

Integrative Analysis of Cellular Morphometric Context Reveals Clinically Relevant Signatures in Lower Grade Glioma

Ju Han^{1,2}, Yunfu Wang^{1,5}, Weidong Cai³, Alexander Borowsky⁴,
Bahram Parvin^{1,2}, and Hang Chang^{1,2}(✉)

¹ Lawrence Berkeley National Laboratory, Berkeley, CA, USA
hchang@lbl.gov

² Department of Electrical and Biomedical Engineering,
University of Nevada, Reno, USA

³ School of Information Technologies, University of Sydney, Sydney, NSW, Australia

⁴ Center for Comparative Medicine, University of California, Davis, CA, USA

⁵ Department of Neurology, Taihe Hospital, Hubei University of Medicine,
Hubei, China

Abstract. Integrative analysis based on quantitative representation of whole slide images (WSIs) in a large histology cohort may provide predictive models of clinical outcome. On one hand, the efficiency and effectiveness of such representation is hindered as a result of large technical variations (e.g., fixation, staining) and biological heterogeneities (e.g., cell type, cell state) that are always present in a large cohort. On the other hand, perceptual interpretation/validation of important multi-variate phenotypic signatures are often difficult due to the loss of visual information during feature transformation in hyperspace. To address these issues, we propose a novel approach for integrative analysis based on cellular morphometric context, which is a robust representation of WSI, with the emphasis on tumor architecture and tumor heterogeneity, built upon cellular level morphometric features within the spatial pyramid matching (SPM) framework. The proposed approach is applied to The Cancer Genome Atlas (TCGA) lower grade glioma (LGG) cohort, where experimental results (i) reveal several clinically relevant cellular morphometric types, which enables both perceptual interpretation/validation and further investigation through gene set enrichment analysis; and (ii) indicate the significantly increased survival rates in one of the cellular morphometric context subtypes derived from the cellular morphometric context.

Keywords: Lower grade glioma · Cellular morphometric context · Cellular morphometric type · Spatial pyramid matching · Consensus clustering · Survival analysis · Gene set enrichment analysis

This work was supported by NIH R01 CA184476 carried out at Lawrence Berkeley National Laboratory.

1 Introduction

Histology sections provide wealth of information about the tissue architecture that contains multiple cell types at different states of cell cycles. These sections are often stained with hematoxylin and eosin (H&E) stains, which label DNA (e.g., nuclei) and protein contents, respectively, in various shades of color. Morphometric aberrations in tumor architecture often lead to disease progression, and it is desirable to quantify indices associated with these aberrations since they can be tested against the clinical outcome, e.g., survival, response to therapy.

For the quantitative analysis of the H&E stained sections, several excellent reviews can be found in [7, 8]. Fundamentally, the trend has been based either on nuclear segmentation and corresponding morphometric representation, or patch-based representation of the histology sections that aids in clinical association. The major challenge for tissue morphometric representation is the large amounts of technical and biological variations in the data. To overcome this problem, recent studies have focused on either fine tuning human engineered features [1, 4, 11, 12], or applying automatic feature learning [5, 9, 15, 16, 19, 20], for robust representation and characterization.

Even though there are inter- and intra- observer variations [6], a trained pathologist always uses rich content (e.g., various cell types, cellular organization, cell state and health), in context, to characterize tumor architecture and heterogeneity for the assessment of disease state. Motivated by the works of [13, 18], we encode cellular morphometric signatures within the spatial pyramid matching (SPM) framework for robust representation (i.e., cellular morphometric context) of WSIs in a large cohort with the emphasis on tumor architecture and tumor heterogeneity, based on which an integrative analysis pipeline is constructed for the association of cellular morphometric context with clinical outcomes and molecular data, with the potential in hypothesis generation regarding the imaging biomarkers for personalized diagnosis or treatment. The proposed approach is applied to the TCGA LGG cohort, where experimental results (i) reveal several clinically relevant cellular morphometric types, which enables both perceptual interpretation/validation and further investigation through gene set enrichment analysis; and (ii) indicate the significantly increased survival rates in one of the cellular morphometric context subtypes derived from the cellular morphometric context.

2 Approaches

The proposed approach starts with the construction of cellular morphometric types and cellular morphometric context, followed by integrative analysis with both clinical and molecular data. Specifically, the nuclear segmentation method in [4] was adopted given its demonstrated robustness in the presence of biological and technical variations, where the corresponding nuclear morphometric

descriptors are described in [3], and the constructed cellular morphometric context representations are released on our website¹.

2.1 Construction of Cellular Morphometric Types and Cellular Morphometric Context

For a set of WSIs and corresponding nuclear segmentation results, let M be the total number of segmented nuclei; N be the number of morphometric descriptors extracted from each segmented nucleus, e.g. nuclear size, and nuclear intensity; and \mathbf{X} be the set of morphometric descriptors for all segmented nuclei, where $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_M]^\top \in \mathbb{R}^{M \times N}$. The construction of cellular morphometric types and cellular morphometric context are described as follows,

1. Construct cellular morphometric types (\mathbf{D}), where $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_K]^\top$ are the K cellular morphometric types to be learned by the following optimization:

$$\min_{\mathbf{D}, \mathbf{Z}} \sum_{m=1}^M \|\mathbf{x}_m - \mathbf{z}_m \mathbf{D}\|^2 \quad (1)$$

subject to $\text{card}(\mathbf{z}_m) = 1, |\mathbf{z}_m| = 1, \mathbf{z}_m \succeq 0, \forall m$

where $\mathbf{Z} = [\mathbf{z}_1, \dots, \mathbf{z}_M]^\top$ indicates the assignment of the cellular morphometric type, $\text{card}(\mathbf{z}_m)$ is a cardinality constraint enforcing only one nonzero element of \mathbf{z}_m , $\mathbf{z}_m \succeq 0$ is a non-negative constraint on the elements of \mathbf{z}_m , and $|\mathbf{z}_m|$ is the $L1$ -norm of \mathbf{z}_m . During training, Eq. 1 is optimized with respect to both \mathbf{Z} and \mathbf{D} ; In the coding phase, for a new set of \mathbf{X} , the learned \mathbf{D} is applied, and Eq. 1 is optimized with respect to \mathbf{Z} only.

2. Construct cellular morphometric context vis SPM. This is done by repeatedly subdividing an image and computing the histograms of different cellular morphometric types over the resulting subregions. As a result, the spatial histogram, H , is formed by concatenating the appropriately weighted histograms of all cellular morphometric types at all resolutions. For more details about SPM, please refer to [13].

In our experiment, K is fixed to be 64. Meanwhile, given the fact that each patient may contain multiple WSIs, SPM is applied at a single scale for the convenient construction of cellular morphometric context as well as the integrative analysis at patient level, where both cellular morphometric types and the subtypes of cellular morphometric context are associated with clinical outcomes, and molecular information.

2.2 Integrative Analysis

The construction of cellular morphometric context at patient level in a large cohort enables the integrative analysis with both clinical and molecular information, which contains the components as follows,

¹ <http://bmihub.org/project/tcgalggcellularmorphcontext>.

1. Identification of cellular morphometric subtypes/clusters: consensus clustering [14] is performed for identifying subtypes/clusters across patients. The input of consensus clustering are the cellular morphometric context at the patient level. Consensus clustering aggregates consensus across multiple runs for a base clustering algorithm. Moreover, it provides a visualization tool to explore the number of clusters in the data, as well as assessing the stability of the discovered clusters.
2. Survival analysis: Cox proportional hazards (PH) regression model is used for survival analysis.
3. Enrichment analysis: Fisher’s exact test is used for the enrichment analysis between cellular morphometric context subtypes and genomic subtypes.
4. Genomic association: linear models are used for assessing differential expression of genes between subtypes of cellular morphometric context, and the correlation between genes and cellular morphometric types.

3 Experiments and Discussion

The proposed approach has been applied on the TCGA LGG cohort, including 215 WSIs from 209 patients, where the clinical annotation of 203 patients are available. For the quality control purpose, background and border portions of each whole slide image were detected and removed from the analysis.

3.1 Phenotypic Visualization and Integrative Analysis of Cellular Morphometric Types

The TCGA LGG cohort consists of ~ 80 million segmented nuclear regions, from which 2 million were randomly selected for construction of cellular morphometric types. As described in Sect. 2, the cellular morphometric context representation for each patient is a 64-dimensional vector, where each dimension represents the normalized frequency of a specific cellular morphometric type appearing in the WSIs of the patient. Initial integrative analysis is performed by linking individual cellular morphometric types to clinical outcomes and molecular data. Each cellular morphometric type is chosen as the predictor variable in the Cox proportional hazards (PH) regression model together with the age of the patient (implemented through the R *survival* package). For each cellular morphometric type, the frequencies are further correlated with the gene expression values across all patients. The top-ranked genes of positive correlation and negative correlation, respectively, are imported into the MSigDB [17] for gene set enrichment analysis. Table 1 summarizes cellular morphometric types that best predict the survival distribution, and the corresponding enriched gene sets. Figure 1 shows the top-ranked examples for these cellular morphometric types.

As shown in Table 1, 8 out of 64 cellular morphometric types are clinically relevant to survival (FDR adjusted p-value < 0.01) with statistical significance. The first four cellular morphometric types in Fig. 1 all have a hazard ratio > 1 , indicating that a higher frequency of these cellular morphometric types may lead

Table 1. Top cellular morphometric types for predicting the survival distribution based on the Cox proportional hazards (PH) regression model, and the corresponding enriched gene sets with respect to genes that best correlate the frequency of the cellular morphometric type appearing in the WSIs of the patient, positively or negatively. Hazard ratio (HR) is the ratio of the hazard rates corresponding to the conditions with a unit difference of an explanatory variable, and higher HR indicates higher hazard of death.

Type	p-value	q-value	Hazard ratio	Enriched gene sets
Worse prognosis				
#5	$7.25e^{-4}$	$7.73e^{-3}$	$3.47e^4$	
#28	$2.05e^{-5}$	$4.37e^{-4}$	$9.32e^3$	Negatively correlated with: <i>genes up-regulated in response to IFNG</i> ; genes up-regulated in response to alpha interferon proteins
#39	$8.57e^{-7}$	$2.74e^{-5}$	$5.07e^3$	Positively correlated with: genes encoding proteins involved in oxidative phosphorylation; genes up-regulated during unfolded protein response, a cellular stress response related to the endoplasmic reticulum; genes involved in DNA repair Negatively correlated with: genes involved in the G2/M checkpoint, as in progression through the cell division cycle; genes important for mitotic spindle assembly; genes defining response to androgens; genes up-regulated by activation of the PI3K/AKT/mTOR pathway
#43	$1.57e^{-9}$	$1.00e^{-7}$	$9.40e^3$	Negatively correlated with: genes up-regulated by activation of Notch signaling
Better prognosis				
#29	$3.01e^{-4}$	$3.85e^{-3}$	$1.74e^{-8}$	Positively correlated with: <i>genes up-regulated by IL6 via STAT3</i> ; genes defining inflammatory response; <i>genes up-regulated in response to IFNG</i> ; <i>genes regulated by NF-kB in response to TNF</i> ; <i>genes up-regulated in response to TGFB1</i> ; genes up-regulated in response to alpha interferon proteins; genes involved in DNA repair; genes mediating programmed cell death (apoptosis) by activation of caspases; genes up-regulated through activation of mTORC1 complex; genes involved in p53 pathways and networks
#31	$1.23e^{-4}$	$1.96e^{-3}$	$5.49e^{-12}$	Positively correlated with: genes encoding components of the complement system, which is part of the innate immune system; genes up-regulated by KRAS activation; <i>genes up-regulated by IL6 via STAT3</i>
#46	$1.17e^{-3}$	$9.84e^{-3}$	$1.07e^{-8}$	Positively correlated with: a subgroup of genes regulated by MYC; genes defining response to androgens; genes involved in DNA repair; genes encoding cell cycle related targets of E2F transcription factors
#52	$1.23e^{-3}$	$9.84e^{-3}$	$6.86e^{-11}$	Positively correlated with: genes up-regulated during transplant rejection; genes up-regulated during formation of blood vessels; <i>genes up-regulated in response to IFNG</i> ; <i>genes regulated by NF-kB in response to TNF</i> ; <i>genes up-regulated in response to TGFB1</i> ; <i>genes up-regulated by IL6 via STAT3</i> ; genes mediating programmed cell death (apoptosis) by activation of caspases

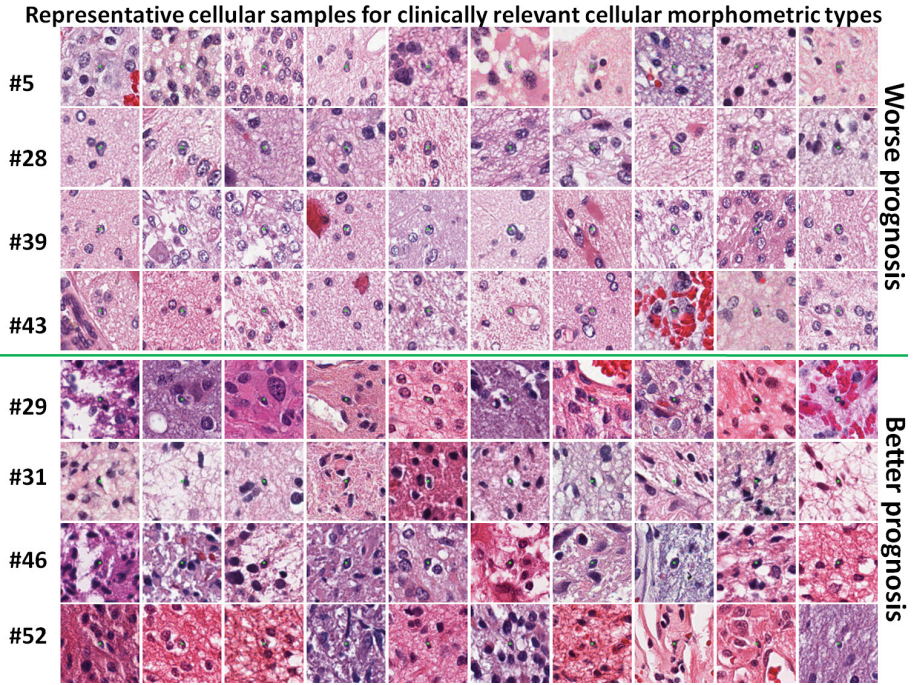


Fig. 1. Top-ranked examples for cellular morphometric types that best predict the survival distribution, as shown in Table 1. Each example is an image patch of 101×101 pixels centered by the retrieved cell marked with the **green** dot. The first four cellular morphometric types (hazard ratio > 1) indicate a worse prognosis and the last four cellular morphometric types (hazard ratio < 1) indicates a protective effect. Note, this figure is best viewed in color at 400% zoom-in.

to a worse prognosis. A common phenotypic property of these cellular morphometric types is the loss of chromatin content in the nuclear regions, which may be associated with poor prognosis of lower grade glioma. The last four cellular morphometric types in Fig. 1 all have a hazard ratio < 1 , indicating that a higher frequency of these cellular morphometric types may lead to a better prognosis.

Table 1 also indicates the enrichment of *genes up-regulated in response to IFNG* in cellular morphometric types #28, #29 and #52. In the glioma microenvironment, tumor cells and local T cells produce abnormally low levels of IFNG. IFNG acts on cell-surface receptors, and activates transcription of genes that offer potentials in the treatment of brain tumors by increasing tumor immunogenicity, disrupting proliferative mechanisms, and inhibiting tumor angiogenesis [10]. The observations of IFNG as a positive survival factor confirms the prognostic effect of these cellular morphometric types: #28 – negative correlation and worse prognosis; #29 and #52 – positive correlation and better prognosis. Other interesting observations include that three cellular morphometric types of better prognosis are enriched with *genes up-regulated by IL6*

via *STAT3*, and two cellular morphometric types of better prognosis are enriched with genes regulated by *NF- κ B* in response to *TNF* and genes up-regulated in response to *TGFB1*, respectively.

3.2 Subtyping and Integrative Analysis of Cellular Morphometric Context

Hierarchical clustering was adopted as the clustering algorithm for consensus clustering, which is implemented via R Bioconductor *ConsensusClusterPlus* package with χ^2 distance as the distance function. The procedure was run for 500 iterations with a sampling rate of 0.8 on 203 patients, and the corresponding consensus clustering matrices with 2 to 9 clusters are shown in Fig. 2, where the matrices with 2 to 5 clusters reveal different levels of similarity among patients and matrices with 6 to 9 clusters provide little further information. Thus, we use the five-cluster result for integrative analysis with clinical outcomes and genomic signatures, where, due to insufficient patients in subtypes #1 (1 patient) and #2 (2 patients), we focus on the remaining three subtypes.

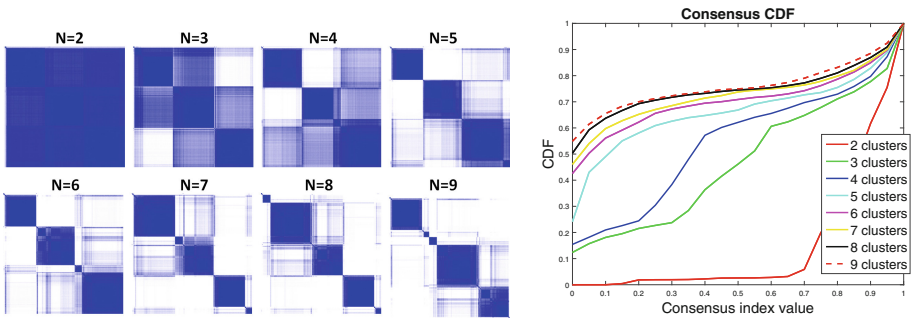


Fig. 2. Consensus clustering matrices and corresponding consensus CDFs of 203 TCGA patients with LGG for cluster number of $N = 2$ to $N = 9$ based on cellular morphometric context.

Figure 3(a) shows the Kaplan-Meier survival plot for three major subtypes of the five-cluster consensus clustering result. The log-rank test p-value of $2.82e^{-5}$ indicates that the difference between survival times of subtype #5 patients and subtypes #3 patients is statistically significant. The integration of genome-wide data from multiple platforms uncovered three molecular classes of lower-grade gliomas that were best represented by IDH and 1p/19q status: wild-type IDH, IDH mutation with 1p/19q codeletion, and IDH mutation without 1p/19q codeletion [2]. Further Fisher's exact test reveals no enrichment between the cellular morphometric subtypes and these molecular subtypes. On the other hand, differential expressed genes between subtype #5 and subtypes #3 (Fig. 3(b)), indicate enrichment of genes that mediate programmed cell death (apoptosis) by activation of caspases, and genes defining epithelial-mesenchymal transition, as in wound healing, fibrosis and metastasis (via MSigDB).

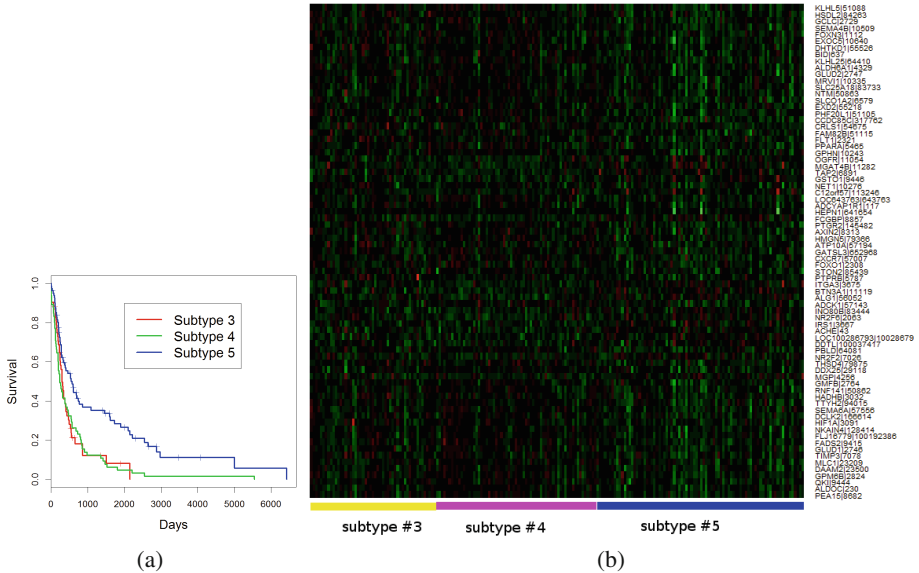


Fig. 3. (a) Kaplan-Meier plot for three major subtypes associated with patient survival, where subtypes #3 (53 patients) #4 (65 patients) and #5 (82 patients) correspond to the three major subtypes from top-left to bottom-right, respectively, in Fig. 2 ($N = 5$). (b) Top genes that are differently expressed between the subtype #5 and subtypes #3.

4 Conclusion and Future Work

In this paper, we encode cellular morphometric signatures within the SPM framework for robust representation (i.e., cellular morphometric context) of WSIs in a large cohort at patient level, based on which an integrative analysis pipeline is constructed for the association of cellular morphometric context with clinical outcomes and molecular data. The integrative analysis, performed on TCGA LGG cohort, reveals clinically relevant cellular morphometric types and morphometric context subtypes, and the corresponding enriched gene sets. We believe that the proposed approach has the potential to contribute to hypothesis generation regarding the imaging biomarkers for personalized diagnosis or treatment, which will be further validated on independent cohort.

References

1. Bhagavatula, R., Fickus, M., Kelly, W., Guo, C., Ozolek, J., Castro, C., Kovacevic, J.: Automatic identification and delineation of germ layer components in *H&E* stained images of teratomas derived from human and nonhuman primate embryonic stem cells. In: IEEE ISBI, pp. 1041–1044 (2010)
2. Cancer Genome Atlas Research Network: Comprehensive, integrative genomic analysis of diffuse lower-grade gliomas. *N. Engl. J. Med.* **372**(26), 2481–2498 (2015)

3. Chang, H., Borowsky, A., Spellman, P.T., Parvin, B.: Classification of tumor histology via morphometric context. In: IEEE CVPR, pp. 2203–2210 (2013)
4. Chang, H., Han, J., Borowsky, A., Loss, L., Gray, J.W., Spellman, P.T., Parvin, B.: Invariant delineation of nuclear architecture in glioblastoma multiforme for clinical and molecular association. *IEEE Trans. Med. Imaging* **32**(4), 670–682 (2013)
5. Chang, H., Zhou, Y., Borowsky, A., Barner, K.E., Spellman, P.T., Parvin, B.: Stacked predictive sparse decomposition for classification of histology sections. *Int. J. Comput. Vis.* **113**(1), 3–18 (2015)
6. Dalton, L., Pinder, S., Elston, C., Ellis, I., Page, D., Dupont, W., Blamey, R.: Histological gradings of breast cancer: linkage of patient outcome with level of pathologist agreements. *Mod. Pathol.* **13**(7), 730–735 (2000)
7. Demir, C., Yener, B.: Automated cancer diagnosis based on histopathological images: a systematic survey (2009)
8. Gurcan, M., Boucheron, L., Can, A., Madabhushi, A., Rajpoot, N., Bulent, Y.: Histopathological image analysis: a review. *IEEE Rev. Biomed. Eng.* **2**, 147–171 (2009)
9. Huang, C.H., Veillard, A., Lomeine, N., Racoceanu, D., Roux, L.: Time efficient sparse analysis of histopathological whole slide images. *Comput. Med. Imaging Graph.* **35**(7–8), 579–591 (2011)
10. Kane, A., Yang, I.: Interferon-gamma in brain tumor immunotherapy. *Neurosurg. Clin. N. Am.* **21**(1), 77–86 (2010)
11. Kong, J., Cooper, L., Sharma, A., Kurk, T., Brat, D., Saltz, J.: Texture based image recognition in microscopy images of diffuse gliomas with multi-class gentle boosting mechanism. In: IEEE ICASSP, pp. 457–460 (2010)
12. Kothari, S., Phan, J.H., Osunkoya, A.O., Wang, M.D.: Biological interpretation of morphological patterns in histopathological whole slide images. In: Proceedings of the ACM Conference on Bioinformatics, Computational Biology and Biomedicine (2012)
13. Lazebnik, S., Schmid, C., Ponce, J.: Beyond bags of features: spatial pyramid matching for recognizing natural scene categories. In: IEEE CVPR, pp. 2169–2178 (2006)
14. Monti, S., Tamayo, P., Mesirov, J., Golub, T.: Consensus clustering: a resampling-based method for class discovery and visualization of gene expression microarray data. *Mach. Learn.* **52**, 91–118 (2003)
15. Romo, D., García-Arteaga, J.D., Arbelez, P., Romero, E.: A discriminant multi-scale histopathology descriptor using dictionary learning. In: SPIE 9041 Medical Imaging (2014)
16. Sirinukunwattana, K., Khan, A.M., Rajpoot, N.M.: Cell words: modelling the visual appearance of cells in histopathology images. *Comput. Med. Imaging Graph.* **42**, 16–24 (2015)
17. Subramanian, A., Tamayo, P., Mootha, V., Mukherjee, S., Ebert, B., Gillette, M., Paulovich, A., Pomeroy, S., Golub, T., Lander, E., Mesirov, J.: Gene set enrichment analysis: a knowledge-based approach for interpreting genome-wide expression profiles. *Proc. Natl. Acad. Sci. USA* **102**(43), 15545–15550 (2005)
18. Yang, J., Yu, K., Gong, Y., Huang, T.: Linear spatial pyramid matching using sparse coding for image classification. In: IEEE CVPR, pp. 1794–1801 (2009)
19. Zhou, Y., Chang, H., Barner, K.E., Parvin, B.: Nuclei segmentation via sparsity constrained convolutional regression. In: IEEE ISBI, pp. 1284–1287 (2015)
20. Zhou, Y., Chang, H., Barner, K.E., Spellman, P.T., Parvin, B.: Classification of histology sections via multispectral convolutional sparse coding. In: IEEE CVPR, pp. 3081–3088 (2014)