

H&E image analysis pipeline for quantifying morphological features



Valeria Ariotta ^{a,1}, Oskari Lehtonen ^{a,1}, Shams Salloum ^{a,c}, Giulia Micoli ^a, Kari Lavikka ^a, Ville Rantanen ^a, Johanna Hynninen ^b, Anni Virtanen ^c, Sampsa Hautaniemi ^{a,*}

^a Research Program in Systems Oncology, Research Programs Unit, Faculty of Medicine, University of Helsinki, 00014 Helsinki, Finland

^b Department of Obstetrics and Gynaecology, University of Turku and Turku University Hospital, 200521 Turku, Finland

^c Department of Pathology, University of Helsinki and HUS Diagnostic Center, Helsinki University Hospital, 00029 Helsinki, Finland

ARTICLE INFO

Keywords:

Digital pathology
Whole-slide images
Instance segmentation
Feature extraction
Ovarian high-grade serous carcinoma
Ploidy

ABSTRACT

Detecting cell types from histopathological images is essential for various digital pathology applications. However, large number of cells in whole-slide images (WSIs) necessitates automated analysis pipelines for efficient cell type detection. Herein, we present hematoxylin and eosin (H&E) Image Processing pipeline (HEIP) for automated analysis of scanned H&E-stained slides. HEIP is a flexible and modular open-source software that performs preprocessing, instance segmentation, and nuclei feature extraction. To evaluate the performance of HEIP, we applied it to extract cell types from ovarian high-grade serous carcinoma (HGSC) patient WSIs. HEIP showed high precision in instance segmentation, particularly for neoplastic and epithelial cells. We also show that there is a significant correlation between genomic ploidy values and morphological features, such as major axis of the nucleus.

Introduction

Histopathological examination of formalin-fixed, paraffin-embedded (FFPE) tissue samples is the cornerstone of cancer diagnosis. The most common staining of the tissue samples is hematoxylin and eosin (H&E), which has been used for more than a century for deducing tumor morphology, cell types, invasion, mitotic activity, and tumor grade.^{1,2} With the development of high-resolution scanners, it has become possible to digitize histopathological samples, which enables the use of machine learning methods on H&E slides. These methods can assist pathologists in diagnostic tasks³ and extract multi-parametric features from the histological phenotype that may not be readily accessible to the human eye.⁴

In recent years, deep learning (DL) methods have been used for various predicting tasks on H&E images without the need to segment and annotate cell types.^{5,6} However, for some computational pathology applications and approaches, such as combining cell morphology to genomics data, it is necessary to extract and annotate cell types from digitalized slides.^{7–10} These approaches enable computational analysis of the morphological features for tens of thousands of cells within a single H&E slide, as well as the spatial distribution of cells.^{11–14} Cell segmentation and annotation tasks are challenging because of the diversity of nuclei characteristics, the presence of overlapping cells, variance in tissue staining, and background noise.

Various methods for cell nuclei classification have been proposed, such as support vector machine¹⁰ and AdaBoost classifiers.¹⁵ Other DL-based approaches have been utilized for nuclei detection, such as the spatially constrained convolutional neural network (CNN)¹⁶ and the multi-task CNN for simultaneous nuclei segmentation and classification.^{17,18} While these approaches perform well on different microscopic image modalities, they lack the necessary flexibility to be trained with a variety of training routines. Additionally, their model architectures lack the flexibility to be adjusted or expanded for inference latency, *i.e.*, the duration between input and output of a model, or segmentation performance gains, making impossible to optimize the latency-performance trade-off of the models. This type of modifiability is necessary in digital pathology, where hundreds of gigapixel-sized whole-slide images (WSIs) are processed.

To address the need of detecting cell types from digitalized H&E slides and extract their morphological features, we developed an open-source computational framework, called H&E Image Processing pipeline (HEIP). HEIP has modular design, which makes it easy to be modified and adjusted to reduce inference latency. The core of HEIP is a modified version of the HoverNet architecture¹⁷ with a post-processing approach that enables the simultaneous segmentation and annotation of cells from digitalized H&E WSIs (subsequently H&E images).

To demonstrate the utility of HEIP, we analyzed H&E images from ovarian high-grade serous carcinoma (HGSC) patients. HGSC is the most

* Corresponding author.

E-mail address: sampsa.hautaniemi@helsinki.fi (S. Hautaniemi).

¹ Equal contribution.

common and aggressive subtype of epithelial ovarian cancer that is typically diagnosed at an advanced stage with widespread metastasis in the peritoneal cavity. Even though most patients have an excellent initial response, the 5-year survival rate in HGSC is less than 40%.¹⁹

Herein, we evaluate HEIP's instance segmentation performance with two HGSC datasets, focusing on cell classification. We also evaluate HEIP's performance in different sites: tubo-ovarian tumors (uterine adnex, ovary, and tubes), and intra-abdominal metastases (omentum and peritoneum). To demonstrate the utility of HEIP, we conducted an exemplifying analysis to explore the association of the morphological nuclear features and the ploidy values, which in a cell correspond to a complete set of chromosomes, computed from whole-genome sequencing data of patients with HGSC.

Material and methods

Patient cohorts

The H&E images used in this study originated from the DECIDER observational clinical trial and PanNuke study.

Firstly, the DECIDER dataset contains image data from HGSC patients participating in the longitudinal, multiregional observational study DECIDER (Multi-layer Data to Improve Diagnosis, Predict Therapy Resistance and Suggest Targeted Therapies in HGSOC; [ClinicalTrials.gov](#) identifier: NCT04846933). The image data used herein consists of scanned images of H&E stained slides from archival formalin-fixed paraffin-embedded (FFPE) tissue blocks collected at the time of diagnosis both for routine diagnostic and research purposes. The archival diagnostic slides were obtained from Auria biobank. The preparation of the research-purpose FFPE block was carried out by the Histology core facility at the Institute of Biomedicine, University of Turku, Finland. All slides were stained at the department of pathology in Turku University Hospital. The scanning of the images was done by Auria Biobank (University of Turku) and the slides were stored in OMERO database.²⁰ The DECIDER data were divided into training and validation datasets (see below).

Secondly, we used the PanNuke dataset²¹ in the training stage. The PanNuke dataset is a publicly available dataset of automatically generated nuclei instance segmentation and classification, from 19 different tissue types and cancer, from more than 20K patches at different magnifications.²¹

Training dataset: For the instance segmentation method, we trained the model using a dataset of 13 H&E images from 13 HGSC patients from the DECIDER cohort. A total of 197 regions of interest (ROIs) were selected from 13 H&E images by a pathologist (A.V.). The ROIs were chosen from various tissue types and had varying dimensions, with a focus on selecting regions that contained different cell types. Subsequently, the cells in the ROIs were annotated by A.V with the train-in-the-loop approach,²² using the software QuPath²³ resulting in 36 093 cell annotations. The cell types included were neoplastic, inflammatory, connective, non-neoplastic epithelial, and dead cells. Additionally, we included 205 343 cell annotations from the PanNuke dataset²¹ in the training dataset.

Validation datasets: The model was validated with 2 subsets of images from the DECIDER cohort. The validation set images were not used in the training stage and were annotated with the train-in-the-loop approach²² by a pathologist (A.V.). The first validation dataset, "CellTypeValidation", was designed to assess instance segmentation performance across the cell types. The second validation dataset, "TumorSiteCellValidation", was designed to assess HEIP performance in different tumor sites.

The CellTypeValidation dataset consisted of 20 human selected ROIs extracted from H&E images of 19 HGSC samples, totaling 9461 train-in-the-loop²² annotated cell instances. The distribution of cell types across the analyzed regions is as follows: 38% of neoplastic cells, 18% of inflammatory cells, 36% of connective cells, 8% of epithelial cells, and 0.1% of dead cells. The majority of neoplastic cells were in ROIs located in the peritoneum, omentum, uterus, mesenterium, and subcutaneous tissue. In contrast, connective cells were more abundant in ROIs from tubo-ovarian

regions, while epithelial cells were more prevalent in ROIs from bowel tissue.

The TumorSiteCellValidation dataset was comprised of 36 ROIs located at the tumor-stroma interface of 18 randomly selected H&E images, including omental (6), peritoneal (6), and tubo-ovarian (6) tumors, from an equal number of HGSC patients. We selected 2 1000 × 1000 pixel ROIs from each H&E image. The distribution of cell types is primarily composed of neoplastic cells (58%), followed by connective cells (26%) and inflammatory cells (16%). Neoplastic cells accounted for over 50% of each tissue type, while connective cells accounted for over 20%, reaching 33% in the case of peritoneum. Dead and epithelial cells were excluded from the analysis as their number in the ROIs was non-existent or too small for reliable analysis.

We also show an example of a possible downstream analysis by calculating correlation between features extracted from images and genomic ploidy values. The ploidy association dataset contains an independent subset of patients in the DECIDER cohort. The samples in ploidy vs. feature correlation analysis were matched, *i.e.*, the H&E image and whole-genome sequencing sample are taken from the adjacent locations of the same tumor piece. We obtained 47 digitalized H&E slides from 23 HGSC patients with this criterion. The H&E images were obtained from omental (18), peritoneal (12), and tubo-ovarian (17) tumors.

Image preprocessing

The H&E images were scanned in MIRAX format with 20× magnification. To prepare the images for analysis, we used a Python library called HistoPrep.²⁴ HistoPrep was employed to identify and segment tissue areas from H&E images into patches. Additionally, patches with insufficient information or a low signal-to-noise ratio were excluded using a series of filters. The H&E images were partitioned into patches with dimensions of 1250 × 1250 pixels, and for each image, the patches were saved in a separate folder in PNG format. The number of patches varied depending on the size of the tissue and the filtering applied, ranging from hundreds to thousands per image.

Deep learning instance segmentation model

A deep learning approach was developed to segment and classify the nuclei. The model is a multi-task CNN, loosely based on the HoVer-net architecture.¹⁷ Similar to HoVer-Net, the architecture comprises a shared encoder and 3 distinct task-specific decoders with distinct output tasks. However, instead of using the post-processing method used by the HoVer-Net model, we opted for the Omnipose post-processing approach²⁵ due to its better overall segmentation performance, as demonstrated in Table S1.

The segmentation and classification performances were evaluated using the following metrics: segmentation quality (SQ), detection quality (DQ), and panoptic quality (PQ). SQ is calculated as the normalized mean of the Intersection over Union (IoU), which measures the quality of the object delineation. DQ, also known as F1-score, is the harmonic mean of precision and recall and measures how well countable objects are detected from the background. PQ is defined as the product of DQ and SQ and quantifies the performance of instance segmentation in a unified manner. The formulas for these metrics are as follows:²⁶

$$DQ = \frac{|TP|}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}$$

$$SQ = \frac{\sum_{(p,g) \in TP} IoU(p,g)}{|TP|}$$

$$PQ = DQ \times SQ = \frac{\sum_{(p,g) \in TP} IoU(p,g)}{|TP| + \frac{1}{2}|FP| + \frac{1}{2}|FN|}$$

where *TP*, *FP*, and *FN* denote the true-positive, false-positive and false-

negative, respectively. *IoU* denotes the intersection-over-union and was set to 0.5.

In order to provide a more precise assessment of the HEIP instance segmentation, we employed estimated confidence intervals (CIs) using a bootstrapping approach.²⁷ Bootstrapping employs resampling the validation dataset multiple times ($n = 200$). We derived confidence bounds by extracting percentiles from the resulting bootstrap distribution.²⁸

After instance segmentation, a json file, in geojson format, is generated containing the coordinates of each detected nucleus as polygons. HEIP then extracts various nuclei features from the json file, including area, volume, solidity, eccentricity, minor axis, major axis, aspect ratio, and perimeter. Additionally, it estimates the percentage of cell types and Shannon index entropy values.²⁹ The definition of each feature is detailed in Table S2.

Whole genome sequencing

We conducted a WGS analysis to investigate the correlation between nuclear cell characteristics, as extracted using HEIP, and ploidy. The approach used for this analysis is consistent with the methodology outlined in the Methods section of Lahtinen et al.³⁰

Copy number calling, ploidy, and purity estimation

Copy number calling was conducted on 23 patients using the Hartwig Medical Foundation toolkit, with genomic breakpoints and breakends extracted using the Genomic Rearrangement Identification Software Suite (GRIDSS).³¹

B-allele frequency (BAF) was calculated with AMBER (<https://github.com/hartwigmedical/hmftools/tree/master/amber/>) using heterozygous single nucleotide polymorphismSNPGATK Mutect2,³² and read depth extracted using COBALT (<https://github.com/hartwigmedical/hmftools/tree/master/cobalt/>). PURITY and PLoidy Estimator (Purple)³³ was used to estimate the copy-number profile, purity, and ploidy by combining BAF, read depth, filtered breakpoints, and somatic mutations.

The model used to calculate purity and ploidy selected the most parsimonious solution among a grid of possible combinations using a fit score. The fit score was determined by a deviation penalty, event penalty multiplier, and somatic deviation penalty. The deviation penalty penalized solutions requiring subclonality to explain copy number patterns, while the event penalty aimed to disfavor the number of alterations required to pass from normal diploid chromosomes to observed minor and major allele

copy numbers. Additionally, combinations of [purity; ploidy] values that violated the rule of somatic variants were penalized.

Statistical analyses

HEIP extracts features for each individual nucleus present in the tissue samples, resulting in data from hundreds or thousands of nuclei features. To summarize the data and provide representative statistical measurements for each sample, we employed the median and variance. Subsequently, the correlation between the median and the variance of each morphological feature (area, volume, major axis, and perimeter) of neoplastic nuclei and the corresponding ploidy value of the samples was computed. The Spearman correlation was used to calculate correlation. Analysis of variance (ANOVA) was used to investigate the correlation between ploidy and the 3 tumor locations: omentum, tubo-ovarian, and peritoneum. All statistical analyses were performed using R software (version 4.2.1).

Results

Overview of the HEIP pipeline

The HEIP pipeline is designed to extract cell nuclei and their morphological nuclear features from H&E images using a DL-based segmentation model as illustrated in Fig. 1. Briefly, the pipeline is based on two customizable steps. The first step processes the H&E images to be amenable for analyses. The second step consists of instance segmentation, which is further divided into cell segmentation and classification steps. Additionally, various Python functions, such as shapely geometry functions, were utilized to extract morphological nuclear features from cell nuclei as well as cell percentages and Shannon Index, which measures entropy. The HEIP pipeline is designed and implemented to be modular and is therefore easy to modify.

Instance segmentation results

Upon visual inspection, the instance segmentation results were very close to the pathologist's ground truth segmentations. Several illustrative cases are shown in Fig. 2. However, we noticed that HEIP tends to make mistakes in detecting very large nuclei, by dividing them into smaller nuclei (Figure S1).

As instance segmentation is arguably the most influential step in the H&E image analysis, we evaluated the HEIP instance segmentation step

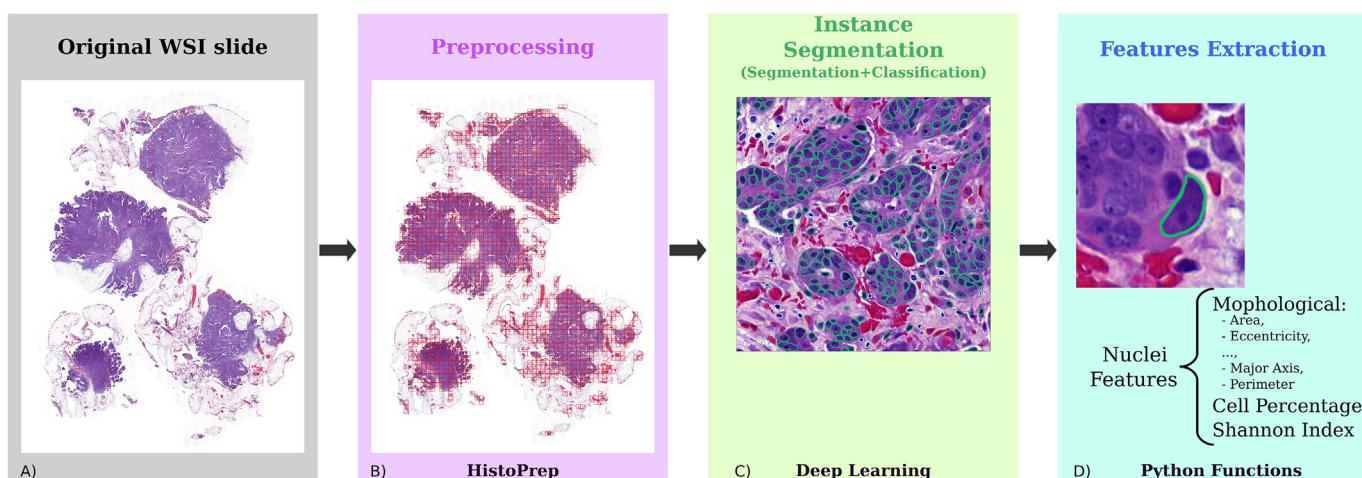


Fig. 1. HEIP schematic workflow. HEIP is a comprehensive software for processing H&E images in order to detect cell nuclei and their morphological features. Panel A: The input to HEIP is a digitized H&E image. Panel B: Preprocessing step is done with HistoPrep. Patches are visible in red. Panel C: Nuclei are detected with deep learning instance segmentation. Panel D: Cell nuclei feature extraction, such as morphological features, cell percentages, and Shannon Index.

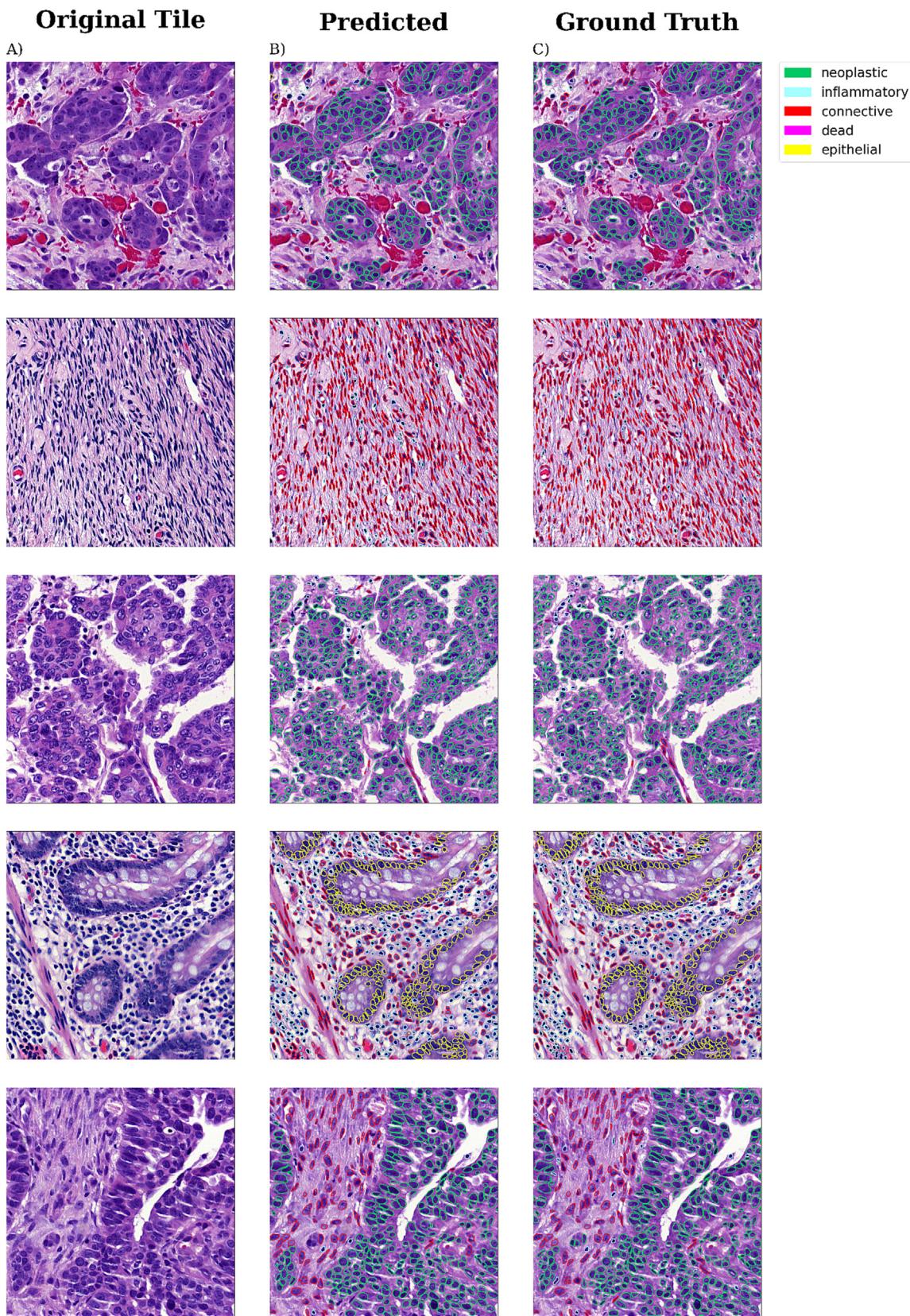


Fig. 2. Instance segmentation examples. Five examples of performance of HEIP, the ROIs were chosen from various tissue types (tubo-ovarian, omentum, bowel, and peritoneum), focusing on different cell types. Panel A: Original tiles from an H&E image. Panel B: Cell classification results by HEIP. Panel C: Ground truth of the nuclei, borders, and types by pathologist.

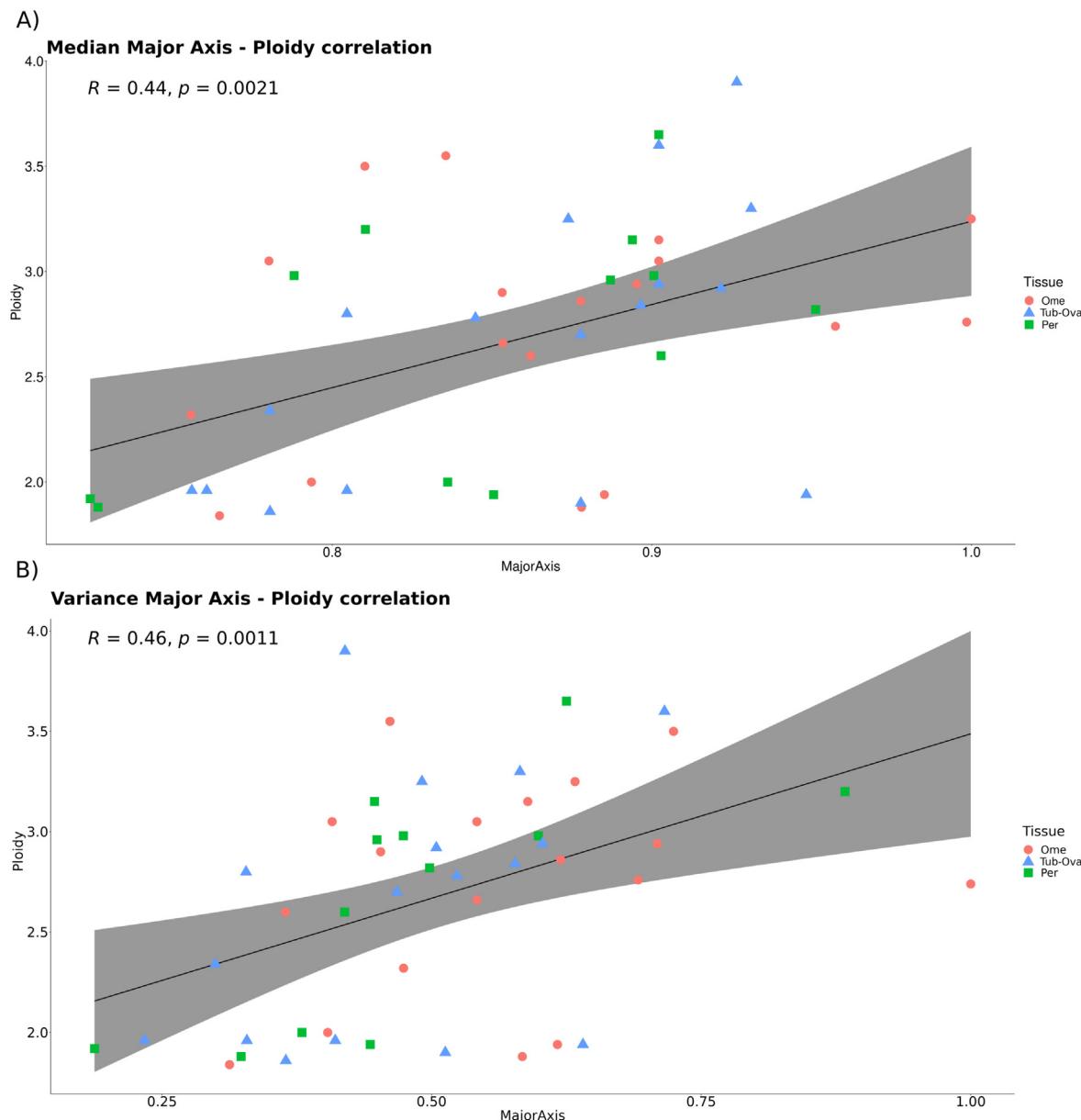


Fig. 3. Correlation between major axis and ploidy values in ploidy association datates. The graph illustrates the correlation between ploidy values and major axis of nuclei across different tissues: omentum (18), peritoneum (12), and tubo-ovarian (17). The ploidy values and H&E image analysis were done using matched samples from the same section. Panel A: Each data point in the plot represents a sample for which we were able to correlate the median value of the major axis of nuclei with its respective ploidy value. Panel B: Each data point in the plot represents a sample for which we were able to correlate the variance value of the major axis of nuclei with its respective ploidy value. A clear positive correlation is observed for both panels. For clarity of understanding the tissue distribution, we have differentiated the three specific tissues in the graph: omentum (Ome), represented by red circles; tubo-ovarian (Tub-Ova), represented by light blue triangles; and peritoneum (Per), represented by green rectangles.

Discussion

Digitalized H&E slides are becoming increasingly important in cancer research.^{1,34} Herein, we have presented HEIP, an automated pipeline for processing H&E images, detecting cell types, and extracting morphological features of the cells, as well as cell percentages and Shannon Index. HEIP is designed and implemented as modular software and trained for HGSC H&E images. Modularity ensures versatility of HEIP to various image analysis tasks with minimal modifications required. Furthermore, the modular design permits easy upgrading to more sophisticated methods as they become available. The output of the nuclei detection is a json file, which contains the polygons with the coordinates of each detected nucleus. By using json files, HEIP reduces the need for memory, compared to the image masks, and storage space, making it efficient.

We showed the utility of HEIP in the analysis of H&E images from histopathological research samples of HGSC patients. Importantly, HEIP estimations for cell type annotations (neoplastic, inflammatory, connective, and epithelial nuclei) agreed well with the pathologist's ground-truth annotations. However, HEIP did not accurately recognize the borders of very large nuclei and tended to divide them into several nuclei. In general, HEIP performance is higher (neoplastic) or on par (connective and inflammatory) with the other nuclei segmentation methods trained with the PanNuke dataset,²¹ whose PQ values range from 0.3 to 0.5. The recognition of dead cells was not involved in the analyses as the number of dead cells in the datasets was negligible.

As an example of a downstream analysis, we explored correlation between ploidy and nuclear morphological features using WGS and histomorphological data from the same tumor piece. Our results indicate

32. McKenna A, Hanna M, Banks E, et al. The Genome Analysis Toolkit: a MapReduce framework for analyzing next-generation DNA sequencing data. *Genome Res* 2010;20(9):1297–1303. <https://doi.org/10.1101/gr.107524.110>.
33. Priestley P, Baber J, Lolkema MP, et al. Pan-cancer whole-genome analyses of metastatic solid tumours. *Nature* 2019;575(7781):210–216. <https://doi.org/10.1038/s41586-019-1689-y>.
34. Bera K, Schalper KA, Rimm DL, Velcheti V, Madabhushi A. Artificial intelligence in digital pathology - new tools for diagnosis and precision oncology. *Nat Rev Clin Oncol* 2019;16(11):703–715. <https://doi.org/10.1038/s41571-019-0252-y>.