |
| | DNN+k-means | 0.578 | 0.458 | 1.377 | 1.575 | 0.601 | 0.485 |
| | HOG+DEC | 0.667 | 0.602 | 1.217 | 1.334 | 0.682 | 0.621 |
| | Our Method | **0.769** | / | **0.977** | / | **0.777** | / |

Especially for Dataset A, Figure 9(a) shows the convergence of $V(D, G)$ (see Equation (8)) and $L_Q$ (see Equation (12)). $V(D, G)$ is used to evaluate how well the generator distribution matches the real data distribution [12]. $L_Q$ approaching zero indicates that mutual information is maximized [14]. Figure 9(b) shows how the purity of clustering increases in the training process.
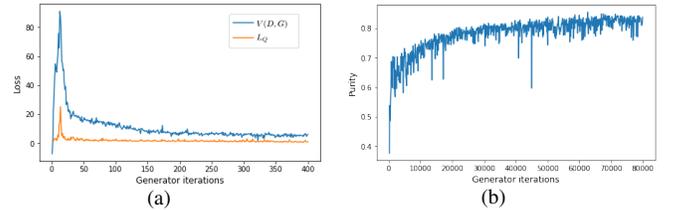


Fig. 9. Visualization of cell-level clustering performed on Dataset A: (a) Training losses converge as the network trains. (b) The purity increases gradually over generator iterations.

**Impacts of the Number of Clusters:** For our method, it is easy to change the number of clusters by sampling the categorical noise from a different dimension. We compare the results of choosing different numbers of clusters shown in Table III, which shows there is no distinct difference between choosing four and five clusters. We choose five clusters (a 5-dimensional categorical random variable) in change for a slightly better performance.

TABLE III. Performance when choosing different numbers of clusters.

| Clusters | 4 | 5 | 6 |
|---|---|---|---|
| F-score | 0.831 | 0.863 | 0.789 |

**Impacts of Uninformative Representations**: The uninformative representations such as the staining color and rotations can be interference factors in the process of classification. Besides using color normalization and data augmentation to ease this problem, we also demonstrate that these features are more likely to be latent encoded in Gaussian random variables which do not influence the classification task. As is shown in Figure 10, we fix the value of the chosen categorical variable $c$ while walking through the random space of the Gaussian noise variable $z$. The result shows that uninformative representations tend to be encoded in noise variables through the process of maximizing the mutual information.
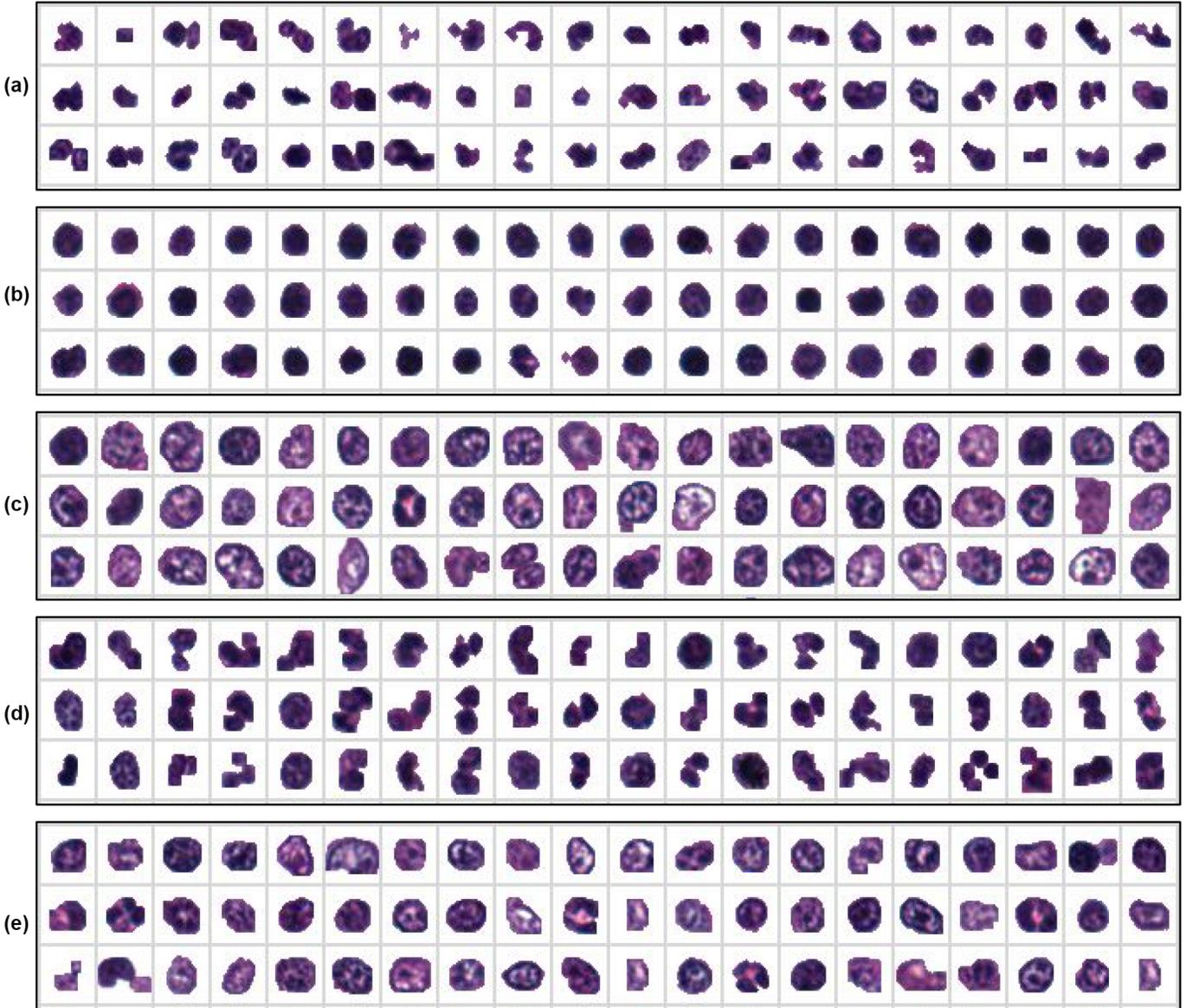
Fig. 8. Visualization of clustering. We randomly select 60 samples from each one of five clusters, displayed as (a) to (e). Instances in the same cluster have a distinct consistency. In (b), cells in marrow with dark, dense, and close phased nuclei tend to be lymphocytes or erythroid precursors. In (c) and (e), cells with dispersed chromatin are most likely granulocytes precursors such as myeloblasts.
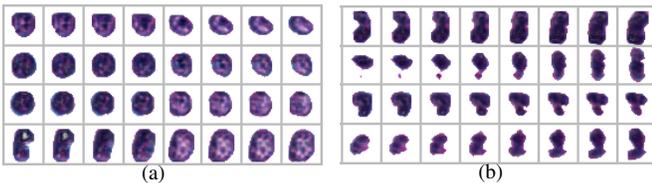


Fig. 10. Examples of how uninformative representations are encoded in Gaussian noise variables $z$. Different columns share the same value of the chosen categorical variable $c$. A random walk is performed between two points in the space of $z$. It can be seen that (a) the staining color and (b) the rotation are both latent encoded in the Gaussian noise variables.

### E. Image-level Classification

We perform image-level classification experiments on Dataset C and Dataset D respectively. Dataset C includes 29 positive and 11 negative images. Dataset D includes 132 positive and 72 negative images. Each dataset is randomly split into four folds for the 4-fold cross-validation. Each score is reported averagely. Each experiment is repeated for four times with different random split for cross-validation. The scores are reported four times to show confidence intervals.

**Comparison:** (1) DNN (cell-level based): We use ResNet-50 features extracted from cell-level instances to perform

cell-level clustering. Then we train an L2-SVM on top of the cell proportions to perform image-level classification. (2) DNN (image-level based): We use ResNet-50 pre-trained on Imagenet-1K as an image-level feature extractor. Images with a resolution of $1500 \times 800$ are normalized and center cropped to $800 \times 800$ pixels, then resized into $224 \times 224$ pixels. An L2-SVM is trained on the feature vectors. We observe this produces a better result than fine-tuning or directly training a ResNet-50 without pre-train. (3) Our method (w/ k-means): We first train our GAN architecture on the training set, then conduct the cell-level clustering on both the training set and test set using the trained model. Cluster centers are calculated given cell proportions of each sample in the training set. The predict label is given by the closest cluster that each sample in the test set belongs to. (4) Our method (w/ SVM): An L2-SVM instead of k-means is used as the final classifier.

**Evaluation:** We use the precision, recall and F-score for evaluation, the details of which have been described in Equation 13. The difference is that the labels are binary in this experiment.

**Results:** Following the proposed pipeline, the GAN architecture is trained on the segmentation output of the split training set. For cell-level clustering task, we achieve 0.791 F-score trained on 12000 training instances of Dataset C and 0.771 F-score trained on 60000 training instances of Dataset D, both evaluated by labeled cells of Dataset A.

Given the cell proportions, when using k-means to perform image-level unsupervised classification, we achieve 0.931 F-score on Dataset C and 0.875 F-score on Dataset D, which is comparative to the DNN method with 0.933 and 0.888 F-score. The advantage is that our model is interpretable. The proportion of which category of cells is irregular is recognizable.

Since there are a large number of cell-level images on both Dataset C and D, it is difficult to test our method under full-supervision with a similar pipeline. We instead train an L2-SVM on cell proportions, taking image-level labels of histopathology images as targets. As the comparison shown in Table IV, our method achieves 0.950 F-score on Dataset C and 0.902 F-score on Dataset D.

On Dataset C, we use Principal Components Analysis (PCA) to perform a dimensionality reduction, cell proportions of each histopathology image are projected onto a two-dimension plane to show that there is a distinct difference between normal and abnormal images, shown in Figure 11.

**Impacts of the Segmentation Parameters:** To validate the impacts of the segmentation performance on the image-level classification result, we change the value of intensity threshold in the segmentation process of experiments on Dataset C. We randomly choose 20 patches with a resolution of $200 \times 200$ pixels in Dataset C for evaluation, which includes 335 nuclei as counted. We use missing instances (nuclei that are missing in outputs), false alarms (mis-segmented background instances), and the F-score for evaluation.

As is shown in Table V, both results of segmentation and classification are the highest when the intensity threshold remains 120. Followed by the decreasing of segmentation performance, the classification performance will stay within
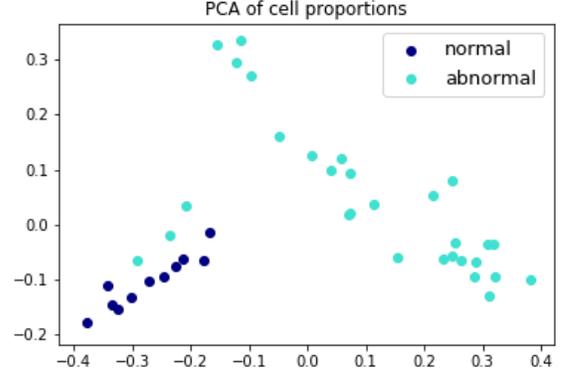


Fig. 11. Visualization of unsupervised classification using cell proportions. It can be observed that the points representing normal and abnormal samples are distinctly distributed in two different clusters.

an acceptable range. Too bad segmentation performance will worsen the classification result since the quality and quantity of the segmentation outputs are not enough to reveal the distinct representation of each image-level instance.

TABLE V.       Performance when changing the segmentation parameters.

| Intensity threshold | 60 | 80 | 100 | 120 | 140 | 160 | 180 |
|---|---|---|---|---|---|---|---|
| Missing Instances | 127 | 48 | 21 | 7 | 14 | 64 | 184 |
| False Alarms | 3 | 4 | 15 | 5 | 20 | 30 | 35 |
| Segmentation F-score | 0.315 | 0.413 | 0.602 | 0.701 | 0.656 | 0.534 | 0.218 |
| Classification F-score | 0.579 | 0.814 | 0.932 | 0.950 | 0.941 | 0.901 | 0.576 |

**Impacts of the Number of Clusters:** For image-level classification of Dataset C, we conduct experiments choosing different number of clusters. Table VI shows that there is no distinct difference of performance between choosing five and six clusters. We still choose five clusters for a better performance.

TABLE VI.       Performance when choosing different numbers of clusters.

| Clusters | 4 | 5 | 6 | 7 |
|---|---|---|---|---|
| Cell-level Classification F-score | 0.711 | 0.791 | 0.762 | 0.710 |
| Image-level Classification F-score | 0.897 | 0.950 | 0.944 | 0.899 |

**Patch-level Classification**: We perform classification based on patches. Using a sliding window with a window size of 224 and a stride of 224, we separately transfer the normalized images from the training set and test set from Dataset C into labeled image patches. This results in 588 positive and 288 negative patches for training, 224 positive and 108 negative patches for testing. If 50% of the patches of an image-level instance are positive, we will consider this instance as positive. In this manner, we achieve 0.851 F-score using DNN feature extractor with SVM and 0.831 F-score using our method, which is not comparative to our image-level classification results.

**Discussion:** Analyzing the results, we find that the cell proportions $\{P_1, P_2, \cdots, P_5\}$ can indicate the presence of blood diseases.

TABLE IV. Performance of image-level classification. Each experiment is repeated for four times with different random split for cross-validation. The scores are reported four times to show confidence intervals.

| Datasets | Methods | Precision | | | | Recall | | | | F-score | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| C | DNN (cell-level based) | 0.539 | 0.598 | 0.688 | 0.524 | 0.711 | 0.723 | 0.734 | 0.678 | 0.636 | 0.678 | 0.701 | 0.621 |
| | DNN (image-level based) | 0.906 | 0.913 | 0.901 | 0.921 | **0.969** | 0.958 | 0.943 | 0.965 | 0.933 | 0.929 | 0.924 | 0.937 |
| | Our Method (w/ k-means) | 0.936 | 0.945 | 0.939 | 0.937 | 0.933 | 0.944 | 0.946 | 0.938 | 0.931 | 0.941 | 0.948 | 0.939 |
| | Our Method (w/ SVM) | **0.950** | **0.948** | **0.940** | **0.946** | **0.969** | **0.968** | **0.950** | **0.966** | **0.950** | **0.949** | **0.940** | **0.949** |
| D | DNN (cell-level based) | 0.469 | 0.579 | 0.498 | 0.581 | 0.697 | 0.654 | 0.643 | 0.665 | 0.558 | 0.612 | 0.583 | 0.621 |
| | DNN (image-level based) | 0.863 | **0.900** | 0.887 | 0.869 | **0.863** | 0.886 | 0.871 | 0.865 | **0.863** | 0.888 | 0.879 | 0.866 |
| | Our Method (w/ k-means) | 0.858 | 0.879 | 0.881 | 0.868 | 0.857 | 0.868 | 0.873 | 0.865 | 0.862 | 0.870 | 0.875 | 0.867 |
| | Our Method (w/ SVM) | **0.864** | 0.897 | **0.901** | **0.882** | 0.858 | **0.892** | **0.898** | **0.878** | **0.863** | **0.891** | **0.902** | **0.880** |

For our experiment, cell-level clustering shows that $\{P_1, P_4\}$ correspond to myeloblasts, $\{P_5\}$ corresponds to lymphocytes and erythroid precursors, and $\{P_2, P_3\}$ correspond to monocytes and glanulocytes. For all normal images, $P_1$ and $P_4$ are relatively lower. This matches the constitution in normal bone marrow where the lymphocytes, glanulocytes and erythroid precursors are in the majority when the percentage of cells with open phased nuclei (such as myeloblasts, under some circumstances plasma cells) is relatively lower (less than 10%). In Figure 11, abnormal images that are confidently discriminated are reflected in the numerous presence of the supposed minority myeloblasts or plasma cells, which in turn is reflected in the sharp increase of $P_1$ and $P_4$.

However, there are three abnormal images that are exceptional. To analyze what causes the failure, we display the example image in Figure 12.
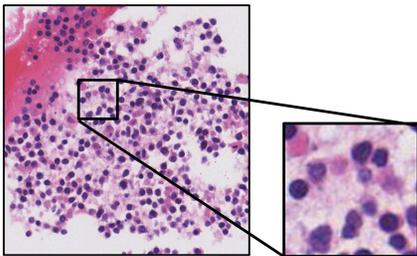


Fig. 12. Example of the failed samples. Too many erythroid precursors indicate the presence of blood disease. The overlap of nuclei and the lousy staining condition add to the difficulties of cell-level classification.

In these images, the irregular proportion of erythroid precursors indicates the presence of blood disease. We find that our model does not correctly classify these cells. The reason could be that the staining condition of these cells is not as good as expected. A typical erythroid precursor should have a close phased, dark-staining nucleus that appears almost black. As Figure 13 shows, the color of nuclei segmented from these images differ from the rest of the dataset. Particularly in these images, our model is still not robust enough to capture the most significant semantic variance in an unsupervised setting. Therefore, acquiring high-quality histopathology images is still a priority.
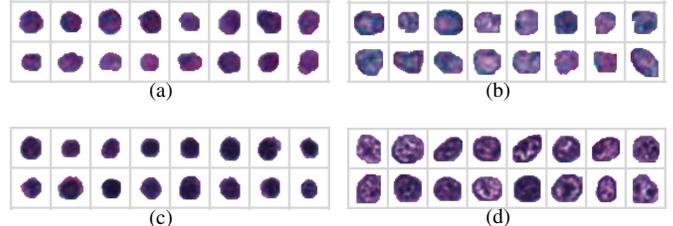


Fig. 13. Variance of staining conditions. $(a)$ and $(b)$ are erythroid precursors and myeloblasts randomly chosen from failed images. $(c)$ and $(d)$ are samples selected from correctly predicted images. Our model mistakes erythroid precursors for myeloblasts particularly in failed images.

## V. CONCLUSION

In this paper, we introduce a unified GAN architecture with a new formulation of the loss function into cell-level visual representation learning of histopathology images. Cell-level unsupervised classification with interpretable visualization is performed by maximizing mutual information. Based on this model, we exploit cell-level information by calculating the cell proportions of histopathology images. Followed by this, we propose a novel pipeline combining cell-level visual representation learning and nuclei segmentation to highlight the varieties of cellular elements, which achieves promising results when tested on bone marrow datasets.

In future work, some improvements can be made to our method. First, the segmentation method and the computational time can be further improved. The gradient penalty added on the network architecture requires the computation of the second order derivative, which is time-consuming in the training process. Secondly, in addition to cell proportions, other information about the patients should be carefully considered, such as clinical trials and gene expression data. By allocating and annotating the relevant genetic variants, the risk can be re-evaluated. In clinical practice, doctors need to consolidate more critical information to make a confident diagnosis. For example, bone marrow cells of children might not be as varied as those of adults'. To classify cells in a more fine-grained manner, the peculiar distribution information such as erythroid cells more likely form clusters (erythroid islands) can be considered.

## VI. ACKNOWLEDGMENT

## REFERENCES

[1] M. N. Gurcan, L. E. Boucheron, A. Can, A. Madabhushi, N. M. Rajpoot, and B. Yener, "Histopathological image analysis: A review," *IEEE Reviews in Biomedical Engineering*, vol. 2, no. 1, pp. 147–171, 2009.

[2] J. M. Bennett, D. Catovsky, M. T. Daniel, G. Flandrin, D. A. Galton, H. R. Gralnick, and C. Sultan, "Proposals for the classification of the acute leukaemias. french-american-british (fab) co-operative group," *British Journal of Haematology*, vol. 33, no. 4, p. 451458, 1976.

[3] Y. Xu, T. Mo, Q. Feng, P. Zhong, M. Lai, and I. C. Chang, "Deep learning of feature representation with multiple instance learning for medical image analysis," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 1626–1630, 2014.

[4] Y. Xu, Z. Jia, Z. Ai, F. Zhang, M. Lai, I. Eric, and C. Chang, "Deep convolutional activation features for large scale brain tumor histopathology image classification and segmentation," in *IEEE International Conference on Acoustics, Speech and Signal Processing*, pp. 947–951, 2015.

[5] J. Xu, X. Luo, G. Wang, H. Gilmore, and A. Madabhushi, "A deep convolutional neural network for segmenting and classifying epithelial and stromal regions in histopathological images," *Neurocomputing*, vol. 191, no. 1, pp. 214–223, 2016.

[6] H. Chen, X. Qi, L. Yu, and P.-A. Heng, "Dcan: Deep contour-aware networks for accurate gland segmentation," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 2487–2496, 2016.

[7] T. Chen and C. Chefdhotel, "Deep learning based automatic immune cell detection for immunohistochemistry images," in *International Workshop on Machine Learning in Medical Imaging*, pp. 17–24, 2014.

[8] D. C. Cireşan, A. Giusti, L. M. Gambardella, and J. Schmidhuber, "Mitosis detection in breast cancer histology images with deep neural networks," in *International Conference on Medical Image Computing and Computer-assisted Intervention*, pp. 411–418, 2013.

[9] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, pp. 2672–2680, 2014.

[10] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.

[11] M. Arjovsky, S. Chintala, and L. Bottou, "Wasserstein gan," *arXiv preprint arXiv:1701.07875*, 2017.

[12] I. Gulrajani, F. Ahmed, M. Arjovsky, V. Dumoulin, and A. Courville, "Improved training of wasserstein gans," *arXiv preprint arXiv:1704.00028*, 2017.

[13] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pp. 770–778, 2016.

[14] X. Chen, Y. Duan, R. Houthooft, J. Schulman, I. Sutskever, and P. Abbeel, "Infogan: Interpretable representation learning by information maximizing generative adversarial nets," in *Advances in neural information processing systems*, pp. 2172–2180, 2016.

[15] S. Nazlibilek, D. Karacor, T. Ercan, M. H. Sazli, O. Kalender, and Y. Ege, "Automatic segmentation, counting, size determination and classification of white blood cells," *Measurement*, vol. 55, no. 3, pp. 58–65, 2014.

[16] Y. Sun and P. A. Sermon, "Methods for nuclei detection, segmentation, and classification in digital histopathology: A review–current status and future potential," *IEEE Reviews in Biomedical Engineering*, vol. 7, no. 1-5, p. 97, 2014.

[17] M. Muthu, Rama Krishnan, C. Chakraborty, R. R. Paul, and A. K. Ray, "Hybrid segmentation, characterization and classification of basal cell nuclei from histopathological images of normal oral mucosa and oral submucous fibrosis," *Expert Systems with Applications*, vol. 39, no. 1, pp. 1062–1077, 2012.

[18] X. Xu, F. Lin, C. Ng, and K. P. Leong, "Dual spatial pyramid on rotation invariant texture feature for hep-2 cell classification," in *International Joint Conference on Neural Networks*, pp. 1–8, 2015.

[19] J. V. Lorenzo-Ginori, W. Curbelo-Jardines, J. D. Lpez-Cabrera, and S. B. Huergo-Surez, *Cervical Cell Classification Using Features Related to Morphometry and Texture of Nuclei*. Springer Berlin Heidelberg, 2013.

[20] M. M. Dundar, S. Badve, G. Bilgin, V. Raykar, R. Jain, O. Sertel, and M. N. Gurcan, "Computerized classification of intraductal breast lesions using histopathological images.," *IEEE Transactions on Biomedical Engineering*, vol. 58, no. 7, pp. 1977–1984, 2011.

[21] K. Nguyen, A. K. Jain, and B. Sabata, "Prostate cancer detection: Fusion of cytological and textural features," *Journal of Pathology Informatics*, vol. 2, no. 1, p. 1, 2011.

[22] W. L. Tai, R. M. Hu, C. W. H. Han, R. M. Chen, and J. J. P. Tsai, "Blood cell image classification based on hierarchical svm," in *IEEE International Symposium on Multimedia*, pp. 129–136, 2011.

[23] L. Putzu, G. Caocci, and C. D. Ruberto, "Leucocyte classification for leukaemia detection using image processing techniques," *Artificial Intelligence in Medicine*, vol. 62, no. 3, pp. 179–191, 2014.

[24] M. C. Su, C. Y. Cheng, and P. C. Wang, "A neural-network-based approach to white blood cell classification," *The scientific world journal*, vol. 2014, no. 4, p. 796371, 2014.

[25] Y. Xu, Z. Jia, L.-B. Wang, Y. Ai, F. Zhang, M. Lai, I. Eric, and C. Chang, "Large scale tissue histopathology image classification, segmentation, and visualization via deep convolutional activation features," *BMC bioinformatics*, vol. 18, no. 1, p. 281, 2017.

[26] J. Xie, R. Girshick, and A. Farhadi, "Unsupervised deep embedding for clustering analysis," in *International Conference on Machine Learning*, pp. 478–487, 2016.

[27] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.

[28] J. Xu, L. Xiang, Q. Liu, H. Gilmore, J. Wu, J. Tang, and A. Madabhushi, "Stacked sparse autoencoder (ssae) for nuclei detection on breast cancer histopathology images.," *IEEE Transactions on Medical Imaging*, vol. 35, no. 1, pp. 119–130, 2016.

[29] A. A. Cruzroa, J. E. Arevalo Ovalle, A. Madabhushi, and F. A. Gonzlez Osorio, "A deep learning architecture for image representation, visual interpretability and automated basal-cell carcinoma cancer detection.," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 403–410, 2013.

[30] X. Zhang, W. Liu, H. Dou, T. Ju, J. Xu, and S. Zhang, "Fusing heterogeneous features from stacked sparse autoencoder for histopathological image analysis,"*IEEE Journal of Biomedical and Health Informatics*, vol. 20, no. 5, pp. 1377–1383, 2016.

[31] N. Dilokthanakul, P. A. Mediano, M. Garnelo, M. C. Lee, H. Salimbeni, K. Arulkumaran, and M. Shanahan, "Deep unsupervised clustering with gaussian mixture variational autoencoders," *arXiv preprint arXiv:1611.02648*, 2016.

[32] Z. Jiang, Y. Zheng, H. Tan, B. Tang, and H. Zhou, "Variational deep embedding: An unsupervised and generative approach to clustering," in *International Joint Conference on Artificial Intelligence*, 2017.

[33] J. T. Springenberg, "Unsupervised and semi-supervised learning with categorical generative adversarial networks," *arXiv preprint arXiv:1511.06390*, 2015.

[34] C. D. Malon and C. Eric, "Classification of mitotic figures with convolutional neural networks and seeded blob features," *Journal of Pathology Informatics*, vol. 4, no. 1, p. 9, 2013.

[35] S. Mohapatra, D. Patra, and S. Satpathy, *An ensemble classifier system for early diagnosis of acute lymphoblastic leukemia in blood microscopic images*. Springer-Verlag, 2014.

[36] J. Zhao, M. Zhang, Z. Zhou, J. Chu, and F. Cao, "Automatic detection and classification of leukocytes using convolutional neural networks," *Medical & biological engineering & computing*, vol. 55, no. 8, pp. 1287–1301, 2017.

[37] K. Sirinukunwattana, S. E. Ahmed Raza, Y. W. Tsang, D. R. Snead, I. A. Cree, and N. M. Rajpoot, "Locality sensitive deep learning for detection and classification of nuclei in routine colon cancer histology images," *IEEE Transactions on Medical Imaging*, vol. 35, no. 5, p. 1196, 2016.

[38] L. Hou, K. Singh, D. Samaras, T. M. Kurc, Y. Gao, R. J. Seidman, and J. H. Saltz, "Automatic histopathology image analysis with cnns," in *Scientific Data Summit (NYSDS)*, pp. 1–6, 2016.

[39] X. Zhang, H. Su, L. Yang, and S. Zhang, "Weighted hashing with multiple cues for cell-level analysis of histopathological images," in *International Conference on Information Processing in Medical Imaging*, pp. 303–314, Springer, 2015.

[40] X. Zhang, W. Liu, M. Dundar, S. Badve, and S. Zhang, "Towards large-scale histopathological image analysis: Hashing-based image retrieval,"*IEEE Transactions on Medical Imaging*, vol. 34, no. 2, pp. 496–506, 2015.

[41] X. Zhang, F. Xing, H. Su, L. Yang, and S. Zhang, "High-throughput histopathological image analysis via robust cell segmentation and hashing,"*Medical image analysis*, vol. 26, no. 1, pp. 306–315, 2015.

[42] X. Shi, F. Xing, Y. Xie, H. Su, and L. Yang, "Cell encoding for histopathology image classification," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 30–38, Springer, 2017.

[43] C. Callau, M. Lejeune, A. Korzynska, M. Garca, G. Bueno, R. Bosch, J. Jan, G. Orero, T. Salvad, and C. Lpez, "Evaluation of cytokeratin-19 in breast cancer tissue samples: a comparison of automatic and manual evaluations of scanned tissue microarray cylinders.," *Biomedical Engineering Online*, vol. 14, no. S2, p. S2, 2015.

[44] S. Wienert, D. Heim, S. Kai, A. Stenzinger, M. Beil, P. Hufnagl, M. Dietel, C. Denkert, and F. Klauschen, "Detection and segmentation of cell nuclei in virtual microscopy images: A minimum-model approach," *Scientific Reports*, vol. 2, no. 7, p. 503, 2012.

[45] L. B. Dorini, R. Minetto, and N. J. Leite, "Semiautomatic white blood cell segmentation based on multiscale analysis.," *IEEE Journal of Biomedical and Health Informatics*, vol. 17, no. 1, p. 250, 2013.

[46] O. Schmitt and M. Hasse, "Morphological multiscale decomposition of connected regions with emphasis on cell clusters," *Computer Vision & Image Understanding*, vol. 113, no. 2, pp. 188–201, 2009.

[47] O. Dzyubachyk, W. A. van Cappellen, J. Essers, W. J. Niessen, and E. Meijering, "Advanced level-set-based cell tracking in time-lapse fluorescence microscopy.," *IEEE Transactions on Medical Imaging*, vol. 29, no. 3, pp. 852–867, 2010.

[48] F. Long, H. Peng, X. Liu, S. K. Kim, and E. Myers, "A 3d digital atlas of c. elegans and its application to single-cell analyses.," *Nature Methods*, vol. 6, no. 9, p. 667, 2009.

[49] S. Hai, F. Xing, J. D. Lee, C. A. Peterson, and Y. Lin, "Automatic myonuclear detection in isolated single muscle fibers using robust ellipse fitting and sparse representation," *IEEE/ACM Transactions on Computational Biology & Bioinformatics*, vol. 11, no. 4, pp. 714–726, 2014.

[50] G. Bueno, R. Gonzlez, O. Dniz, M. Garcarojo, J. Gonzlezgarca, M. M. Fernndezcarrobles, N. Vllez, and J. Salido, "A parallel solution for high resolution histological image analysis.," *Computer Methods & Programs in Biomedicine*, vol. 108, no. 1, pp. 388–401, 2012.

[51] H. Chang, J. Han, A. Borowsky, L. Loss, J. W. Gray, P. T. Spellman, and B. Parvin, "Invariant delineation of nuclear architecture in glioblastoma multiforme for clinical and molecular association," *IEEE Transactions on Medical Imaging*, vol. 32, no. 4, pp. 670–682, 2013.

[52] S. Arslan, T. Ersahin, R. Cetin-Atalay, and C. Gunduz-Demir, "Attributed relational graphs for cell nucleus segmentation in fluorescence microscopy images," *IEEE Transactions on Medical Imaging*, vol. 32, no. 6, pp. 1121–1131, 2013.

[53] D. Nie, R. Trullo, J. Lian, C. Petitjean, S. Ruan, Q. Wang, and D. Shen, "Medical image synthesis with context-aware generative adversarial networks," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 417–425, 2017.

[54] Z. Li, Y. Wang, and J. Yu, "Reconstruction of thin-slice medical images using generative adversarial network," in *International Workshop on Machine Learning in Medical Imaging*, pp. 325–333, Springer, 2017.

[55] D. Mahapatra, B. Bozorgtabar, S. Hewavitharanage, and R. Garnavi, "Image super resolution using generative adversarial networks and local saliency maps for retinal image analysis," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pp. 382–390, 2017.

[56] J. M. Wolterink, T. Leiner, M. A. Viergever, and I. Isgum, "Generative adversarial networks for noise reduction in low-dose ct.," *IEEE Transactions on Medical Imaging*, vol. PP, no. 99, p. 1, 2017.

[57] E. Reinhard, M. Ashikhmin, B. Gooch, and P. Shirley, "Color transfer between images," *IEEE Computer Graphics & Applications*, vol. 21, no. 5, pp. 34–41, 2001.

[58] M. Macenko, M. Niethammer, J. S. Marron, D. Borland, J. T. Woosley, X. Guan, C. Schmitt, and N. E. Thomas, "A method for normalizing histology slides for quantitative analysis," in *IEEE International Conference on Symposium on Biomedical Imaging: From Nano To Macro*, pp. 1107–1110, 2009.

[59] P. Kainz, M. Urschler, S. Schulter, P. Wohlhart, and V. Lepetit, "You should use regression to detect cells," in *International Conference on Medical Image Computing and Computer Assisted Intervention*, pp. 276–283, 2015.

[60] D. Arthur and S. Vassilvitskii, "k-means++: The advantages of careful seeding," in *Proceedings of the eighteenth annual ACM-SIAM symposium on Discrete algorithms*, pp. 1027–1035, Society for Industrial and Applied Mathematics, 2007.

[61] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[62] D. G. Lowe, "Distinctive image features from scale-invariant keypoints," *International journal of computer vision*, vol. 60, no. 2, pp. 91–110, 2004.

[63] T. Ojala, M. Pietikainen, and T. Maenpaa, "Multiresolution gray-scale and rotation invariant texture classification with local binary patterns," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 24, no. 7, pp. 971–987, 2002.