

# Cracking Instance Jigsaw Puzzles: An Alternative to Multiple Instance Learning for Whole Slide Image Analysis

Xiwen Chen<sup>1\*</sup> Peijie Qiu<sup>2\*</sup> Wenhui Zhu<sup>3\*</sup> Hao Wang<sup>1</sup> Huayu Li<sup>4</sup> Xuanzhao Dong<sup>3</sup>

Xiaotong Sun<sup>5</sup> Xiaobing Yu<sup>2</sup> Yalin Wang<sup>3</sup> Abolfazl Razi<sup>1</sup> Aristeidis Sotiras<sup>2</sup>

1. Clemson University 2. Washington University in St. Louis 3. Arizona State University

4. University of Arizona 5. University of Arkansas

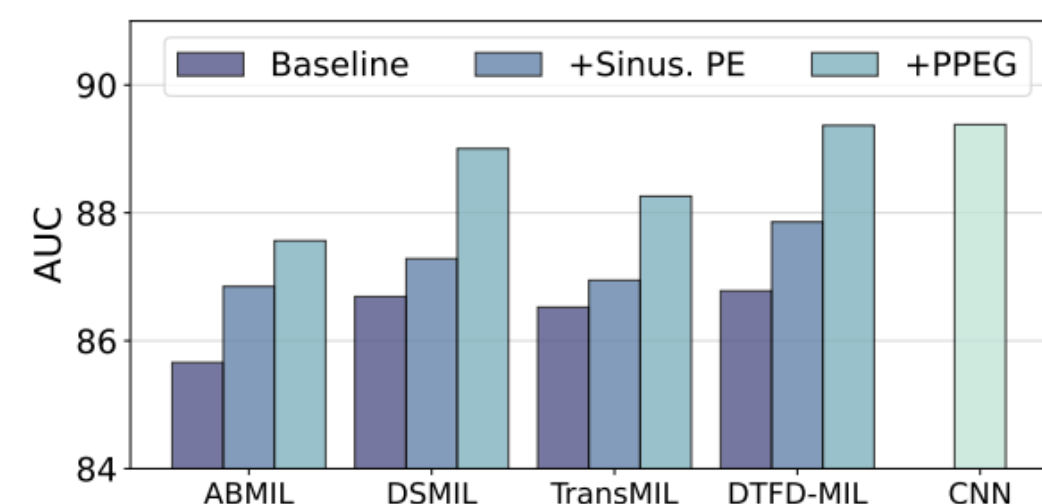
\* equal contributions



Paper

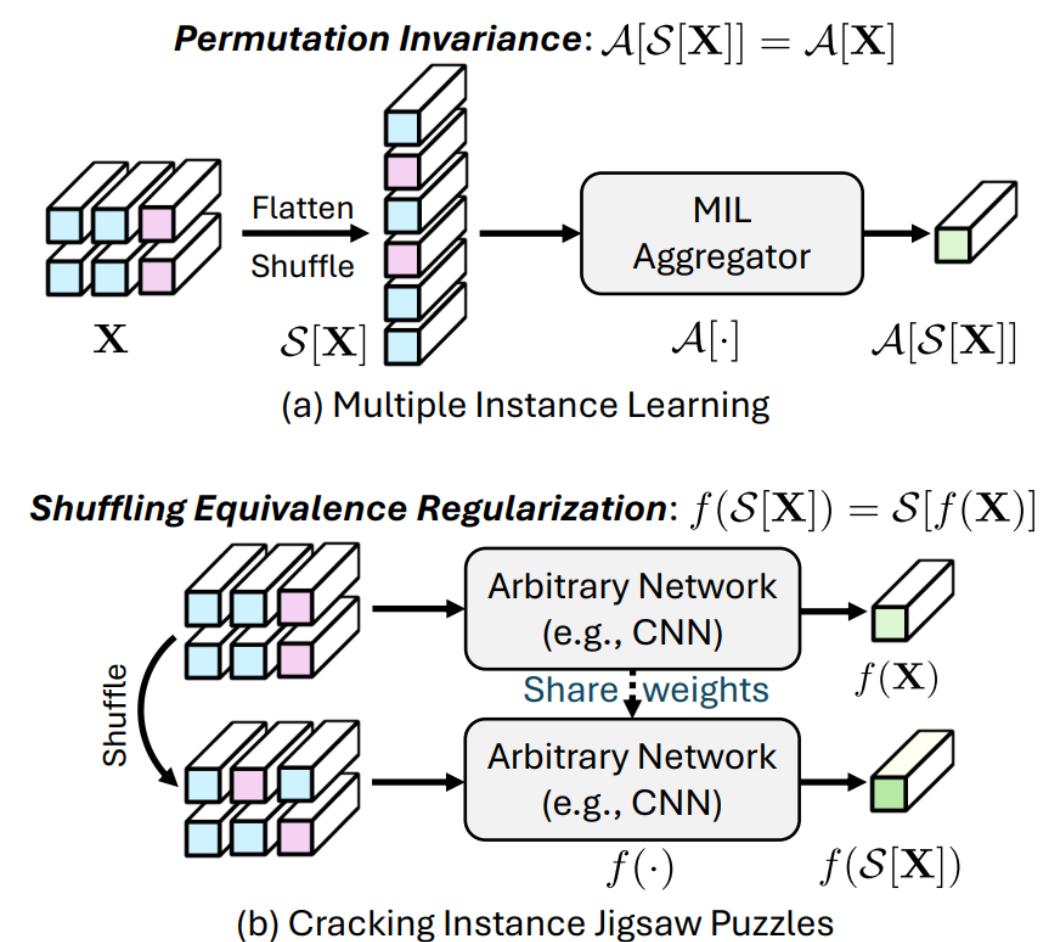
## Motivation

- Permutation invariance** is a core assumption in MIL but it overlooks spatial correlations between tiles in WSIs, which are essential for image-based analysis. **Neighboring tiles** in WSIs are usually spatially and semantically related, often belonging to the same tissue category.
- This creates a **dilemma** between maintaining permutation invariance and preserving spatial correlations.
- Convolutional operations are effective for modeling spatial structure but they are not permutation-invariant because rearranging the inputs changes the output. **Adding convolutional layers** to MIL aggregators often leads to clear performance improvements. **Replacing traditional MIL aggregators with simple CNNs** can even achieve better performance, a direction that has received little attention in the literature.
- These facts lead us to question whether we really need an MIL for WSI analysis. Instead, we argue that, *despite violating the permutation-invariant constraint, modeling the spatial correlations between tiles/instances is necessary.*



## Methods

(a) the traditional MIL method and (b) the proposed method for solving instance jigsaw puzzles. Compared to MIL, the proposed method is advantageous in uncovering semantic correlations between instances.



## Methods

Rather than achieving better modeling instance correlation through different attention mechanisms, naive position encoding e.g., sin-cos, PPEG [1]), or accessing additional information, we discover **that learning to restore the original order from shuffled tiles can lead to better performance by exploiting their semantic correlations.** We term this process as solving an **instance jigsaw puzzle**.

### A Siamese Network Solution

To solve the above instance jigsaw puzzles problem, we propose a Siamese network solution. Specifically, given a shuffling operator  $\mathcal{S}[\cdot]$ , which randomly shuffles input instances, the following equivalence should hold for solving the instance jigsaw puzzles:

$$f_{\theta}(\mathcal{S}[\mathbf{X}]) = \mathcal{S}[f_{\theta}(\mathbf{X})], \quad (2)$$

where  $f_{\theta}$  is parameterized by a neural network, the rationale is that if a network can learn to restore the correct arrangement from shuffled instances, applying the inverse shuffling operation  $\mathcal{S}^{-1}[\cdot]$  to the network's output should recover the original arrangement:  $\mathcal{S}^{-1}[f_{\theta}(\mathcal{S}[\mathbf{X}])] = f_{\theta}(\mathbf{X})$ . Accordingly, we design a shuffling equivalence regularization loss as follows:

$$\mathcal{L}_{\text{Equiv}}(\mathbf{X}) = \frac{1}{2n} \|\mathcal{S}^{-1}[f_{\theta}(\mathcal{S}[\mathbf{X}])] - f_{\theta}(\mathbf{X})\|_2^2, \quad (3)$$

which penalizing the mean squared error between  $f(\mathcal{S}[\mathbf{X}])$  and  $\mathcal{S}^{-1}[f(\mathbf{X})]$ . This equivalence loss is then implemented as a Siamese network with two branches that share the same weights (see Fig. 1(b)). The first branch takes as input the unshuffled instances  $\mathbf{X}$ , while the second branch takes as input the shuffled instances  $\mathcal{S}[\mathbf{X}]$ . The final objective for WSI classification is then a weighted combination of the equivalence loss and the binary cross-entropy loss ( $\mathcal{L}_{\text{BCE}}$ ):

$$\mathcal{L}_{\text{final}}(\mathbf{X}, \mathbf{Y}) = \mathcal{L}_{\text{BCE}}(\mathbf{X}, \mathbf{Y}) + \lambda \mathcal{L}_{\text{Equiv}}(\mathbf{X}), \quad (4)$$

**Theorem 3.** *When approximating the optimal transport plan  $\mathbf{T}_{\#}$  with the inverse shuffling operation  $\mathcal{S}^{-1}$ , the proposed shuffling equivalence regularization is the solution to the inverse optimal transport problem.*

## Experiment and Ablation

Table 1. Main results on the CAMELYON16 dataset and TCGA-NSCLC dataset by using different feature extractors. Our method significantly outperforms all MIL-based competitors (see [Appendix E](#) for the statistical test).

	CAMELYON16			TCGA-NSCLC		
	Accuracy	F1	AUC	Accuracy	F1	AUC
	<b>Swin-ViT ImageNet Pretrained</b>					
ABMIL ( <i>ICML'18</i> )	84.73 <sub>0.85</sub>	83.20 <sub>0.81</sub>	85.66 <sub>1.76</sub>	91.07 <sub>1.08</sub>	91.27 <sub>1.23</sub>	95.88 <sub>1.18</sub>
DSMIL ( <i>CVPR'21</i> )	84.42 <sub>1.12</sub>	82.72 <sub>1.15</sub>	86.69 <sub>2.33</sub>	90.98 <sub>1.49</sub>	90.97 <sub>1.49</sub>	95.71 <sub>0.18</sub>
TransMIL ( <i>NeurIPS'21</i> )	85.04 <sub>1.70</sub>	83.72 <sub>1.29</sub>	88.26 <sub>0.88</sub>	89.73 <sub>0.40</sub>	89.93 <sub>0.62</sub>	95.66 <sub>0.99</sub>
MaxS ( <i>CVPR'22</i> )	84.57 <sub>1.22</sub>	78.87 <sub>1.13</sub>	89.69 <sub>1.25</sub>	87.33 <sub>1.00</sub>	87.05 <sub>1.31</sub>	93.09 <sub>0.85</sub>
AFS ( <i>CVPR'22</i> )	79.61 <sub>2.22</sub>	72.18 <sub>0.95</sub>	83.88 <sub>1.68</sub>	90.79 <sub>1.52</sub>	90.36 <sub>1.80</sub>	96.17 <sub>0.89</sub>
MaxMinS ( <i>CVPR'22</i> )	83.80 <sub>1.01</sub>	76.73 <sub>1.29</sub>	86.78 <sub>1.59</sub>	89.83 <sub>0.87</sub>	89.44 <sub>1.22</sub>	95.70 <sub>0.57</sub>
ILRA-MIL ( <i>ICLR'23</i> )	84.96 <sub>1.05</sub>	83.60 <sub>0.86</sub>	87.76 <sub>1.45</sub>	90.69 <sub>1.13</sub>	90.68 <sub>1.13</sub>	95.56 <sub>0.97</sub>
MHIM-MIL ( <i>ICCV'23</i> )	86.24 <sub>1.68</sub>	84.35 <sub>2.15</sub>	86.12 <sub>1.95</sub>	89.64 <sub>1.66</sub>	89.61 <sub>1.67</sub>	93.93 <sub>0.84</sub>
DGR-MIL ( <i>ECCV'24</i> )	87.60 <sub>2.39</sub>	86.47 <sub>2.39</sub>	88.19 <sub>1.73</sub>	90.88 <sub>1.83</sub>	90.85 <sub>1.84</sub>	95.81 <sub>1.25</sub>
AC-MIL ( <i>ECCV'24</i> )	86.24 <sub>1.01</sub>	84.94 <sub>1.31</sub>	87.77 <sub>1.61</sub>	90.50 <sub>1.29</sub>	90.63 <sub>1.23</sub>	95.61 <sub>0.79</sub>
<b>Ours [Trans.]</b>	<b>89.53<sub>1.40</sub></b>	<b>88.57<sub>1.53</sub></b>	<b>92.17<sub>0.49</sub></b>	<b>92.32<sub>1.25</sub></b>	<b>92.31<sub>1.26</sub></b>	<b>96.40<sub>0.77</sub></b>
<b>Ours [CNN]</b>	<b>88.11<sub>0.58</sub></b>	<b>87.11<sub>0.63</sub></b>	<b>91.80<sub>0.26</sub></b>	<b>92.51<sub>0.80</sub></b>	<b>92.49<sub>0.81</sub></b>	<b>96.32<sub>0.67</sub></b>
	<b>ResNet-18 ImageNet Pretrained</b>					
ABMIL ( <i>ICML'18</i> )	85.74 <sub>0.99</sub>	84.21 <sub>1.11</sub>	85.91 <sub>1.53</sub>	88.10 <sub>0.80</sub>	88.18 <sub>0.82</sub>	93.88 <sub>1.11</sub>
DSMIL ( <i>CVPR'21</i> )	84.19 <sub>2.25</sub>	82.21 <sub>2.82</sub>	84.84 <sub>1.74</sub>	88.58 <sub>1.02</sub>	88.61 <sub>1.06</sub>	93.73 <sub>0.87</sub>
TransMIL ( <i>NeurIPS'21</i> )	82.79 <sub>1.89</sub>	76.63 <sub>1.86</sub>	87.71 <sub>1.84</sub>	84.65 <sub>1.11</sub>	84.20 <sub>0.90</sub>	90.71 <sub>1.20</sub>
MaxS ( <i>CVPR'22</i> )	84.81 <sub>2.09</sub>	83.62 <sub>2.20</sub>	87.22 <sub>1.78</sub>	88.39 <sub>0.81</sub>	88.54 <sub>1.00</sub>	93.43 <sub>0.84</sub>
AFS ( <i>CVPR'22</i> )	81.94 <sub>1.55</sub>	77.85 <sub>1.45</sub>	89.23 <sub>1.07</sub>	88.48 <sub>0.81</sub>	88.27 <sub>1.16</sub>	94.83 <sub>0.92</sub>
MaxMinS ( <i>CVPR'22</i> )	82.02 <sub>1.86</sub>	76.11 <sub>0.88</sub>	88.04 <sub>1.84</sub>	87.81 <sub>0.86</sub>	87.51 <sub>0.00</sub>	94.19 <sub>0.95</sub>
ILRA-MIL ( <i>ICLR'23</i> )	87.08 <sub>2.31</sub>	86.19 <sub>2.56</sub>	89.30 <sub>2.98</sub>	88.77 <sub>0.98</sub>	88.81 <sub>0.99</sub>	94.25 <sub>0.68</sub>
MHIM-MIL ( <i>ICCV'23</i> )	86.05 <sub>1.64</sub>	84.48 <sub>1.82</sub>	86.17 <sub>1.76</sub>	87.43 <sub>1.37</sub>	87.41 <sub>1.35</sub>	93.65 <sub>0.62</sub>
DGR-MIL ( <i>ECCV'24</i> )	86.63 <sub>0.85</sub>	85.25 <sub>0.96</sub>	88.20 <sub>1.30</sub>	87.43 <sub>1.18</sub>	87.43 <sub>1.14</sub>	93.88 <sub>0.41</sub>
AC-MIL ( <i>ECCV'24</i> )	87.02 <sub>1.49</sub>	85.55 <sub>1.77</sub>	87.56 <sub>2.37</sub>	88.58 <sub>0.69</sub>	88.58 <sub>0.69</sub>	94.31 <sub>1.12</sub>
<b>Ours [Trans.]</b>	<b>87.47<sub>2.12</sub></b>	<b>86.30<sub>2.34</sub></b>	<b>90.44<sub>1.41</sub></b>	<b>88.96<sub>0.97</sub></b>	<b>89.02<sub>0.98</sub></b>	<b>94.98<sub>0.81</sub></b>
<b>Ours [CNN]</b>	<b>88.37<sub>0.45</sub></b>	<b>87.16<sub>0.36</sub></b>	<b>92.92<sub>0.87</sub></b>	<b>90.40<sub>0.98</sub></b>	<b>90.39<sub>0.98</sub></b>	<b>94.93<sub>1.15</sub></b>
	<b>CTransPath Self-supervised Pretrained</b>					
ABMIL ( <i>ICML'18</i> )	94.80 <sub>0.50</sub>	94.39 <sub>0.55</sub>	96.50 <sub>0.67</sub>	93.38 <sub>1.10</sub>	93.36 <sub>1.11</sub>	96.81 <sub>0.63</sub>
DSMIL ( <i>CVPR'21</i> )	94.49 <sub>0.64</sub>	94.08 <sub>0.69</sub>	95.64 <sub>0.56</sub>	94.24 <sub>1.25</sub>	94.22 <sub>1.26</sub>	97.85 <sub>0.69</sub>
TransMIL ( <i>NeurIPS'21</i> )	94.42 <sub>0.58</sub>	92.44 <sub>0.68</sub>	97.34 <sub>0.19</sub>	90.79 <sub>0.72</sub>	90.39 <sub>0.69</sub>	96.22 <sub>0.79</sub>
MaxS ( <i>CVPR'22</i> )	94.96 <sub>1.21</sub>	94.77 <sub>1.15</sub>	97.33 <sub>0.30</sub>	93.86 <sub>0.95</sub>	93.85 <sub>0.94</sub>	97.84 <sub>0.52</sub>
AFS ( <i>CVPR'22</i> )	94.42 <sub>0.68</sub>	92.42 <sub>0.86</sub>	97.14 <sub>0.27</sub>	93.28 <sub>1.04</sub>	92.95 <sub>1.17</sub>	97.81 <sub>0.46</sub>
MaxMinS ( <i>CVPR'22</i> )	95.19 <sub>0.47</sub>	93.38 <sub>0.54</sub>	97.66 <sub>0.44</sub>	93.66 <sub>0.70</sub>	93.34 <sub>0.74</sub>	97.78 <sub>0.39</sub>
ILRA-MIL ( <i>ICLR'23</i> )	94.83 <sub>1.77</sub>	94.45 <sub>1.94</sub>	95.85 <sub>1.00</sub>	93.57 <sub>0.75</sub>	93.56 <sub>0.75</sub>	97.44 <sub>0.56</sub>
MHIM-MIL ( <i>ICCV'23</i> )	94.57 <sub>0.55</sub>	94.16 <sub>0.65</sub>	96.38 <sub>0.61</sub>	93.95 <sub>1.21</sub>	93.94 <sub>1.21</sub>	97.87 <sub>0.53</sub>
DGR-MIL ( <i>ECCV'24</i> )	95.73 <sub>1.16</sub>	95.41 <sub>1.26</sub>	96.30 <sub>0.47</sub>	94.53 <sub>1.26</sub>	94.52 <sub>1.27</sub>	97.87 <sub>0.53</sub>
AC-MIL ( <i>ECCV'24</i> )	95.15 <sub>0.64</sub>	94.88 <sub>0.73</sub>	97.00 <sub>0.69</sub>	94.72 <sub>0.68</sub>	94.72 <sub>0.68</sub>	97.76 <sub>0.76</sub>
<b>Ours [Trans.]</b>	<b>96.64<sub>0.37</sub></b>	<b>96.39<sub>0.38</sub></b>	<b>98.00<sub>0.17</sub></b>	<b>95.20<sub>1.23</sub></b>	<b>95.19<sub>1.23</sub></b>	<b>97.99<sub>0.67</sub></b>
<b>Ours [CNN]</b>	<b>96.25<sub>0.83</sub></b>	<b>95.99<sub>0.88</sub></b>	<b>98.10<sub>0.31</sub></b>	<b>95.11<sub>0.83</sub></b>	<b>95.10<sub>0.84</sub></b>	<b>97.55<sub>0.77</sub></b>

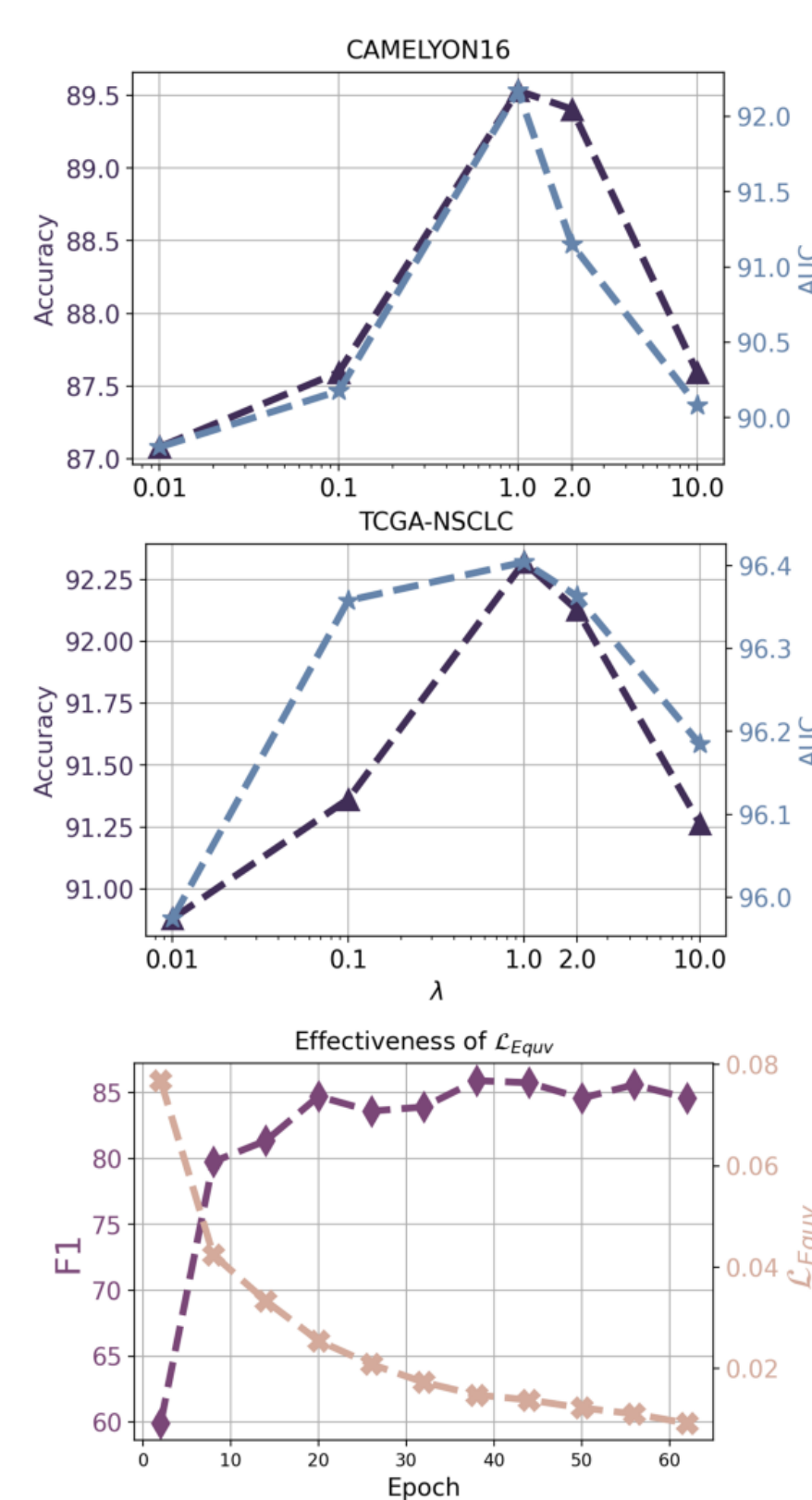


Figure 3. **Top:** The effectiveness of  $\lambda$  on CAMELYON16 datasets, **Middle:** the effectiveness of  $\lambda$  on TCGA-NSCLC datasets, and **Bottom:** Training dynamics for  $\mathcal{L}_{\text{Equiv}}$  and performance.

## Conclusion

We question the necessity of traditional permutation-invariant MIL for whole slide image analysis and highlight the critical role of spatial and semantic correlations among instances. To address this, we propose **“Cracking Instance Jigsaw Puzzles”**, a new paradigm that explicitly learns semantic relationships by restoring the spatial order of shuffled tiles. Our **Siamese network with shuffling equivalence regularization**, theoretically supported by optimal transport principles, consistently outperforms state-of-the-art MIL models in both **WSI classification** and **survival prediction** tasks. Despite a minor computational cost, the approach broadens the design space by enabling flexible architectures such as CNNs, paving the way for more spatially aware and semantically enriched WSI analysis.

## Reference

- [1] Shao, Zhuchen, et al. "Transmil: Transformer based correlated multiple instance learning for whole slide image classification." *Advances in neural information processing systems* 34 (2021): 2136-2147.