

# Accepted Manuscript

Supervised Graph Hashing for Histopathology Image Retrieval and Classification

Xiaoshuang Shi, Fuyong Xing, Kadi Xu, Yuanpu Xie, Hai Su, Lin Yang

PII: S1361-8415(17)30123-8  
DOI: [10.1016/j.media.2017.07.009](https://doi.org/10.1016/j.media.2017.07.009)  
Reference: MEDIMA 1285



To appear in: *Medical Image Analysis*

Received date: 4 November 2016  
Revised date: 25 July 2017  
Accepted date: 31 July 2017

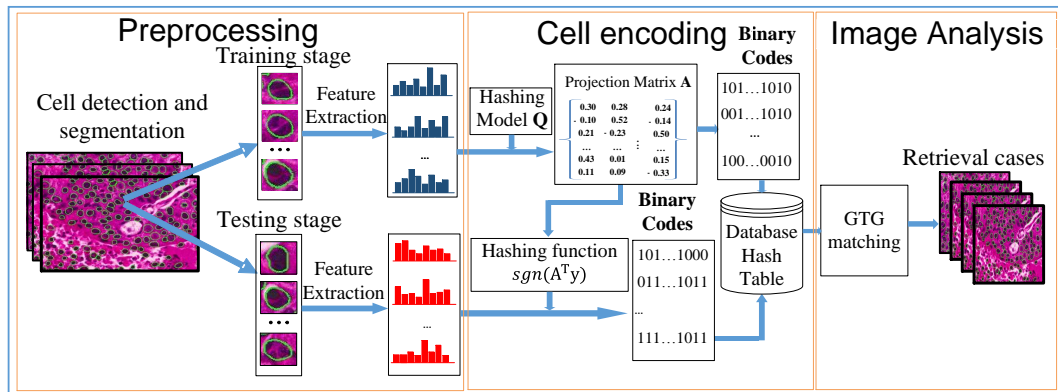
Please cite this article as: Xiaoshuang Shi, Fuyong Xing, Kadi Xu, Yuanpu Xie, Hai Su, Lin Yang, Supervised Graph Hashing for Histopathology Image Retrieval and Classification, *Medical Image Analysis* (2017), doi: [10.1016/j.media.2017.07.009](https://doi.org/10.1016/j.media.2017.07.009)

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

**Research highlights**

- (1) An framework based on cell encoding for large-scale histopathological image analysis is proposed.
- (2) A supervised graph-based model via asymmetric relaxation and its scalable version are proposed.
- (3) A group-to-group matching method to retrieve images based on binary codes of cells is proposed.

ACCEPTED MANUSCRIPT



# Supervised Graph Hashing for Histopathology Image Retrieval and Classification

Xiaoshuang Shi<sup>a</sup>, Fuyong Xing<sup>b</sup>, Kadi Xu<sup>c</sup>, Yuanpu Xie<sup>a</sup>, Hai Su<sup>a</sup>, Lin Yang<sup>a,\*</sup>

<sup>a</sup>*J. Crayton Pruitt Family Department of Biomedical Engineering  
University of Florida, Gainesville, FL, 32611-6130, U.S.A*

<sup>b</sup>*Department of Electrical and Computer Engineering,  
University of Florida, Gainesville, FL, 32611-6130, U.S.A*

<sup>c</sup>*Department of Computer Science and Engineering,  
University of Florida, Gainesville, FL, 32611-6130, U.S.A*

---

## Abstract

In pathology image analysis, morphological characteristics of cells are critical to grade many diseases. With the development of cell detection and segmentation techniques, it is possible to extract cell-level information for further analysis in pathology images. However, it is challenging to conduct efficient analysis of cell-level information on a large-scale image dataset because each image usually contains hundreds or thousands of cells. In this paper, we propose a novel image retrieval based framework for large-scale pathology image analysis. For each image, we encode each cell into binary codes to generate image representation using a novel graph based hashing model and then conduct image retrieval by applying a group-to-group matching method to similarity measurement. In order to improve both computational efficiency and memory requirement, we further introduce matrix factorization into the hashing model for scalable image retrieval. The proposed framework is extensively validated with thousands of lung cancer images, and it achieves 97.98% classification accuracy and 97.50% retrieval precision with all cells of each query image used.

*Keywords:* Image retrieval, large-scale images, hashing, histopathology image analysis.

---

\*Corresponding author

*Email address:* lin.yang@bme.ufl.edu (Lin Yang)

## 1. Introduction

Histopathology plays an important role in the early diagnosis of different cancers, such as lung and breast cancers. However, manual examination of histopathological images is labor intensive, time consuming and error-prone due to high-resolution and subjective assessment of doctors. To reduce the workload of pathologists and provide more reliable and consistent analysis of histopathological images, image process techniques and modern machine learning algorithms have been widely used for medical diagnosis, disease detection and decision support (Petushi et al., 2006), (Yang et al., 2007), (Caicedo et al., 2009), (Basavanhally et al., 2010), (Dundar et al., 2011), (Tabesh et al., 2007), (Xing and Yang, 2016). Compared to classifier-based computer-aided diagnosis (CAD) systems that directly provide diagnosis results or grading scores, content-based image retrieval (CBIR) (Comaniciu et al., 1999), (Schnorrenberg et al., 2000), (Zheng et al., 2003), (El-Naqa et al., 2004), (Liu et al., 2016) is able to provide more clinical evidence to support the diagnosis, because CBIR methods can be used to not only classify the query image but also retrieve and visualize the images with morphological profiles most relevant (Greenspan and Pinhas, 2007), (Akakin and Gurcan, 2012), (Zhang et al., 2014) (Zhang et al., 2015a), (Zhang et al., 2015b), (Jiang et al., 2015), (Zhang and Metaxas, 2016), (Jiang et al., 2016).

Cell-level information, including shape, area, and nuclear and cytoplasm appearances, plays a significant role in disease grading. In clinical practice, they do have extensive applications: 1) Two major types of non-small cell lung cancer, adenocarcinoma and squamous carcinoma, often contain cells that exhibit a mixture of representative morphologies, and the disease classification can be achieved using a majority voting of different cells. 2) Gliomas highly depends on the cellular-level information. For example, 1p/19q co-deletion greatly relies on the shape of the cell (roundness), and searching for cells that exhibit the similar morphology has important prognosis values because this would allow the clinician researchers to check whether there exists important known clinical information in the existing image database. 3) Bladder cancer is another example. Finding the nucleus with prominent nucleolus inside and their content-wise similar nucleus in the database is very important to differentiate low- and high-grade bladder cancer. Therefore, rigorously measuring and analyzing each individual cell can significantly assist

36 pathologists for diagnosis and disease detection (Zhang et al., 2015c). How-  
37 ever, it is a challenging task because there often exist hundreds of thousands  
38 of cells in one single digitized image. When using high-dimensional features  
39 for large-scale cell images, traditional CBIR methods usually exhibit low  
40 computational efficiency. As a result, most previous methods (Doyle et al.,  
41 2008), (Zhang et al., 2015a), (Jiang et al., 2015) encode the whole image as  
42 a holistic high-dimensional features by representing the statistics of cell-level  
43 information, and then compress the high-dimensional features for compu-  
44 tational efficiency. Nevertheless, some significant information might lose in  
45 such holistic representation.

46 To enable large-scale image analysis, recently hashing-based retrieval  
47 methods are considerably attractive since they can significantly improve the  
48 requirement of computer memory and query time cost (Zhang et al., 2015c),  
49 (Wang et al., 2016), thereby facilitating fast image retrieval in a large-scale  
50 database. Hashing encodes the data using a set of discrete binary codes,  
51 which can be *easily stored and quickly searched*. Thus it is often used to re-  
52 trieve nearest neighbors based on a certain similarity measurement in large-  
53 scale databases (Wang et al., 2012). Many hashing methods (Weiss et al.,  
54 2009), (Liu et al., 2012), (Shen et al., 2015), (Jiang et al., 2015), (Shi et al.,  
55 2017) have been proposed, and they have achieved great success in computer  
56 vision and data mining. In general, these methods can be roughly classified  
57 into two categories: (1) unsupervised hashing (Weiss et al., 2009), which  
58 aims to explore the intrinsic structure of data to preserve the similarity of  
59 neighbors without any supervision; (2) supervised hashing (Liu et al., 2012),  
60 (Shen et al., 2015), which utilizes the semantic information to assist search-  
61 ing and retrieval. Due to the semantic gap (Wang et al., 2012), unsupervised  
62 hashing models usually exhibit inferior performance to supervised hashing  
63 approaches. Therefore, supervised hashing methods are more preferred in  
64 histopathological image analysis.

65 Unfortunately, when hashing encodes each cell image into binary codes,  
66 there are two major concerns: (1) Directly learning discrete binary codes is an  
67 NP-hard problem (Liu et al., 2014), (Shen et al., 2015). To address this prob-  
68 lem, most of hashing methods utilize the strategy of symmetric relaxation  
69 followed by a threshold to obtain binary codes (Weiss et al., 2009). However,  
70 this strategy would generate accumulated quantization errors between the  
71 discrete binary code matrix and its relaxed continuous matrix, which would  
72 decrease the retrieval and classification accuracy (Shen et al., 2015). (2) It is  
73 a non-trivial task to effectively utilize binary codes of cells in a query image

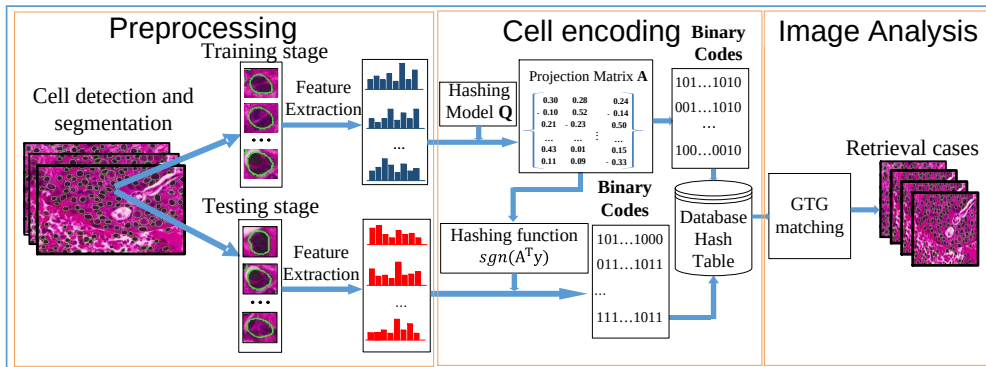


Fig 1: The flowchart of the proposed framework for pathology image analysis.

74 to retrieve relevant images. Usually, hashing methods retrieve images by a  
 75 point-to-point matching. In other words, they encode the entire images into  
 76 a set of binary codes and then calculate the Hamming distance between two  
 77 images. This strategy is not suitable for cell based medical image retrieval  
 78 since one pathology image contains more than one cell.

79 In this paper, we propose a novel framework for histopathological image  
 80 analysis via encoding large-scale cells (see Fig 1). Specifically, instead of using  
 81 symmetric relaxation followed by a threshold, we propose a supervised graph-  
 82 based hashing (GH) algorithm via asymmetric relaxation, which preserves  
 83 the hashing function in the objective function and thus effectively reduces  
 84 the accumulated quantization errors between the discrete and continuous  
 85 matrices (Shi et al., 2016b). Additionally, the GH model jointly learns binary  
 86 codes and a projection matrix, usually exhibiting better and more robust  
 87 performance than that learning them individually (Shi et al., 2015). In order  
 88 to reduce the cost of memory storage and computational time, we further  
 89 improve the GH model to enable scalable graph-based hashing (SGH), which  
 90 selects a subset of cells from training cells to build an asymmetric graph to  
 91 preserve the similarity of neighbors. The complexity of building a graph in  
 92 SGH is significantly lower than that in GH. Next, we propose a novel method,  
 93 namely group-to-group matching, to conduct image retrieval by using the  
 94 cellular information in each image. In summary, our contributions are listed  
 95 as follows:

- 96 • We propose a novel image retrieval framework for large-scale histopatho-  
 97 logical image analysis based on cell encoding.

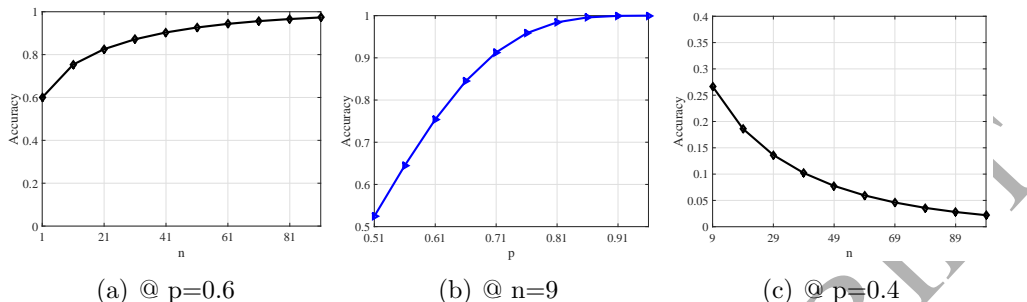


Fig 2: The change of image Classification accuracy with different  $n$  and  $p$ . (a) Accuracy vs  $n$  when  $p > 0.5$ , (b) Accuracy vs  $p$ , (a) Accuracy vs  $n$  when  $p < 0.5$

- 98      • We propose a supervised graph-based model via asymmetric relaxation  
 99      to learn binary codes and meanwhile reduce the accumulated quantiza-  
 100      tion errors between the discrete and continuous matrices. Furthermore,  
 101      we improve the model to allow scalable image retrieval by reducing both  
 102      space and time complexity.
- 103      • We propose a group-to-group matching method to retrieve images based  
 104      on the binary codes of cells.

105      The rest of this paper is structured as follows: Section 2 explains the rea-  
 106      son why cell examination is helpful for histopathological image classification.  
 107      Section 3 introduces the proposed graph-based model and its scalable ver-  
 108      sion. Section 4 shows the novel strategy, group-to-group matching. Section 5  
 109      presents and interprets experimental results on lung cancer images. Finally,  
 110      section 6 concludes this paper and points out the future work.

## 111 2. Why cell examination is helpful for histopathological image clas- 112 sification

113      In this section, we explain the reason why cell examination is helpful for  
 114      histopathological image classification from a statistical viewpoint. Given two  
 115      types of cells: adenocarcinoma (class 1) and squamous carcinoma (class 2)  
 116      in lung cancer images, suppose that an adenocarcinoma image  $I$  contains  $n$   
 117      (without loss of generality,  $n$  is assumed to be an odd constant) cells, with  $p$   
 118      representing the probability of each cell belonging to adenocarcinoma type,  
 119      where  $p \in (0.5, 1)$ . Based on the popular major voting strategy (Penrose,

120 1946), (Gans and Smart, 1996), the probability of the image  $I$  belongs to the  
 121 adenocarcinoma type is:

$$P_I = E(l_I = 1|p, n) = \sum_{i=\frac{n+1}{2}}^n C_n^i p^i (1-p)^{n-i}, \quad (1)$$

122 where  $l_I$  represents the predicted label of the image  $I$ .

123 With  $n$  fixed, the probability  $P_I$  increases when  $p$  grows; with  $p$  fixed, a  
 124 larger  $n$  means a higher  $P_I$ . For clarity, we present these two cases in Fig 2a-  
 125 b, which suggest that in a query image, when  $p \in (0.5, 1)$ , a larger number of  
 126 cells means a better image classification accuracy (Fig 2a), and a better cell  
 127 classification accuracy will also lead to a better image classification accuracy  
 128 (Fig 2b). These two observations show the relations between cell examination  
 129 and pathological image classification, which is one major motivation of our  
 130 framework. Note that when  $p \in (0, 0.5)$ , a larger number of cells will lead to  
 131 a worse image classification accuracy (Fig 2c).

### 132 3. Supervised Graph Hashing

#### 133 3.1. Graph-based Hashing

##### 134 3.1.1. Problem formulation

135 Each cell is cropped as an image patch (see Fig 1), from which feature  
 136 representation is calculated. Given a set of  $N$  cells  $\mathbf{X} = [\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N] \in$   
 137  $\mathbb{R}^{d \times N}$ , the goal of hashing is to find  $c$  hash functions  $H = [h_1, h_2, \dots, h_c]$   
 138 leading to  $c$ -bit Hamming embedding of  $\mathbf{X}$ . The hash function can be written  
 139 as:

$$h_k(\mathbf{x}_i) = \text{sgn}(\mathbf{a}_k^T \mathbf{x}_i + b_k), \quad (2)$$

140 where  $h_k$  is the  $k$ -th hash function with learnable parameters  $\mathbf{a}_k \in \mathbb{R}^d$  and  $b_k$ .  
 141 Usually  $b_k$  is equivalent to  $-\frac{1}{N} \sum_{i=1}^n \mathbf{a}_k^T \mathbf{x}_i$  and will be zero if  $\mathbf{X}$  is normalized  
 142 to have zero-mean.

143 Intuitively, adjacent cells in the original feature space should have similar  
 144 binary codes after embedding, meaning that they should still be close to each  
 145 other in the embedding space. Graph embedding is an attractive technique by  
 146 considering the intrinsic dimensionality. Traditional graph embedding (Yan  
 147 et al., 2007), (Nie et al., 2011) aims to find an optimal low-dimensional pro-  
 148 jection matrix  $\mathbf{A} = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_c] \in \mathbb{R}^{d \times c}$  to preserve the similarity among

149 data points, and the optimization problem is formulated as:

$$\begin{aligned} \min_{\mathbf{A}} \quad & Tr \{ \mathbf{A}^T \mathbf{X} \mathbf{L} \mathbf{X}^T \mathbf{A} \}, \\ \text{s.t.} \quad & \mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} = \mathbf{I}_c, \end{aligned} \quad (3)$$

150 where  $\mathbf{L} = \mathbf{D} - \mathbf{W}$ ,  $\mathbf{D}$  is a diagonal matrix with the  $i$ -th diagonal element  $d_{ii} =$   
 151  $\sum_{j=1}^n w_{ij}$ , and  $\mathbf{W}$  is a symmetric matrix with the weight  $w_{ij}$  between data  
 152 points  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . In order to produce binary codes for hashing encoding, the  
 153 linear projection function  $\mathbf{A}^T \mathbf{X}$  is usually replaced by the hashing function  
 154  $sgn(\mathbf{A}^T \mathbf{X})$ , and thus Eq. (3) becomes:

$$\begin{aligned} \min_{\mathbf{A}} \quad & Tr \{ sgn(\mathbf{A}^T \mathbf{X}) \mathbf{L} sgn(\mathbf{X}^T \mathbf{A}) \}, \\ \text{s.t.} \quad & sgn(\mathbf{A}^T \mathbf{X}) sgn(\mathbf{X}^T \mathbf{A}) = \mathbf{N} \mathbf{I}_c. \end{aligned} \quad (4)$$

155 By using the fact that  $Tr \{ sgn(\mathbf{A}^T \mathbf{X}) \mathbf{D} sgn(\mathbf{X}^T \mathbf{A}) \}$  is a constant and  $\mathbf{L} =$   
 156  $\mathbf{D} - \mathbf{W}$  as well as symmetric relaxation of the hashing function  $sgn(\mathbf{A}^T \mathbf{X})$   
 157 (Weiss et al., 2009), (Jiang et al., 2015), Eq. (4) can be written as:

$$\begin{aligned} \max_{\mathbf{A}} \quad & Tr \{ \mathbf{A}^T \mathbf{X} \mathbf{W} \mathbf{X}^T \mathbf{A} \}, \\ \text{s.t.} \quad & \mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} = \mathbf{N} \mathbf{I}_c, \end{aligned} \quad (5)$$

158 where the constraint can reduce the redundancy among data points (Shi  
 159 et al., 2016a).

160 One major problem of Eq. (5) is that the accumulated quantization error  
 161 between the hashing function (binary code matrix) and the linear projection  
 162 function might largely affect the retrieval accuracy, especially for a large  
 163 number of training data (Shen et al., 2015). In order to reduce accumulated  
 164 quantization errors, we utilize an asymmetric relaxation strategy (Shi et al.,  
 165 2016b) and propose the following optimization model:

$$\begin{aligned} \max_{\mathbf{A}, \mathbf{B} \in \{-1, 1\}^{c \times N}} \quad & Tr \{ \mathbf{B} \mathbf{W} \mathbf{X}^T \mathbf{A} \} - \frac{\alpha}{2} \|\mathbf{B} - \mathbf{A}^T \mathbf{X}\|_F^2, \\ \text{s.t.} \quad & \|\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} - \mathbf{N} \mathbf{I}_c\|_F^2 \leq \epsilon, \end{aligned} \quad (6)$$

166 where  $\mathbf{B}$  is a discrete matrix representing the binary codes of training data,  
 167 and the Frobenius norm regularization term aims to reduce the accumulated  
 168 quantization error between the binary code matrix  $\mathbf{B}$  and the linear projec-  
 169 tion function  $\mathbf{A}^T \mathbf{X}$ . Note that we utilize the constraint  $\|\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} - \mathbf{N} \mathbf{I}_c\|_F^2 \leq$   
 170  $\epsilon$  rather than  $\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} = \mathbf{N} \mathbf{I}_c$  to improve the robustness of the projection

171 matrix  $\mathbf{A}$ . In addition, jointly learning the binary codes  $\mathbf{B}$  and projection  
 172 matrix  $\mathbf{A}$  in Eq. (6) can further improve the robustness of the projection  
 173 matrix  $\mathbf{A}$  (Shi et al., 2015). Eq. (6) seems to be similar to the objective  
 174 function in (Liu et al., 2014), however, they are intrinsically different in the  
 175 following two aspects: (1) Eq. (6) is to learn the projection matrix  $\mathbf{A}$  for  
 176 supervised hashing, while (Liu et al., 2014) focuses on unsupervised hashing;  
 177 (2) The discrete constraint in the objective function of (Liu et al., 2014) is  
 178 symmetric, but an asymmetric discrete constraint is used in Eq. (6), which  
 179 can be solved with lower time costs.

### 180 3.1.2. Optimization procedure

181 It is difficult to simultaneously calculate  $\mathbf{B}$  and  $\mathbf{A}$  in Eq. (6). By intro-  
 182 ducing an auxiliary  $\mathbf{C}$ , the optimization problem in Eq. (6) can be optimized  
 183 by solving the following three subproblems:

184 **B-subproblem:**

$$185 \max_{\mathbf{B}^{c \times N} \in \{-1,1\}} \text{Tr} \{ \mathbf{B}(\mathbf{W} + \alpha \mathbf{I}_N) \mathbf{X}^T \mathbf{A} \}. \quad (7)$$

186 Eq. (7) can be obtained from Eq. (6) because  $\text{Tr} \{ \mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} \}$  and  $\text{Tr} \{ \mathbf{B} \mathbf{B}^T \}$   
 187 are constants. The solution to Eq. (7) is  $\mathbf{B} = \text{sign}(\mathbf{A}^T \mathbf{X}(\mathbf{W} + \alpha \mathbf{I}_N))$ .

188 **C-subproblem:**

$$189 \begin{aligned} & \max_{\mathbf{C} \in \{-1,1\}^{N \times c}} \text{Tr} \{ \mathbf{B}(\mathbf{W} + \alpha \mathbf{I}_N) \mathbf{C} \}, \\ & \text{s.t. } \mathbf{C}^T \mathbf{C} = \mathbf{N} \mathbf{I}_c, \mathbf{1}_N \mathbf{C} = 0, \end{aligned} \quad (8)$$

190 where  $\mathbf{I}_c \in \mathbb{R}^{c \times c}$  is a unit matrix,  $\mathbf{1}_N \in \mathbb{R}^N$  is a row vector with all elements  
 191 being one, and  $\mathbf{1}_N \mathbf{C} = 0$  because of  $\mathbf{X} \mathbf{1}_N^T = 0$ . We can obtain  $\mathbf{C}$  based on  
 192 the singular value decomposition (SVD) of  $\mathbf{B}(\mathbf{W} + \alpha \mathbf{I}_N) \mathbf{J} = \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$ , where  
 193  $\mathbf{J} = \mathbf{I}_N - \frac{1}{N} \mathbf{1}_N^T \mathbf{1}_N$ , and  $\mathbf{C} = \sqrt{N} \mathbf{V}(:, 1:c) \mathbf{U}^T$ .

194 **A-subproblem:**

$$195 \min_{\mathbf{A}} \|\mathbf{C} - \mathbf{X}^T \mathbf{A}\|_F^2 + \gamma \|\mathbf{A}\|_F^2, \quad (9)$$

196 whose solution is  $\mathbf{A} = (\mathbf{X} \mathbf{X}^T + \gamma \mathbf{I}_d)^{-1} \mathbf{X} \mathbf{C}$ . When  $\gamma = 0$ ,  $\mathbf{C} = \mathbf{X}^T \mathbf{A}$  and  
 197 then  $\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} = \mathbf{N} \mathbf{I}_c$ ; when  $\gamma > 0$ ,  $\|\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} - \mathbf{N} \mathbf{I}_c\|_F^2 \leq \epsilon$ , and a larger  
 198  $\gamma$  means a larger  $\epsilon$ .

199 The detailed procedure to solve Eq. (6), namely graph-based hashing  
 200 (GH), is summarized in Algorithm 1. Since each step in Algorithm 1 has a

**Algorithm 1:** Graph-based Hashing (GH)

**Input:** Data matrix  $\mathbf{X} \in \mathbb{R}^{d \times N}$ , weight matrix  $\mathbf{W} \in \mathbb{R}^{N \times N}$ , regularization parameter  $\alpha, \gamma$ , the number of bits  $c$ , the maximum iterations  $k$ .

**Output:** Binary code matrix  $\mathbf{B} \in \mathbb{R}^{c \times N}$ , projection matrix  $\mathbf{A} \in \mathbb{R}^{d \times c}$ .

1. Calculate  $\mathbf{A}$  corresponding to the  $c$  largest eigenvalues of  $\mathbf{X}(\mathbf{W} + \alpha \mathbf{I}_N)\mathbf{X}^T$ .

2. Calculate the matrix  $\mathbf{B} = \text{sgn}(\mathbf{A}^T \mathbf{X})$ .

3. Loop until convergence or reach maximum iterations

3.1 Update  $\mathbf{B} = \text{sign}(\mathbf{A}^T \mathbf{X}(\mathbf{W} + \alpha \mathbf{I}_N))$ ;

3.2 Update  $\mathbf{C} = \sqrt{N} \mathbf{V}(:, 1:c) \mathbf{U}^T$ ,

where  $\mathbf{U} \Sigma \mathbf{V}^T = \mathbf{B}(\mathbf{W} + \alpha \mathbf{I}_N) \mathbf{J}$ ;

3.3 Update  $\mathbf{A} = (\mathbf{X} \mathbf{X}^T + \gamma \mathbf{I}_d)^{-1} \mathbf{X} \mathbf{C}$ .

201 closed form solution, empirically we can obtain approximately optimal solu-  
202 tions within a few iterations (e.g.  $f=5$ ), where  $f$  is the number of iterations.

### 203 3.1.3. Time complexity analysis of GH

204 In Algorithm 1, the calculation of the matrix  $\mathbf{A}$  in step 1 requires  $\mathcal{O}(N^2 d)$   
205 operations because of  $N > d$ , and step 2 needs  $\mathcal{O}(Ndc)$  operations to calcu-  
206 late the matrix  $\mathbf{B}$ . In step 3, the calculation of the matrices  $\mathbf{B}$ ,  $\mathbf{C}$  and  $\mathbf{A}$   
207 requires  $\mathcal{O}(fN^2 d)$ ,  $\mathcal{O}(fN^2 c)$  and  $\mathcal{O}(fNd^2)$  operations, respectively. Usually  
208  $c < d < N$ , and thus the total computational complexity of Algorithm 1 is  
209  $\mathcal{O}(fN^2 d)$ .

### 210 3.2. Scalable Graph-based Hashing

211 In GH, the complexity of both storage and computational time to con-  
212 struct a graph with  $N$  data points is  $\mathcal{O}(N^2)$ , and thus GH might not adapt to  
213 large-scale datasets (a large  $N$ ). To improve memory requirement and com-  
214 putational efficiency, we present a fast hashing algorithm, namely scalable  
215 graph-based hashing (SGH), to approximate the model in Eq. (6).

#### 216 3.2.1. Problem formulation

217 Dimensionality reduction methods, like principal component analysis (PCA)  
218 (Zou et al., 2006) and nonnegative matrix factorization (NMF) (Recht et al.,  
219 2012), suggest that a matrix with high-dimensionality yet low rank can be  
220 represented or approximated by a linear combination of basis vectors. In our

221 problem, the weight matrix  $\mathbf{W}$  can be approximated by  $\mathbf{W} \approx \mathbf{TP}$ , where  
 222  $\mathbf{T} \in \mathbb{R}^{N \times M}$ ,  $\mathbf{T} \subset \mathbf{W}$  and  $\mathbf{P} \in \mathbb{R}^{M \times N}$  is a weight matrix. Therefore, Eq. (6)  
 223 can be reformulated as:

$$\begin{aligned} & \max_{\mathbf{A}, \mathbf{B} \in \{-1, 1\}^{c \times N}} Tr \{ \mathbf{BTPX}^T \mathbf{A} \} - \frac{\alpha}{2} \| \mathbf{B} - \mathbf{A}^T \mathbf{X} \|_F^2, \\ & s.t. \quad \| \mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} - N \mathbf{I}_c \|_F^2 \leq \epsilon. \end{aligned} \quad (10)$$

224 Although  $\mathbf{T}$  can be a subset of  $\mathbf{W}$ , computing  $\mathbf{P}$  is expensive. To avoid the  
 225 calculation of  $\mathbf{P}$ , we suppose  $\mathbf{Z} = \mathbf{XP}^T$  be generated new data (anchors)  
 226 that are used to construct graph, and  $\mathbf{T}$  can be seen as a weight matrix  
 227 characterizing the relationship between the training data  $\mathbf{X}$  and anchors  $\mathbf{Z}$ .  
 228 Since the projection matrix mainly depends on the data matrix  $\mathbf{X}$  due to  
 229  $M \ll N$ , we ignore the effect of the new anchors on the projection matrix  
 230 and thus neglect the calculation of  $\mathbf{P}$ . Let  $\mathbf{D} = \mathbf{Z}^T \mathbf{A}$ , Eq. (10) can be  
 231 rewritten as:

$$\begin{aligned} & \max_{\mathbf{A}, \mathbf{D}, \mathbf{B} \in \{-1, 1\}^{c \times N}} Tr \{ \mathbf{BTD} \} - \frac{\alpha}{2} \| \mathbf{B} - \mathbf{A}^T \mathbf{X} \|_F^2, \\ & s.t. \quad \mathbf{D}^T \mathbf{D} = M \mathbf{I}_c, \| \mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} - N \mathbf{I}_c \|_F^2 \leq \epsilon, \end{aligned} \quad (11)$$

232 where the constraint  $\mathbf{D}^T \mathbf{D} = M \mathbf{I}_c$  is derived from the constraint  $\mathbf{A}^T \mathbf{X} \mathbf{X}^T \mathbf{A} =$   
 233  $N \mathbf{I}_c$ .

### 234 3.2.2. Optimization procedure

235 Similar to solving Eq. (6), we adopt an iterative and alternative strategy  
 236 to compute the discrete binary codes  $\mathbf{B}$  and the projection matrix  $\mathbf{A}$  in Eq.  
 237 (11). We divide the optimization problem in Eq. (11) into four subproblems  
 238 as follows and present the details of scalable graph-based hashing (SGH) in  
 239 Algorithm 2.

240 **B-subproblem:**

$$241 \max_{\mathbf{B} \in \{-1, 1\}^{c \times N}} Tr \{ \mathbf{BTD} + \alpha \mathbf{B} \mathbf{X}^T \mathbf{A} \}, \quad (12)$$

242 whose solution is  $\mathbf{B} = \text{sign}(\mathbf{D}^T \mathbf{T}^T + \alpha \mathbf{A}^T \mathbf{X})$ .

243 **D-subproblem:**

$$244 \max_{\mathbf{D}} Tr \{ \mathbf{BTD} \}, s.t. \quad \mathbf{D}^T \mathbf{D} = M \mathbf{I}_c. \quad (13)$$

245 The solution to Eq. (13) is  $\mathbf{D} = \sqrt{M} \mathbf{V}_D(:, 1:c) \mathbf{U}_D^T$ , where  $\mathbf{U}_D$  and  $\mathbf{V}_D$  can  
 246 be obtained by the SVD of  $\mathbf{BT} = \mathbf{U}_D \mathbf{\Sigma}_D \mathbf{V}_D^T$ .

**Algorithm 2:** Scalable Graph-based Hashing (SGH)

**Input:** Data matrix  $\mathbf{X} \in \mathbb{R}^{d \times N}$ , graph matrix  $\mathbf{T} \in \mathbb{R}^{N \times M}$ , parameters  $\alpha, \gamma$ , the number of bits  $c$ , the number of iterations  $k$ .

**Output:** Binary code matrix  $\mathbf{B} \in \mathbb{R}^{c \times N}$ , projection matrix  $\mathbf{A} \in \mathbb{R}^{d \times c}$ .

1. Calculate  $\mathbf{A}$  corresponding to the  $c$  largest eigenvalues of  $\mathbf{X}\mathbf{T}\mathbf{T}^T\mathbf{X}^T$ ;
2. Calculate the matrix  $\mathbf{B} = \text{sgn}(\mathbf{A}^T\mathbf{X})$ ;
3. Loop until converge or reach maximum iterations
  - 3.1 Update  $\mathbf{B} = \text{sgn}(\mathbf{D}^T\mathbf{T}^T + \alpha\mathbf{A}^T\mathbf{X})$ ;
  - 3.2 Update  $\mathbf{D} = \sqrt{M}\mathbf{V}_D(:, 1:c)\mathbf{U}_D^T$ , where  $\mathbf{B}\mathbf{T} = \mathbf{U}_D\mathbf{\Sigma}_D\mathbf{V}_D^T$ ;
  - 3.3 Update  $\mathbf{C} = \sqrt{N}\mathbf{V}_C(:, 1:c)\mathbf{U}_C^T$ , where  $\mathbf{B}\mathbf{J} = \mathbf{U}_C\mathbf{\Sigma}_C\mathbf{V}_C^T$ ;
  - 3.4 Update  $\mathbf{A} = (\mathbf{X}\mathbf{X}^T + \gamma\mathbf{I}_d)^{-1}\mathbf{X}\mathbf{C}$ ;

247 **C-subproblem:**

248

$$\max_{\mathbf{C}} \text{Tr} \{ \mathbf{B}\mathbf{C} \}, \text{ s.t. } \mathbf{C}^T\mathbf{C} = N\mathbf{I}_c, \mathbf{1}_N\mathbf{C} = 0. \quad (14)$$

249 The solution to Eq. (14) is  $\mathbf{C} = \sqrt{N}\mathbf{V}_C(:, 1:c)\mathbf{U}_C^T$ , where  $\mathbf{U}_C$  and  $\mathbf{V}_C$  can  
250 be obtained by the SVD of  $\mathbf{B}\mathbf{J} = \mathbf{U}_C\mathbf{\Sigma}_C\mathbf{V}_C^T$ .

251 **A-subproblem:** This subproblem is the same to Eq. (9).

252 Since supervised hashing methods usually have better performance than  
253 unsupervised hashing, in this paper, we focus on supervised hashing method  
254 to retrieve cells, and the weight matrix  $\mathbf{W} \in \mathbb{R}^{N \times N}$  that utilizes semantic  
255 information is defined as follows:

$$w_{ij} = \begin{cases} \frac{1}{N_k} & \text{if } x_i, x_j \in k\text{-th class,} \\ 0 & \text{otherwise,} \end{cases} \quad (15)$$

256 where  $N_k$  is the number of samples in the  $k$ -th class,  $1 \leq k \leq K$  and  $K$  is  
257 the number of classes. Since  $\mathbf{T} \in \mathbb{R}^{N \times M}$  is the subset of  $\mathbf{W}$ , we can attain  $\mathbf{T}$   
258 by randomly selecting  $M_k$  data points as anchors to construct  $\mathbf{T}$ , where  $M_k$   
259 is the selected number of anchors from the  $k$ -th class.

### 260 3.2.3. Time complexity analysis

261 In Algorithm 2, the computational complexity of calculating  $\mathbf{A}$  in step 1  
262 is  $\max(\mathcal{O}(NMd), \mathcal{O}(d^3))$ . The calculation of the matrix  $\mathbf{B}$  in step 2 requires  
263  $\mathcal{O}(Ndc)$  operations. In step 3, calculating the matrices  $\mathbf{B}$ ,  $\mathbf{D}$ ,  $\mathbf{C}$  and  $\mathbf{A}$  needs

264  $\max(\mathcal{O}(fNMd), \mathcal{O}(fNdc)), \mathcal{O}(fNMd), \mathcal{O}(fNc^2)$  and  $\mathcal{O}(fNd^2)$  operations,  
 265 respectively. Since usually  $fc < d$ , the total time complexity of Algorithm 2  
 266 is  $\max(\mathcal{O}(fNd^2), \mathcal{O}(NMd))$ .

### 267 3.3. Relations to other related hashing algorithms

268 There are many hashing algorithms proposed in the literature, and to  
 269 better explain the proposed methods, in this section we discuss the relations  
 270 between ours and several popular and related hashing algorithms.

271 Spectral hashing (SH) (Weiss et al., 2009) is a popular graph based hash-  
 272 ing algorithm, but it has two major limitations: (1) The complexity of both  
 273 storage and computational time to construct a graph with  $N$  data points  
 274 is  $\mathcal{O}(N^2)$ ; (2) Symmetric relaxation is used to solve the NP-hard optimiza-  
 275 tion problem, which would generate accumulated errors between the discrete  
 276 binary code matrix and its relaxed continuous matrix. To reduce the com-  
 277 plexity of building a graph, anchor graph hashing (AGH) (Liu et al., 2011)  
 278 constructs an approximative symmetric graph with a small number of an-  
 279 chors to preserve the similarity between neighbors. However, AGH utilizes  
 280 the symmetric relaxation strategy to solve the NP-hard optimization prob-  
 281 lem. To reduce the accumulated quantitative errors, discrete graph hashing  
 282 (DGH) (Liu et al., 2014) preserves the symmetric discrete constraint. Com-  
 283 pared to GH that preserves the asymmetric discrete constraint, the optimiza-  
 284 tion procedure of DGH is more complex and takes more training time. In  
 285 addition, both AGH and DGH focus on learning binary codes in an unsu-  
 286 pervised manner such that it is difficult to directly utilize semantic (label)  
 287 information to learn binary codes.

288 KSH (Liu et al., 2012) is a kernel based hashing algorithm applied to  
 289 pathological image analysis. It utilizes symmetric relaxation followed by  
 290 greedy optimization to solve its non-differentiable objective function. This  
 291 strategy is not able to greatly reduce the accumulated errors and usually costs  
 292 a large amount of training time. Joint kernel graph hashing (JKGH) (Jiang  
 293 et al., 2015) also uses symmetric relaxation to solve its NP-hard optimization  
 294 problem. It constructs the kernel of each data with weighting each sub-  
 295 kernel constructed by each feature and thus achieves better performance  
 296 than KSH. However, JKGH usually requires considerable training and test  
 297 time to construct the kernel. Hence it is not suitable for tackling large-scale  
 298 data with high-dimensional features.

299 Although kernel-based supervised discrete hashing (KSDH) (Shi et al.,  
 300 2016b) also utilizes the asymmetric relaxation strategy to directly learn bi-

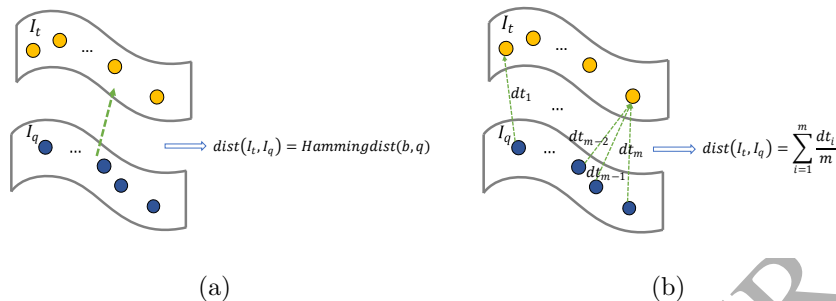


Fig 3: The ideas of point-to-point and group-to-group. (a) point-to-point, (b) group-to-group. (Each point in a plane represents a cell in one image, the discrete vectors  $\mathbf{b}$  and  $\mathbf{q}$  represent binary codes of the entire target and query images, respectively.)

301 nary codes, it does not learn binary codes and a projection matrix simulta-  
 302 neously, which might decrease its performance (Shi et al., 2015).

303 Compared to KSH, JKSH and kernel-based supervised discrete hashing  
 304 (KSDH) (Shi et al., 2016b), the proposed algorithms have three major advan-  
 305 tages: (1) GH and SGH are linear methods and thus they do not require  
 306 kernel selection and construction; (2) GH and SGH utilize asymmetric relax-  
 307 ation to solve the NP-hard optimization problem, and thus they can reduce  
 308 the accumulated quantitative errors; (3) KSH, JKSH and KSDH require  
 309  $\mathcal{O}(N^2)$  operations to construct their weight matrix, while SGH needs only  
 310  $\mathcal{O}(MN)$  ( $M \ll N$ ) operations.

#### 311 4. Group-to-group matching for pathological image retrieval

312 One single histopathology image usually contains multiple cells, and this  
 313 would fail the traditional point-to-point matching (Fig 3a) in image re-  
 314 trieval, which calculates the Hamming distance between two entire images  
 315 that are encoded into a set of binary codes. To address this problem, we  
 316 propose a novel strategy: group-to-group (GTG) matching, for image re-  
 317 trieval using cell-level information. Given one target image  $I_t$  containing  
 318  $n$  cells and one query image  $I_q$  including  $m$  cells, denote their binary vec-  
 319 tors by  $[\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n]$  and  $[\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_m]$ , respectively, where  $\mathbf{b}_i \in \{0, 1\}^c$ ,  
 320  $\mathbf{q}_j \in \{0, 1\}^c$ ,  $1 \leq i \leq n$  and  $1 \leq j \leq m$ . We can regard that the cells in one  
 321 target image  $\mathbf{I}_t$  form a plane  $\mathcal{P}_t$ , and each query cell can be regarded as a  
 322 point outside of the plane  $\mathcal{P}_t$ . Then the distance between the point  $\mathbf{q}_i$  and

**Algorithm 3:** Group-to-group (GTG)

**Input:** Training binary matrix  $[\mathbf{b}_1, \mathbf{b}_2, \dots, \mathbf{b}_n] \in \mathbb{R}^{n \times c}$ ,  
 query binary matrix  $[\mathbf{q}_1, \mathbf{q}_2, \dots, \mathbf{q}_m] \in \mathbb{R}^{m \times c}$ .

**Output:**  $dist(I_t, I_q)$ .

---

```

for  $j = 1$  to  $m$  do
  for  $i = 1$  to  $n$  do
     $dt(i) = Hammingdist(\mathbf{b}_i, \mathbf{q}_j)$ ,
  end;
 $dist(j, I_t) \leftarrow min(dt)$ ,
end;
 $dist(I_t, I_q) \leftarrow \frac{1}{m} \sum_{j=1}^m dist(j, I_t)$ .

```

---

323 the plane  $\mathcal{P}_t$  is equivalent to  $dist(\mathbf{q}_j, \mathcal{P}_t) = \min_{1 \leq i \leq n} dist(\mathbf{q}_j, \mathbf{b}_i)$ , and the distance  
 324 between these two images is defined as follows:

$$dist(I_t, I_q) = \frac{1}{m} \sum_{j=1}^m \min_{1 \leq i \leq n} dist(\mathbf{q}_j, \mathbf{b}_i). \quad (16)$$

325 For better illustration, we explain this idea in Fig 3b. The details of the  
 326 GTG matching to calculate the distance between two images are shown in  
 327 Algorithm 3. Note that compared to the major voting strategy, GTG can be  
 328 viewed as a weighted voting strategy, because the cells in one query image  
 329 has different distances to the target image.

## 330 5. Experiments

331 To evaluate the proposed hashing algorithms (GH and SGH), we con-  
 332 duct extensive experiments on the lung cancer image dataset with two types  
 333 of diseases: squamous cell carcinoma and adenocarcinoma. The lung can-  
 334 cer dataset is collected from The Cancer Genome Atlas (TCGA) (Institute:,  
 335 2013), which consists of 1240 images (630 squamous and 610 adenocarcinoma  
 336 images). These images contains around 1248K squamous and 589K adeno-  
 337 carcinoma cells. All cells are detected and segmented using the method in  
 338 (Xing et al., 2014). Then all cells are cropped as image patches with corre-  
 339 sponding labels. For each cell, we extract a 1024-dimensional feature vector  
 340 by GIST (Oliva and Torralba, 2001) and HOG (Dalal and Triggs, 2005), re-  
 341 spectively. Additionally, we also utilize LeNet (LeCun et al., 2015) to extract  
 342 a 300-dimensional feature vector from each cell. The sample images and their  
 343 segmented cells are shown in Fig 4.

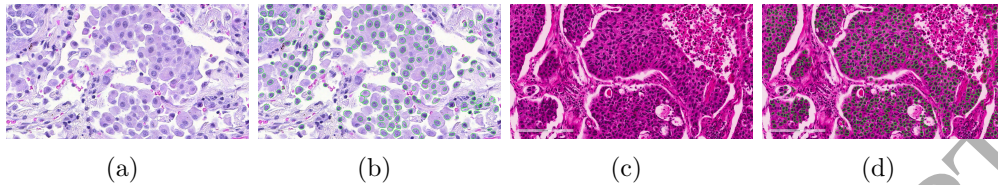


Fig 4: Sample images and segmented cells. (a) and (b) Adenocarcinoma image and cell segmentation results, (c) and (d) squamous cell carcinoma image and cell segmentation results.

### 344 5.1. Experimental setting:

345 We compare GH and SGH against two state-of-the-art supervised hash-  
 346 ing algorithms: KSH (Liu et al., 2012) and SDH (Shen et al., 2015). After  
 347 using hashing algorithms to encode each cell into binary codes, the GTG  
 348 matching is used for image retrieval. To better show the characteristics of  
 349 GTG, we also present the classification accuracy of the hashing algorithms  
 350 with major voting (MV). In addition, we compare with two popular classi-  
 351 fiers: support vector machine (SVM) (Fan et al., 2008) and nearest neighbors  
 352 (NN). We utilize SVM and NN to classify cells and then applying major vot-  
 353 ing to image classification, and also show the baseline results obtained by  
 354 directly applying SVM and NN on the holistic high-dimensional features ex-  
 355 tracted from the whole image. Specifically, similar to (Jiang et al., 2015),  
 356 (Zhang et al., 2015a), we first detect scale-invariant keypoints from the whole  
 357 pathology image, and then employ SIFT (Lowe, 2004) and HOG to extract  
 358 features around these keypoints. Next, both descriptors are encoded as 2000-  
 359 dimensional histograms using bag-of-words (BoW) methods (Caicedo et al.,  
 360 2009), (Caicedo et al., 2011). Finally, SVM and NN are applied to these  
 361 histograms for image classification.

362 We adopt five criteria including classification accuracy, precision, mean  
 363 average precision (MAP), training and test time, to evaluate the hashing  
 364 algorithms. Classification accuracy evaluates the classification performance;  
 365 precision measures the retrieval accuracy (returned corrected cells over total  
 366 returned samples), MAP evaluates the ranking performance of total returned  
 367 cells, training and test time are used to measure the time cost on training  
 368 and query, respectively. For GH and SGH, we set  $\alpha = 10^{-4}$  and  $\alpha = 10^{-3}$ ,  
 369 respectively, and  $\gamma = 10^{-3}$  in the following experiments, and we set  $M =$   
 370  $0.1N$  in SGH. For KSH, SDH, KSDH\_H and KSDH\_B, we randomly select

371 500 training cells to construct the kernels.

372 We first randomly select 100 images from each category for training and  
 373 the remaining images are used for testing. We randomly select 100 cells from  
 374 each training image to constitute a training cell dataset. In total, there are  
 375 20K cells in the training dataset. To form a test cell dataset we randomly  
 376 select  $m$  cells from each test image. After obtaining the binary codes of these  
 377 test cells, we use them to classify test images and retrieve relevant images.  
 378 Since  $m$  is an important hyperparameter, we present the performance of all  
 379 hashing algorithms with respect to different  $m$ 's. We repeat this process 10  
 380 times and calculate the mean accuracy of different methods. All experiments  
 381 are conducted using MATLAB on a 3.50GHz Intel Xeon CPU with 128GB  
 382 memory.

### 383 5.2. Experimental results

384 Table 1 shows image classification results of two popular classifiers, SVM  
 385 and NN, using the holistic high-dimensional HOG and SIFT features ex-  
 386 tracted from the whole image, and cell classification results of eight algo-  
 387 rithms using HOG and GIST features extracted from cell images. We use  
 388 GIST rather than SIFT to extract cell features because SIFT is much slower  
 389 to extract features from a larger number of cell images. Additionally, we also  
 390 present cell classification results of eight algorithms using CNN features. As  
 391 shown in Table 1, GH achieves the best cell classification accuracy with HOG,  
 392 GIST and CNN features. Although SGH has slightly worse cell classification  
 393 accuracy than SVM, SDH, KSDH.H and KSDH.B with HOG features, it  
 394 significantly outperforms NN and KSH. When GIST and CNN features are  
 395 used, SGH attains better performance than NN, KSH, SDH, KSDH.H and  
 396 KSDH.B. The results in Table 1 are used as the baseline for comparison.

397 Table 2 presents the classification results of two classifiers SVM, NN  
 398 and six retrieval algorithms KSH, SDH, KSDH.H, KSDH.B, GH and SGH  
 399 with MV, and the classification and retrieval results of the six retrieval algo-  
 400 rithms with GTG. Table 2 suggests that the retrieval algorithms with GTG  
 401 achieve better classification accuracy than those with MV in most of cases.  
 402 Meanwhile, when 9 cells are selected from each query image, KSH+MV and  
 403 SGH+MV obtain worse image classification performance, since the classi-  
 404 fication accuracy of KSH (using HOG, GIST or CNN features) and SGH  
 405 (using HOG features) on squamous cells is very low ( $< 50\%$ ). On the con-  
 406 trary, KSH+GTG and SGH+GTG achieve much better image classification

Table 1: Baseline classification accuracy (%) on images and cells. (Hashing algorithms encode each cell into 10-bit binary codes.)

(a)

	Method	Adeno	Squam	Mean	Adeno	Squam	Mean
		HOG			SIFT		
Image	SVM	40.39	81.51	61.35	74.90	81.51	78.27
	NN	61.57	59.25	60.38	74.71	83.02	78.94

(b)

Cell ( $m = 9$ )	Method	HOG			GIST			CNN		
		Adeno	Squam	Mean	Adeno	Squam	Mean	Adeno	Squam	Mean
Cell ( $m = 9$ )	SVM	64.60	60.88	62.70	64.99	64.63	64.81	78.54	71.20	74.60
	NN	42.57	63.88	53.43	56.21	56.67	56.44	70.07	75.05	72.61
	KSH	83.66	24.19	53.35	91.42	12.54	51.22	99.11	0.01	49.08
	SDH	59.98	68.76	62.46	69.04	59.16	64.01	75.75	72.50	74.10
	KSDH.H	71.74	53.98	62.69	73.90	55.24	64.39	80.15	67.95	73.94
	KSDH.B	69.63	56.39	62.88	72.57	56.48	64.37	78.96	69.63	74.20
	GH	64.77	61.28	<b>62.99</b>	68.39	61.74	<b>65.00</b>	76.75	73.50	<b>75.10</b>
	SGH	75.99	49.27	62.37	74.14	55.22	64.50	79.98	68.33	74.27

407 performance. This indicates the superior performance of GTG to MV. In ad-  
 408 dition, with 9 cells selected from each query image, SVM, NN and six hashing  
 409 algorithms with CNN features can attain better performance than that with  
 410 HOG or GIST features; among these algorithms, SGH+GTG achieves the  
 411 best classification accuracy (88.83%), precision (88.50%) and MAP (88.73%),  
 412 GH+GTG with CNN features achieves similar performance to the other hash-  
 413 ing algorithms except SGH+GTG.

Table 2: Classification accuracy, precision and MAP (%) of different algorithms on lung cancer images. ('Accuracy' represents the classification accuracy in table and 'All' means that all cells in the query image are used for retrieval. Hashing algorithms encode each cell into 10-bit binary codes.)

	Method	HOG				
		Accuracy			Precision (Top 100)	MAP (Top 100)
		Adeno	Squam	Mean		
$m = 9$	SVM+MV	76.73	81.96	71.70	-	-
	NN+MV	30.59	75.09	53.27	-	-
	KSH+MV	98.43	5.28	50.96	-	-
	SDH+MV	71.75	82.96	77.46	-	-
	KSDH.H+MV	90.91	60.00	75.75	-	-
	KSDH.B+MV	89.27	62.36	76.07	-	-
	GH+MV	82.35	71.32	76.73	-	-
	SGH+MV	99.22	10.19	53.85	-	-

$m = 9$	KSH+GTG	79.98	80.70	<b>80.35</b>	66.34	72.14
	SDH+GTG	74.71	79.81	77.31	77.31	77.31
	KSDH.H+GTG	80.76	75.42	78.04	77.89	83.78
	KSDH.B+GTG	89.20	76.74	77.94	78.09	78.06
	<b>GH+GTG</b>	82.35	73.21	77.69	77.89	77.99
	<b>SGH+GTG</b>	85.69	73.21	79.33	<b>79.02</b>	<b>85.54</b>
All	SVM+MV	98.04	91.89	94.90	-	-
	KSH+GTG	93.12	99.00	<b>96.12</b>	79.47	87.83
	SDH+GTG	90.59	98.68	94.71	94.71	94.71
	KSDH.H+GTG	96.08	94.15	95.10	94.64	95.03
	KSDH.B+GTG	96.08	94.15	95.10	95.10	95.10
	<b>GH+GTG</b>	97.65	94.34	95.96	<b>96.02</b>	<b>95.91</b>
	<b>SGH+GTG</b>	98.24	93.02	95.58	95.23	94.93
<b>GIST</b>						
$m = 9$	SVM+MV	81.25	80.00	80.26	-	-
$m = 9$	NN+MV	71.18	61.70	66.35	-	-
	KSH+MV	100	0.19	49.13	-	-
	SDH+MV	87.84	73.21	80.38	-	-
	KSDH.H+MV	92.80	48.00	70.40	-	-
	KSDH.B+MV	91.20	56.80	74.00	-	-
	GH+MV	87.45	71.17	82.22	-	-
	SGH+MV	95.10	57.74	76.06	-	-
$m = 9$	KSH+GTG	90.39	72.45	81.25	65.84	72.04
	SDH+GTG	88.04	73.58	80.67	80.67	80.67
	KSDH.H+GTG	90.00	73.58	81.63	81.34	84.07
	KSDH.B+GTG	90.20	73.77	81.83	81.83	81.83
	<b>GH+GTG</b>	88.24	74.91	81.44	81.37	82.12
	<b>SGH+GTG</b>	91.76	76.98	<b>84.23</b>	<b>84.10</b>	<b>86.81</b>
All	SVM+MV	98.61	96.82	97.75	-	-
	KSH+GTG	96.65	99.00	97.85	81.04	89.82
	SDH+GTG	100	92.83	96.35	96.35	96.35
	KSDH.H+GTG	100	94.91	97.40	96.83	97.20
	KSDH.B+GTG	99.80	94.89	97.30	97.30	96.05
	<b>GH+GTG</b>	99.61	95.47	97.56	97.56	<b>97.56</b>
	<b>SGH+GTG</b>	99.80	96.23	<b>97.98</b>	<b>97.50</b>	96.52
<b>CNN</b>						
$m = 9$	SVM+MV	88.43	83.77	86.06	-	-
	NN+MV	88.24	86.23	87.21	-	-
	KSH+MV	100	0	49.04	-	-
	SDH+MV	91.57	82.08	86.73	-	-
	KSDH.H+MV	91.50	76.23	85.48	-	-
	KSDH.B+MV	92.94	81.32	87.02	-	-
	GH+MV	89.02	84.34	86.63	-	-
	SGH+MV	88.82	86.79	<b>87.79</b>	-	-

$m = 9$	KSH+GTG	92.55	63.58	77.79	68.31	73.14
	SDH+GTG	91.26	83.46	87.29	87.28	87.19
	KSDH.H+GTG	92.55	82.64	87.39	87.39	87.39
	KSDH.B+GTG	90.76	85.04	87.85	87.85	87.85
	<b>GH+GTG</b>	90.22	84.47	88.31	88.31	88.31
	<b>SGH+GTG</b>	91.41	86.34	<b>88.83</b>	<b>88.50</b>	<b>88.73</b>
All	SVM+MV	94.71	93.40	94.04	-	-
	KSH+GTG	97.25	93.40	95.29	73.54	82.15
	SDH+GTG	96.86	92.83	94.81	94.81	94.81
	KSDH.H+GTG	97.65	91.32	94.42	94.12	94.25
	KSDH.B+GTG	96.86	92.82	94.80	94.80	94.80
	<b>GH+GTG</b>	97.06	93.21	95.10	95.10	95.10
	<b>SGH+GTG</b>	97.45	93.40	<b>95.38</b>	<b>95.38</b>	<b>95.38</b>

414

415 When all cells are selected, both GH+GTG and SGH+GTG outperform  
416 SDH on classification accuracy, precision and MAP; GH+GTG has similar  
417 performance to KSDH.H and KSDH.B, while SGH+GTG has slightly bet-  
418 ter performance; GH+GTH, SGH+GTG and KSH+GTG achieve almost the  
419 same classification accuracy, however, GH+GTH and SGH+GTG have sig-  
420 nificantly better precision and MAP than KSH. Compared to the case with  
421 only 9 cells selected, SVM+MV, KSH+GTG, SDH+GTG, KSDH.H+GTG,  
422 KSDH.B+GTG, GH+GTG and SGH+GTG using all cells of each query  
423 image can achieve better performance, especially for SGH+GTG, which us-  
424 ing GIST features can achieve the best classification accuracy (97.98%) and  
425 precision (97.50%), and the sub-best MAP (96.52%). Note that we do not  
426 show the results of NN+MV when all cells are selected, because it is very  
427 computationally expensive. Compared to Table 1, Table 2 shows that image  
428 classification using features extracted from cells can obtain better accuracy  
429 than that using the holistic high-dimensional features extracted from the  
430 whole image. To better understand our framework and hashing algorithms,  
431 we also present a retrieval example in Fig 5, which further exhibits better  
432 performance of GH and SGH than the other hashing methods. Note that the  
433 similarity between the query and training images is calculated by using GTG,  
434 which determines the weight of each query cell. Usually, the smaller Ham-  
435 ming distance between a query cell and a training image, the larger weight of  
436 the query cell. If the weight of query cells changes, different returned images  
437 might be obtained. Moreover, the number of query cells in one image also  
438 affects the retrieved images.

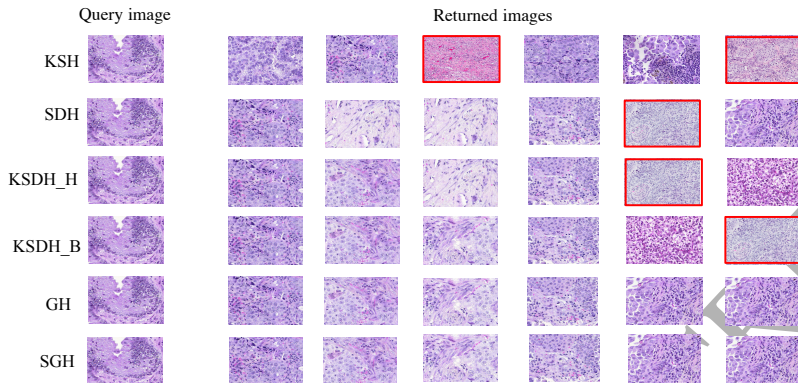


Fig 5: A query image and its top six returned images obtained by using the six hashing algorithms and GIST features. (The error images are marked by a red rectangle block. Additionally, since all hashing algorithms use the strategy GTG, we do not show its name in this figure for simplicity.)

Table 3: Training and query time (second) comparison of different retrieval algorithms (Since all hashing algorithms use the strategy GTG, we do not show its name in this table for simplicity. Training time is the time with all training data learning binary codes, and query time is the querying time of one test image used. ‘All’ means that all cells in the query image are used for retrieval. In addition, hashing algorithms encode each cell into 10-bit binary codes.)

Time	KSH	SDH	KSDH_H	KSDH_B	GH	SGH
Training	$1.7 \times 10^3$	1.6	$3.0 \times 10^1$	$2.0 \times 10^1$	$1.9 \times 10^1$	7.6
Query ( $m = 9$ )	$8.8 \times 10^{-3}$	$8.7 \times 10^{-3}$	$8.8 \times 10^{-3}$	$8.8 \times 10^{-3}$	$8.5 \times 10^{-3}$	$8.5 \times 10^{-3}$
Query (All)	3.1	3.1	3.1	3.1	3.0	3.0

439 Table 3 lists the training time of KSH, SDH, KSDH\_H, KSDH\_B, GH and  
 440 SGH, and their query time as well. KSH has the highest training time cost,  
 441 while SDH takes the lowest. SGH costs lower training time than KSDH\_H,  
 442 KSDH\_B and GH, and in contrast with SDH, its training time cost is accept-  
 443 ably lower. When 9 cells are selected, all algorithms spend less query time  
 444 than all cells being selected. Note that in Table 3, we utilize only one hash  
 445 table to retrieve relevant images, the query time can be reduced if multiple  
 446 hashing tables are used.

### 447 5.3. Retrieval performance vs. query cells

448 Fig 6 presents the classification accuracy, precision and MAP of different  
 449 hashing algorithms with respect to different number  $m$  of cells selected from  
 450 each query image. Fig 6a-c, Fig 6d-f and Fig 6g-i correspond to the cases

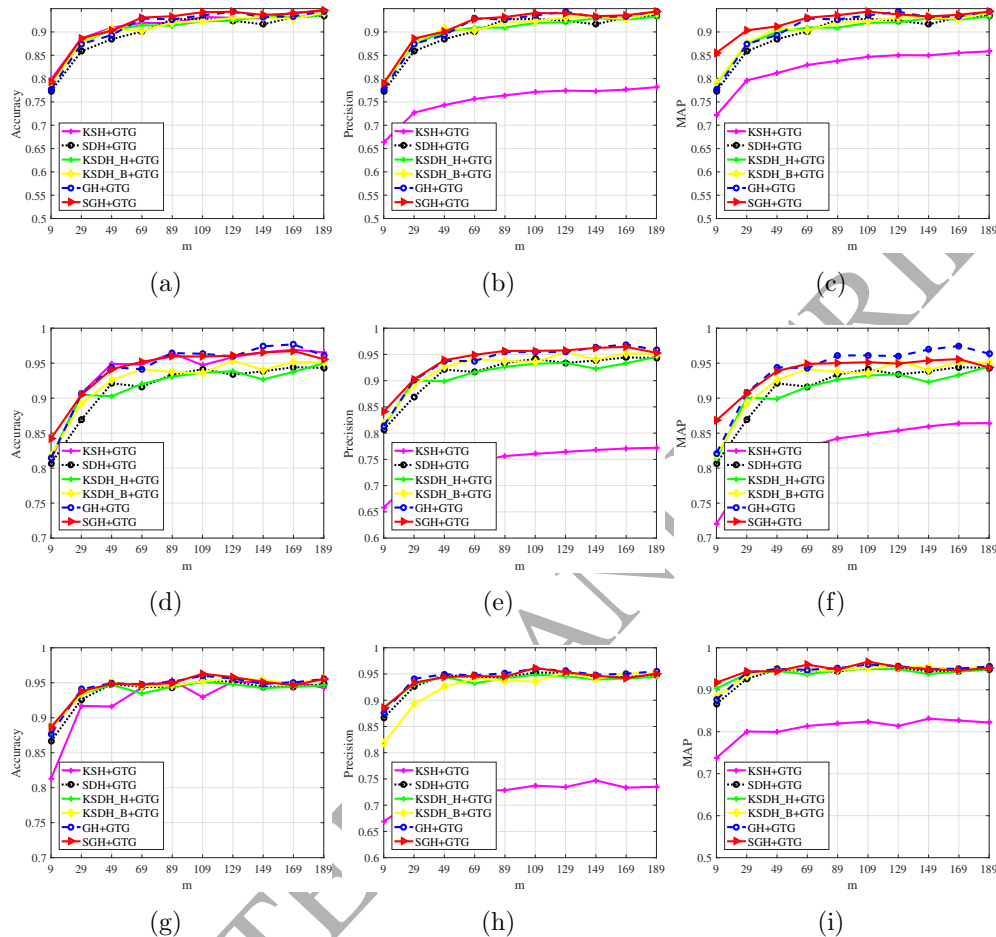


Fig 6: Retrieval accuracy vs the number of selected cells  $m$ . (a)-(c) Accuracy, precision and MAP vs  $m$ , respectively, with HOG features; (d)-(f) accuracy, precision and MAP vs  $m$ , respectively, with GIST features; (g)-(i) accuracy, precision and MAP vs  $m$ , respectively, with CNN features.

451 using HOG, GIST and CNN features, respectively. In Fig 6a-c, the clas-  
 452 sification accuracy, precision and MAP improve with the increasing of  $m$ ,  
 453 and GH+GTG and SGH+GTG perform better than KSH+GTG, KSDH\_H,  
 454 KSDH\_B and SDH+GTG in almost all cases. When  $m=189$ , GH+GTG  
 455 achieves 94.54% classification accuracy, 94.34% precision and 94.42% MAP,  
 456 respectively. SGH+GTG can also obtain the similar results. In Fig 6d-f  
 457 hashing algorithms with GIST features have the similar trends to those with

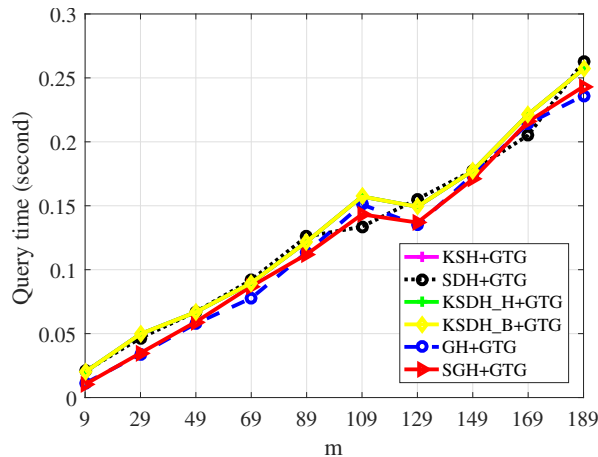


Fig 7: Query time vs the number  $m$  of selected anchors.

458 HOG features, and GH+GTG and SGH+GTG outperform KSH+GTG and  
 459 SDH+GTG. When  $m = 189$ , GH+GTG attains the best classification accu-  
 460 racy (96.06%), precision (95.84%) and MAP (96.36%) among all hashing al-  
 461 gorithms. In Fig 6g-i, the hashing algorithms with CNN features have similar  
 462 trends to that with HOG or GIST features. When  $m = 189$ , GH+GTG at-  
 463 tains the best classification accuracy (95.02%), precision (95.02%) and MAP  
 464 (95.02%) among all hashing algorithms. It is worthy noting that hashing al-  
 465 gorithms with CNN features can achieve their best or sub-best performance  
 466 using less number of selected cells than that with HOG or GIST features.  
 467 Although the hashing algorithms can achieve better results with a larger  
 468 number of cells selected, Fig 7 suggests that a larger  $m$  means a higher query  
 469 time cost. Therefore, in practice, we make a trade-off between the retrieval  
 470 accuracy and query time.

#### 471 5.4. Retrieval performance vs. dimension

472 In this section, we show the retrieval accuracy of various hashing al-  
 473 gorithms with different number  $c$  of dimensions. Fig 8 presents the per-  
 474 formance of various hashing algorithms with  $c$  during  $[2, 20]$ , respectively.  
 475 Fig 8 shows that SDH+GTG and KSDH\_B can achieve robust performance  
 476 are more robust on different number of dimensions with HOG, GIST or  
 477 CNN features; Fig 8 also displays that SGH+GTG and KSDH\_H+GTG ob-  
 478 tain robust performance with HOG or GIST features, while SGH+GTG and  
 479 KSDH\_H+GTG with CNN features achieve bad performance at 2- and 4-bit,  
 480 respectively. GH+GTG achieves robust performance on CNN features, while

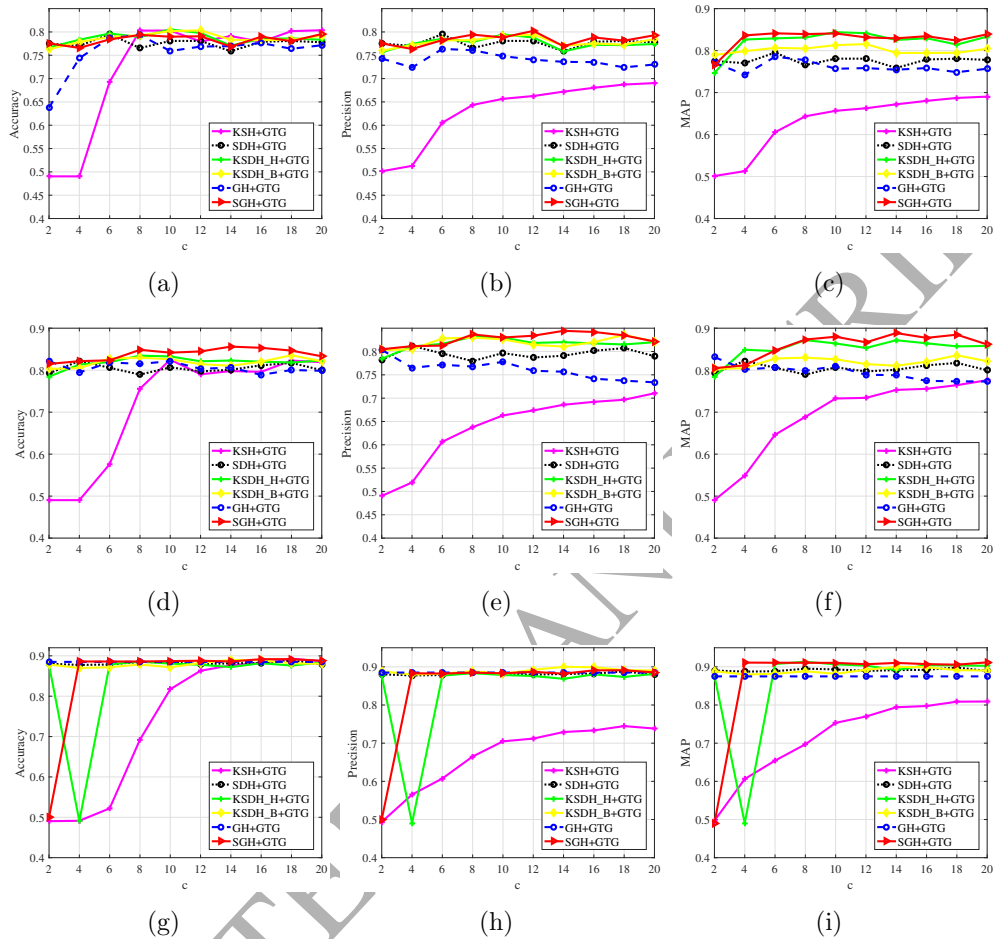


Fig 8: Retrieval accuracy vs the number of dimensions  $c$ . (a)-(c) Accuracy, precision and MAP vs  $c$ , respectively, with HOG features; (d)-(f) accuracy, precision and MAP vs  $c$ , respectively, with GIST features; (g)-(i) accuracy, precision and MAP vs  $c$ , respectively, with CNN features.

481 its performance is slightly changing with different number of dimensions on  
 482 HOG or GIST features. KSH has worse retrieval accuracy on the three types  
 483 of features when  $c$  is small, probably because it utilizes the greedy strategy to  
 484 attain the approximated optimal solutions. When encoding cells into 10-bit  
 485 binary codes, all hashing algorithms except KSH can achieve their best or  
 486 sub-best performance on HOG, GIST and CNN features.

## 487 5.5. Parameters analysis

488 In the proposed algorithm GH, two essential parameters  $\alpha$  and  $\gamma$  affect  
 489 its classification and retrieval performance. In SGH, besides  $\alpha$  and  $\gamma$ ,  $M$  also  
 490 largely affects its performance. We show their classification accuracy on dif-  
 491 ferent parameters with GIST features used in Fig 9. Since  $M$  is determined  
 492 by the number of  $N$ , we use *Rate* to describe the relationship between  $M$   
 493 and  $N$  in Fig 9, where  $Rate = \frac{M}{N}$ . As shown in Fig 9, the range of both  
 494  $\alpha$  and  $\gamma$  is  $[10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}, 1, 10]$  and the range of *Rate* is  
 495 from 2% to 40% with the step 2%. Fig 9a shows the accuracy of GH with  
 496 different  $\alpha$  and  $\gamma$ , which can achieve the best or sub-best accuracy during  
 497 the range  $[10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}, 10^{-1}]$  and  $[10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}]$ , re-  
 498 spectively; Fixing  $Rate = 10\%$ , Fig 9b displays the accuracy of SGH with  
 499 different  $\alpha$  and  $\gamma$ , which can achieve the best or sub-best accuracy during the  
 500 range  $[10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}]$  and  $[10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}]$ , respec-  
 501 tively; Fixing  $\gamma = 10^{-3}$ , Fig 9c presents the accuracy of SGH with different  
 502 *Rate* and  $\alpha$ , which can achieve the best or sub-best accuracy during the range  
 503 from 4% to 28%, and  $[10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}]$ , respectively; Fixing  $\alpha = 10^{-3}$ ,  
 504 Fig 9d shows the accuracy of SGH with different *Rate* and  $\gamma$ , which can  
 505 achieve the best or sub-best accuracy during the range from 4% to 40%, and  
 506  $[10^{-6}, 10^{-5}, 10^{-4}, 10^{-3}, 10^{-2}]$ , respectively. Similar findings can be observed  
 507 on HOG and CNN features. In our experiments, we choose  $\gamma = 10^{-3}$  for both  
 508 GH and SGH, set  $\alpha = 10^{-4}$  for GH, and select  $\alpha = 10^{-3}$  and  $Rate = 10\%$   
 509 for SGH. Using these parameters, the proposed methods GH and SGH can  
 510 obtain the approximately best and robust performance.

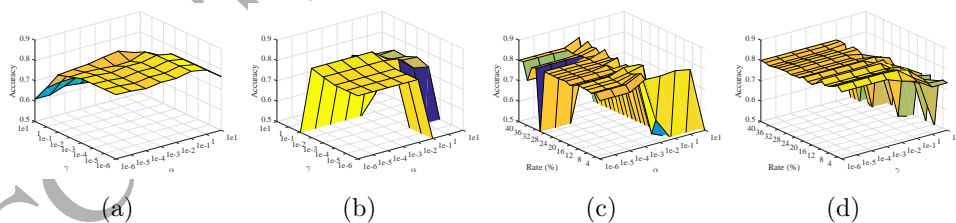


Fig 9: Classification accuracy on different parameters with GIST features. (a) Classification accuracy of GH on parameters  $\alpha$  and  $\gamma$ ; (b) Accuracy of SGH on parameters  $\alpha$  and  $\gamma$  with  $M = 0.1N$ ; (c) Accuracy of SGH on parameters *Rate* and  $\alpha$  with  $\gamma = 10^{-3}$ ; (d) Accuracy of SGH on parameters *Rate* and  $\gamma$  with  $\alpha = 10^{-3}$ .

511 *5.6. Discussion and future work*

512 Based on experiments on the lung cancer database, we can observe that:  
 513 (1) SVM+MV, SDH+MV and GH+MV obtain reasonable image classifica-  
 514 tion accuracy, but NN+MV and KSH+MV sometimes perform worse than  
 515 NN and KSH. (2) The performance of all eight algorithms is affected by  
 516 extracted features. GIST obtains better performance than HOG and CNN  
 517 features when all cells are selected, although CNN features can achieve bet-  
 518 ter performance than GIST and HOG features when 9 cells are selected. (3)  
 519 The GTG matching can be effectively combined with the hashing algorithm  
 520 for image retrieval and exhibits superior performance to MV. (4) Among six  
 521 hashing algorithms, SGH and GH achieve superior retrieval performance to  
 522 KSH, SDH, KSDH.H and KSDH.B in many cases; (5) Although all hashing  
 523 algorithms can quickly retrieve the most relevant images when 9 cells are  
 524 selected, the query time is relatively high when all cells are selected. The  
 525 main reasons are summarized as follows:

- 526 • SVM+MV, SDH+MV and GH+MV provide promising performance,  
 527 because their classification accuracy on each type of cells is larger than  
 528 50%. Based on Eq. 1 and Fig 2, for  $p \in (0.5, 1)$ , the classification  
 529 accuracy will increase when the number of cells grows. By contrast,  
 530 the classification accuracy will decrease if  $p \in (0, 0.5)$ , and this is the  
 531 main reason why NN+MV and KSH+MV has inferior performance to  
 532 NN and KSH.
- 533 • GIST and HOG are much faster than some other popular feature ex-  
 534 traction methods like SIFT, and thus they are more suitable for tackling  
 535 the large-scale cell problem. Usually, GIST extracting features via de-  
 536 scribing the entire cell image achieves better performance than HOG  
 537 extracting cell texture and edge information. It suggests that global in-  
 538 formation of cell images is more important than their texture and edge  
 539 information on cell and image classification, probably because global  
 540 information contains the shape and area of the cell, which is beneficial  
 541 to the cell and image classification. CNN features show superior per-  
 542 formance to GIST and HOG features when a small number of cells are  
 543 selected, while with all cells selected, they exhibit slightly worse perfor-  
 544 mance than GIST and even HOG features, probably because: (1) CNN  
 545 features contain high-level information, while HOG and GIST extract  
 546 low-level information. This leads to better performance of CNN fea-  
 547 tures when few cells are selected; (2) When all cells are selected, a large

548 number of cells in one query image will decrease the gap of accuracy  
 549 obtained by GIST, HOG and CNN features, and CNN features might  
 550 lead to overfitting to some extent.

551 • For GTG, when most of cells of one query image have smaller distances  
 552 to one certain training image than that to the other training images,  
 553 the distance of these two images will be the smallest in most of cases.  
 554 Thus, GTG is effective for image classification and retrieval. Compared  
 555 to MV, GTG is a weighted voting strategy which can better measure  
 556 the similarity between the query image and those in the database, be-  
 557 cause the cells in one query image might have different distance to  
 558 a specific target image. This explains why KSH+GTG, SDH+GTG,  
 559 KSDH\_H+GTG, KSDH\_B+GTG, GH+GTG and SGH+GTG outper-  
 560 form KSH+MV, SDH+MV, KSDH\_H+MV, KSDH\_B+MV, GH+MV  
 561 and SGH+MV, especially for KSH, whose classification accuracy on  
 562 squamous cells is smaller than 50%.

563 • Compared with KSH that uses the symmetric relaxation + greedy strat-  
 564 egy to learn binary codes, GH and SGH directly preserve the discrete  
 565 binary code matrices to reduce the accumulated quantization error be-  
 566 tween the discrete and its relaxed matrices. This might be the main  
 567 reason for the superior performance of GH and SGH to KSH. Addi-  
 568 tionally, GH and SGH can be easily solved and the optimization proce-  
 569 dure quickly converges, and thus they require lower training time costs  
 570 than KSH. Compared with SDH, GH and SGH produce better perfor-  
 571 mance in almost all cases, probably because the performance of SDH  
 572 is largely affected by the selection of kernels, and the graph structure  
 573 in our methods can better preserve the similarity (label) information  
 574 than the regression model used by SDH. Compared with KSDH\_H and  
 575 KSDH\_B, SGH and GH can have superior performance in many cases,  
 576 probably because (1) KSDH\_H and KSDH\_B do not jointly learn bi-  
 577 nary codes and projection matrices, and (2) their performance highly  
 578 relies on the kernels; however, SGH and GH formulate binary code and  
 579 projection matrix learning into a joint optimization problem, and the  
 580 parameters  $\alpha$  and  $\gamma$  can achieve their best performance during a wide  
 581 range.

582 • When all query cells are selected, the time cost of image searching is  
 583 high. To alleviate this issue, currently we can utilize multiple hash-

584 ing tables to reduce time costs. In our future work, we will design a  
585 novel weighted voting algorithm with lower querying time costs and  
586 meanwhile obtaining similar or even better performance. Moreover,  
587 our framework utilizes the segmentation method (Xing et al., 2014) to  
588 crop cell patches, which takes a high relatively time cost for large-scale  
589 image data. One potential solution is to design an end-to-end deep  
590 learning model for efficient cell segmentation, which will be our future  
591 work.

- 592 • In addition to lung images, we also conduct some experiments on other  
593 pathological applications including breast and brain cancers diagnosis,  
594 and our proposed framework can produce good performance. Moreover,  
595 they might be also suitable for 3D/4D pathological images with effective  
596 features (e.g., CNN features), since the proposed framework including  
597 the hashing algorithms are general methods.

## 598 6. Conclusion

599 In this paper, we present a hashing-based image retrieval framework for  
600 pathology image analysis. Firstly, we explain the reason why image classifi-  
601 cation can benefit from cell-level information. Then we present a graph-based  
602 hashing algorithm via asymmetric relaxation to encode each cell into a set of  
603 binary codes. To improve the scalability of the proposed graph hashing al-  
604 gorithm, we further propose a novel algorithm, namely scalable graph-based  
605 hashing. Next, we propose a group-to-group matching method to retrieve  
606 images based on binary codes of cells. Experimental results on lung cancer  
607 images demonstrate the effectiveness and efficiency of our framework. Since  
608 the method group-to-group matching takes relatively high searching time for  
609 the large number of cells in one query image. In the future, we will design  
610 a novel weighted voting method with lower time costs and meanwhile ob-  
611 taining similar or even better performance. Additionally, since training and  
612 query images usually contain many noise cells, which usually affect the per-  
613 formance of the framework, in the future we will focus on selecting robust  
614 cells to further improve the robustness of the proposed framework. Moreover,  
615 based on the proposed framework we will design an end-to-end deep learning  
616 model for efficient and fast cell segmentation so that the total time cost is  
617 decreased.

618 **References**

- 619 Akakin, H. C., Gurcan, M. N., 2012. Content-based microscopic image re-  
620 trieval system for multi-image queries. *IEEE Trans. Inf. Tech. Biomed.*  
621 16 (4), 758–769.
- 622 Basavanhally, A. N., Ganesan, S., Agner, S., Monaco, J. P., Feldman, M. D.,  
623 Tomaszewski, J. E., Bhanot, G., Madabhushi, A., 2010. Computerized  
624 image-based detection and grading of lymphocytic infiltration in her2+  
625 breast cancer histopathology. *IEEE Trans. Biomed. Eng.* 57 (3), 642–653.
- 626 Caicedo, J. C., Cruz, A., Gonzalez, F. A., 2009. Histopathology image classi-  
627 fication using bag of features and kernel functions. In: *Proc. Conf. Artificial*  
628 *Intell. Med. Euro.* Springer, pp. 126–135.
- 629 Caicedo, J. C., González, F. A., Romero, E., 2011. Content-based  
630 histopathology image retrieval using a kernel-based semantic annotation  
631 framework. *J. Biomed. Inf.* 44 (4), 519–528.
- 632 Comaniciu, D., Meer, P., Foran, D. J., 1999. Image-guided decision support  
633 system for pathology. *Mach. Vision App.* 11 (4), 213–224.
- 634 Dalal, N., Triggs, B., 2005. Histograms of oriented gradients for human de-  
635 tection. In: *Proc. Int. Conf. Comput. Vision Pattern Recog.* Vol. 1. IEEE,  
636 pp. 886–893.
- 637 Doyle, S., Agner, S., Madabhushi, A., Feldman, M., Tomaszewski, J., 2008.  
638 Automated grading of breast cancer histopathology using spectral cluster-  
639 ing with textural and architectural image features. In: *IEEE Int. Sym. on*  
640 *Biomed. Imag.* IEEE, pp. 496–499.
- 641 Dundar, M. M., Badve, S., Bilgin, G., Raykar, V., Jain, R., Sertel, O.,  
642 Gurcan, M. N., 2011. Computerized classification of intraductal breast  
643 lesions using histopathological images. *IEEE Trans. Biomed. Eng.* 58 (7),  
644 1977–1984.
- 645 El-Naqa, I., Yang, Y., Galatsanos, N. P., Nishikawa, R. M., Wernick, M. N.,  
646 2004. A similarity learning approach to content-based image retrieval: ap-  
647 plication to digital mammography. *IEEE Trans. Med. Imag.* 23 (10), 1233–  
648 1244.

- 649 Fan, R.-E., Chang, K.-W., Hsieh, C.-J., Wang, X.-R., Lin, C.-J., 2008. Liblin-  
650 ear: A library for large linear classification. *J. Mach. Learn. Res.* 9 (Aug),  
651 1871–1874.
- 652 Gans, J. S., Smart, M., 1996. Majority voting with single-crossing prefer-  
653 ences. *J. Public Economics* 59 (2), 219–237.
- 654 Greenspan, H., Pinhas, A. T., 2007. Medical image categorization and re-  
655 trieval for pacs using the gmm-kl framework. *IEEE Trans. Inf. Tech.*  
656 *Biomed.* 11 (2), 190–202.
- 657 Institute:, N. C., 2013. The cancer genome atals retrieved from. [https://tcga-](https://tcga-data.nci.nih.gov)  
658 [data.nci.nih.gov](https://tcga-data.nci.nih.gov).
- 659 Jiang, M., Zhang, S., Huang, J., Yang, L., Metaxas, D. N., 2015. Joint kernel-  
660 based supervised hashing for scalable histopathological image analysis. In:  
661 *Med. Image Comput. Comput. Assist. Interv.* pp. 366–373.
- 662 Jiang, M., Zhang, S., Huang, J., Yang, L., Metaxas, D. N., 2016. Scalable  
663 histopathological image analysis via supervised hashing with multiple fea-  
664 tures. *Medical Image Analysis* 34, 3–12.
- 665 LeCun, Y., Bengio, Y., Hinton, G., 2015. Deep learning. *Nature* 521 (7553),  
666 436–444.
- 667 Liu, J., Zhang, S., Liu, W., Deng, C., Zheng, Y., Metaxas, D. N., 2016. Scal-  
668 able mammogram retrieval using composite anchor graph hashing with it-  
669 erative quantization. *IEEE Transactions on Circuits and Systems for Video*  
670 *Technology*.
- 671 Liu, W., Mu, C., Kumar, S., Chang, S., 2014. Discrete graph hashing. In:  
672 *Proc. Neur. Inf. Process. Syst.* pp. 3419–3427.
- 673 Liu, W., Wang, J., Ji, R., Jiang, Y., Chang, S., 2012. Supervised hashing  
674 with kernels. In: *Proc. Int. Conf. Comput. Vision Pattern Recog.* pp. 2074–  
675 2081.
- 676 Liu, W., Wang, J., Kumar, S., Chang, S.-F., 2011. Hashing with graphs. In:  
677 *Proc. Int. Conf. Mach. Learn.* pp. 1–8.
- 678 Lowe, D. G., 2004. Distinctive image features from scale-invariant keypoints.  
679 *Int. J. Comput. Vision* 60 (2), 91–110.

- 680 Nie, F., Wang, H., Huang, H., Ding, C., 2011. Unsupervised and semi-  
681 supervised learning via 1-norm graph. In: Proc. Int. Conf. Comput. Vision.  
682 IEEE, pp. 2268–2273.
- 683 Oliva, A., Torralba, A., 2001. Modeling the shape of the scene: A holistic  
684 representation of the spatial envelope. *Int. J. Comput. Vision* 42 (3), 145–  
685 175.
- 686 Penrose, L. S., 1946. The elementary statistics of majority voting. *J. Royal*  
687 *Statist. Society* 109 (1), 53–57.
- 688 Petushi, S., Garcia, F. U., Haber, M. M., Katsinis, C., Tozeren, A., 2006.  
689 Large-scale computations on histology images reveal grade-differentiating  
690 parameters for breast cancer. *Bio. Med. Comput. Med. Imag.* 6 (1), 1.
- 691 Recht, B., Re, C., Tropp, J., Bittorf, V., 2012. Factoring nonnegative matrices  
692 with linear programs. In: Proc. Neur. Inf. Process. Syst. pp. 1214–1222.
- 693 Schnorrenberg, F., Pattichis, C., Schizas, C., Kyriacou, K., 2000. Content-  
694 based retrieval of breast cancer biopsy slides. *Tech. Health Care* 8 (5),  
695 291–297.
- 696 Shen, F., Shen, C., Liu, W., Shen, H., 2015. Supervised discrete hashing. In:  
697 Proc. Int. Conf. Comput. Vision Pattern Recog. pp. 37–45.
- 698 Shi, X., Guo, Z., Lai, Z., Yang, Y., Bao, Z., Zhang, D., 2015. A framework of  
699 joint graph embedding and sparse regression for dimensionality reduction.  
700 *IEEE Trans. Imag. Process.* 24 (4), 1341–1355.
- 701 Shi, X., Guo, Z., Nie, F., Yang, L., You, J., Tao, D., 2016a. Two-dimensional  
702 whitening reconstruction for enhancing robustness of principal component  
703 analysis. *IEEE Trans. Pattern Anal. Mach. Intell.* 38 (10), 2130–2136.
- 704 Shi, X., Xing, F., Cai, J., Zhang, Z., Xie, Y., Yang, L., 2016b. Kernel-  
705 based supervised discrete hashing for image retrieval. In: Proc. Euro. Conf.  
706 Comput. Vision. pp. 419–433.
- 707 Shi, X., Xing, F., Xu, K., Sapkota, M., Yang, L., 2017. Asymmetric discrete  
708 graph hashing. In: Thirty-First AAAI Conference on Artificial Intelligence.

- 709 Tabesh, A., Teverovskiy, M., Pang, H.-Y., Kumar, V. P., Verbel, D., Kot-  
710 sianti, A., Saidi, O., 2007. Multifeature prostate cancer diagnosis and glea-  
711 son grading of histological images. *IEEE Trans. Med. Imag.* 26 (10), 1366–  
712 1378.
- 713 Wang, J., Kumar, S., Chang, S., 2012. Semi-supervised hashing for large-  
714 scale search. *IEEE Trans. Pattern Anal. Mach. Intell.* 34 (12), 2393–2406.
- 715 Wang, J., Liu, W., Kumar, S., Chang, S.-F., 2016. Learning to hash for  
716 indexing big dataa survey. *Proceedings of the IEEE* 104 (1), 34–57.
- 717 Weiss, Y., Torralba, A., Fergus, R., 2009. Spectral hashing. In: *Proc. Neur.*  
718 *Inf. Process. Syst.* pp. 1753–1760.
- 719 Xing, F., Su, H., Neltner, J., Yang, L., 2014. Automatic ki-67 counting using  
720 robust cell detection and online dictionary learning. *IEEE Trans. Bio. Med.*  
721 *Eng.* 61 (3), 859–870.
- 722 Xing, F., Yang, L., 2016. Robust nucleus/cell detection and segmentation in  
723 digital pathology and microscopy images: A comprehensive review. *IEEE*  
724 *Reviews in Biomedical Engineering PP* (99), 1–1.
- 725 Yan, S., Xu, D., Zhang, B., Zhang, H., Yang, Q., Lin, S., 2007. Graph em-  
726 bedding and extensions: a general framework for dimensionality reduction.  
727 *IEEE Trans. Pattern Anal. Mach. Intell.* 29 (1), 40–51.
- 728 Yang, L., Chen, W., Meer, P., Salaru, G., Feldman, M. D., Foran, D. J., 2007.  
729 High throughput analysis of breast cancer specimens on the grid. In: *Proc.*  
730 *Int. Conf. Medical Imag. Comput. Computer-Assist. Interv.* Springer, pp.  
731 617–625.
- 732 Zhang, S., Metaxas, D., 2016. Large-scale medical image analytics: Recent  
733 methodologies, applications and future directions. *Medical Image Analysis*  
734 33, 98–101.
- 735 Zhang, X., Liu, W., Dundar, M., Badve, S., Zhang, S., 2015a. Towards large-  
736 scale histopathological image analysis: Hashing-based image retrieval.  
737 *IEEE Trans. Med. Imag.* 34 (2), 496–506.
- 738 Zhang, X., Su, H., Yang, L., Zhang, S., 2015b. Fine-grained histopathological  
739 image analysis via robust segmentation and large-scale retrieval. In: *Proc.*  
740 *Comput. Vision Pattern Recog.* pp. 5361–5368.

- 741 Zhang, X., Su, H., Yang, L., Zhang, S., 2015c. Weighted hashing with multi-  
742 ple cues for cell-level analysis of histopathological images. In: *Med. Imag.*  
743 *Analysis*. pp. 303–314.
- 744 Zhang, X., Yang, L., Liu, W., Su, H., Zhang, S., 2014. Mining histopatholog-  
745 ical images via composite hashing and online learning. In: *MICCAI* (2).  
746 pp. 479–486.
- 747 Zheng, L., Wetzel, A. W., Gilbertson, J., Becich, M. J., 2003. Design and  
748 analysis of a content-based pathology image retrieval system. *IEEE Trans.*  
749 *Inf. Tech. Biomed.* 7 (4), 249–255.
- 750 Zou, H., Hastie, T., Tibshirani, R., 2006. Sparse principal component anal-  
751 ysis. *J. Comput. Graph. Statist.* 15 (2), 265–286.