

Artifact Evaluation: Privacy-Preserving and Cost-Effective Scheduling for Parallelizing the Large Medical Image Processing Workflow over Hybrid Clouds (PriCE)
Euro-Par 2024 paper 136

Yuandou Wang

14 May, 2024

Question 1: Where to download the gigapixel medical data?

Download a Whole Slide Image (WSI). Download the data from https://surfdrive.surf.nl/files/index.php/apps/files/?dir=/Research%20Datasets/WSI_dataset&fileid=14843054309 and place it at “/PriCE/dataset/iWSI/data/”

Question 2: How to setup experimental environment

Use the terminal for the following steps:

Part (a): Create the environment from the environment.yml file

```
conda env create -f environment.yml
```

Part (b): Activate the new environment

```
conda activate myenv
```

Part (c): Verify that the new environment was installed correctly

```
conda env list
```

Question 3: Folder explanations?

There are the following five main folders in the ZIP file.

- dataset: storing the datasets, e.g., the original WSI example and its intermediate data files, etc.
- inference: storing the CNN inference models.
- pipeline-example for artifact detection: storing the application using CNN inference models for artifact detection in a WSI.
- PriCE-exps: storing the experimental workflows and /or Jupyter notebooks of the PriCE experiments and simulations

Question 4: How it works?

To cope with the diverse image samples of the privacy-preserving data-splitting procedure, we abstract the entire image as a grid graph where different patches with pixel size $p \times p$ are cropped from the original image \mathcal{D} , containing sensitive image labels and objects. The image label contains sensitive coordinate information to reconstruct the image and guide the outcome.

Let $G = (V, E)$ be a graph extracted from the entire patch dataset D cropped from the original image \mathcal{D} . Each patch is represented as a vertex $v \in V$. Two vertices v and μ of V such that $(v, \mu) \in E$ are called to be adjacent. Let $v = (x_i, y_i)$ and $\mu = (x_{i+1}, y_{i+1})$, we denote all possible adjacent relationships between v and μ as: (1) horizontal: $|x_{i+1} - x_i| = p$; (2) vertical: $|y_{i+1} - y_i| = p$; and (3) diagonal: $\sqrt{(x_{i+1} - x_i)^2 + (y_{i+1} - y_i)^2} = \sqrt{2} \times p$. With these characteristics, the positions of the patches can be identified in the original image.

Based on the assumption, if more adjacent patches are placed in the same sub-dataset, the higher the probability that the adversary will restore the entire image. We study different split strategies to scramble these identifications and reduce the risk of restoring the original dataset from the image fragments by the adversary. On the one hand, we adopt the graph-coloring-based split strategies, including 'largest_first', 'random_sequential', 'smallest_last', 'independent_set', 'connected_sequential', 'saturation_largest_first', to split the entire dataset D into different sub-datasets $d_{p,1}, \dots, d_{p,N}$, such that no two adjacent vertices share the same color or dataset. On the other hand, we introduce a random data perturbation to preserve the sensitive coordinates on split datasets' labels by inserting random noise.

PriCE Method

Privacy-preserving Image Splitting with Graph-coloring

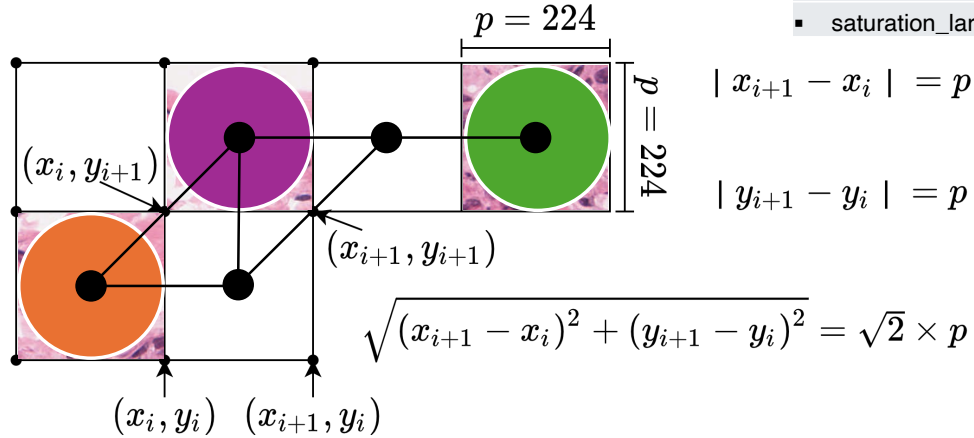


Figure 1: Abstraction of the graph-coloring-based-image-splitting: (a) crop the patches from the original image, e.g., a WSI, (b) identify the position $((x, y)_{\text{coord}})$ of each cropped patch, and (c) identify the adjacent edges: .

Part (a): how to split a gigapixel medical image?

We extract $(x, y)_{\text{coord}}$ as a data matrix A_p of size $(a \times b)$, $a < b$, from $d_{p,k} \subset D$. After normalization, we compute the covariance matrix of the normalized matrix $A_{x,c}$, and then computed the eigenvalues λ and eigenvectors \vec{V} so that we can get the top-k eigenvectors \vec{V}_k to calculate A_e . Moreover, we transform the data into a new coordinate system and encrypt it into datasets $\{d_{e,1}, d_{e,2}, \dots, d_{e,N}\}$. From the perturbed data, since we know the noise variance, we obtain the

estimate coordinates \hat{Y} from decryption by inversely transforming the eigenvector matrix \vec{V} and A_e . Besides, since we know the mappings of original labels and their corresponding encrypted labels, it is easy to measure the output utility.

1. Check the code cells and markdowns in the Jupyter Notebook named `Price/Price-exps/graph_coloring_based_image_splitting.ipynb`
2. Check the code cells and markdowns in the Jupyter Notebook named `Price/Price-exps/evenly_split_w_wo_shuffle.ipynb`

Part (b): how to encrypt/decrypt sensitive information of medical images? How to quantify the privacy-preserving goals?

Check the code cells and markdowns in the Jupyter Notebook named `Price/Price-exps/pertubedata_privacy_risk_evaluation.ipynb` (data perturbation and its privacy-preserving algorithm evaluation)

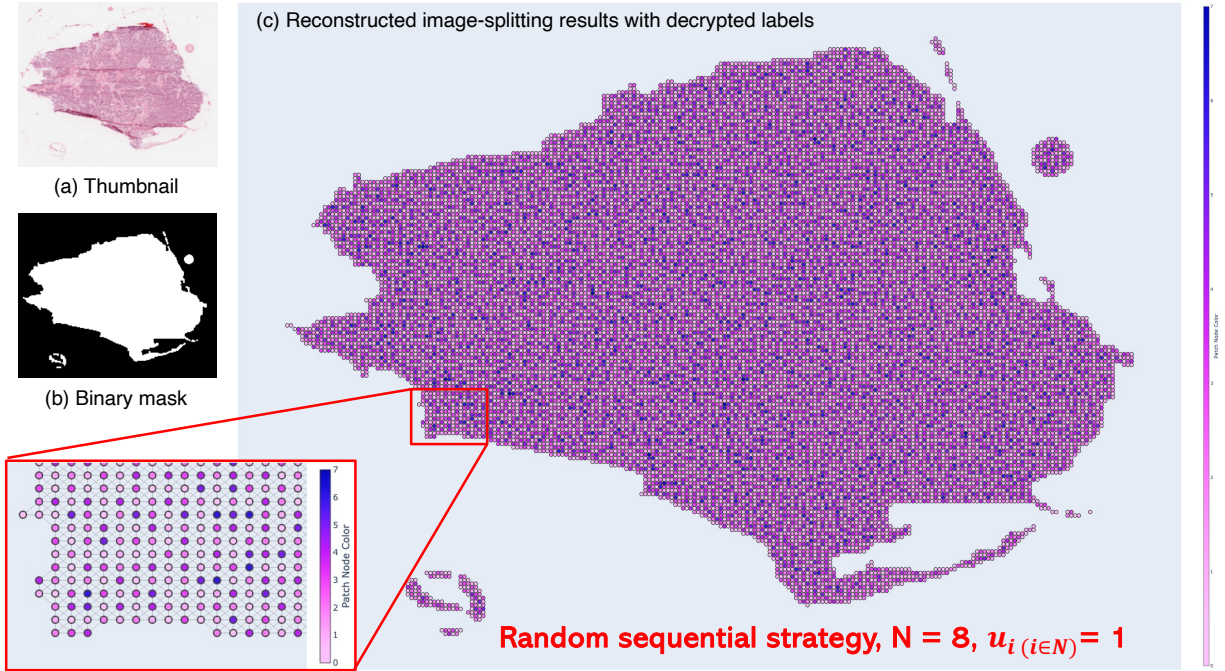


Figure 2: Visualization of the image-splitting: (a) the thumbnail picture of the original medical image, (b) the binary mask picture, and (c) the reconstructed graph from estimated coordinates after decryption.

Part (c): how to seek the 3D Pareto optimal resource planning?

Check the code cells and markdowns in the Jupyter Notebook named `Price/Price-exps/Pareto_3D_evaluation.ipynb`