

# SGPMIL: Sparse Gaussian Process Multiple Instance Learning

## 1. Additional Experimental Results

	ACC	AUC	ACE
ABMIL	.964 <sup>*</sup> <sub>.010</sub>	<b>.990</b> <sub>.005</sub>	.032 <sup>*</sup> <sub>.011</sub>
CLAM	<u>.978</u> <sub>.007</sub>	<u>.986</u> <sub>.007</sub>	<u>.021</u> <sub>.007</sub>
TransMIL	.962 <sup>*</sup> <sub>.009</sub>	.980 <sup>*</sup> <sub>.004</sub>	.029 <sup>*</sup> <sub>.014</sub>
DGRMIL	.960 <sup>*</sup> <sub>.012</sub>	.980 <sup>*</sup> <sub>.010</sub>	.045 <sup>*</sup> <sub>.016</sub>
BayesMIL	.976 <sub>.007</sub>	.981 <sup>*</sup> <sub>.006</sub>	<b>.020</b> <sub>.007</sub>
MixMIL	.960 <sup>*</sup> <sub>.009</sub>	.978 <sup>*</sup> <sub>.007</sub>	.430 <sup>*</sup> <sub>.002</sub>
AGP	.883 <sup>*</sup> <sub>.021</sub>	.954 <sup>*</sup> <sub>.019</sub>	.069 <sup>*</sup> <sub>.013</sub>
<b>SGPMIL</b>	<b>.980</b> <sub>.007</sub>	<u>.986</u> <sub>.005</sub>	<u>.021</u> <sub>.005</sub>

Table 1. Slide-level performance on CAMELYON16. \* indicates significance based on a one-sided paired t-test across folds ( $p < .05$ ).

	ACC ( $p$ )	AUC ( $p$ )	ACE ( $p$ )
ABMIL	<b>1.19e-3</b>	4.68e-1	<b>2.72e-3</b>
CLAM	8.59e-2	2.40e-1	2.60e-1
TransMIL	<b>6.40e-5</b>	<b>5.66e-3</b>	<b>2.11e-2</b>
DGRMIL	<b>3.15e-4</b>	<b>9.82e-3</b>	<b>1.17e-3</b>
BayesMIL	<b>4.43e-2</b>	<b>3.07e-5</b>	3.79e-1
MixMIL	<b>4.12e-5</b>	<b>4.28e-3</b>	<b>1.23e-18</b>
AGP	<b>7.56e-8</b>	<b>1.09e-4</b>	<b>2.91e-7</b>

Table 2.  $p$ -values from one-sided paired t-tests comparing SGP-MIL with baseline models across folds on CAMELYON16 (slide level). Bold values indicate  $p < .05$ .

	ACC	AUC	ACE
ABMIL	.953 <sub>.003</sub>	.973 <sub>.009</sub>	.039 <sub>.008</sub>
CLAM	.934 <sup>*</sup> <sub>.014</sub>	.953 <sup>*</sup> <sub>.004</sub>	.056 <sub>.016</sub>
TransMIL	.950 <sub>.017</sub>	.970 <sub>.012</sub>	.046 <sub>.019</sub>
DGRMIL	.947 <sub>.024</sub>	<u>.974</u> <sub>.011</sub>	<u>.038</u> <sub>.022</sub>
BayesMIL	<u>.953</u> <sub>.023</sub>	.973 <sub>.021</sub>	<b>.033</b> <sub>.017</sub>
MixMIL	.925 <sup>*</sup> <sub>.015</sub>	.963 <sub>.014</sub>	.410 <sup>*</sup> <sub>.025</sub>
AGP	.948 <sup>*</sup> <sub>.026</sub>	<b>.976</b> <sub>.014</sub>	.048 <sub>.025</sub>
<b>SGPMIL</b>	<b>.955</b> <sub>.021</sub>	.973 <sub>.014</sub>	.047 <sub>.027</sub>

Table 3. Slide-level performance on TCGA-NSCLC. \* indicates significance based on a one-sided paired t-test across folds ( $p < .05$ ).

	ACC ( $p$ )	AUC ( $p$ )	ACE ( $p$ )
ABMIL	2.08e-1	3.14e-1	3.66e-1
CLAM	<b>3.93e-2</b>	<b>2.18e-2</b>	1.50e-1
TransMIL	5.15e-2	1.26e-1	2.83e-1
DGRMIL	6.06e-2	2.59e-1	3.54e-1
BayesMIL	1.61e-1	2.89e-1	4.71e-1
MixMIL	<b>1.19e-2</b>	5.38e-2	<b>1.32e-4</b>
AGP	<b>4.40e-2</b>	4.97e-1	2.22e-1

Table 4.  $p$ -values from one-sided paired t-tests comparing SGP-MIL with baseline models across folds on TCGA-NSCLC (slide level). Bold values indicate  $p < .05$ .

	ACC	$\kappa$	ACE
ABMIL	.834 <sup>*</sup> <sub>.064</sub>	.910 <sup>*</sup> <sub>.028</sub>	.044 <sup>*</sup> <sub>.015</sub>
CLAM	<u>.867</u> <sub>.061</sub>	.927 <sup>*</sup> <sub>.025</sub>	.031 <sup>*</sup> <sub>.018</sub>
TransMIL	.827 <sup>*</sup> <sub>.074</sub>	.911 <sup>*</sup> <sub>.030</sub>	.043 <sup>*</sup> <sub>.021</sub>
DGRMIL	.843 <sup>*</sup> <sub>.097</sub>	<u>.933</u> <sub>.047</sub>	.036 <sup>*</sup> <sub>.025</sub>
BayesMIL	.850 <sup>*</sup> <sub>.060</sub>	.926 <sup>*</sup> <sub>.031</sub>	.031 <sup>*</sup> <sub>.016</sub>
MixMIL	.690 <sup>*</sup> <sub>.054</sub>	.870 <sup>*</sup> <sub>.028</sub>	.180 <sup>*</sup> <sub>.010</sub>
AGP	.802 <sup>*</sup> <sub>.086</sub>	.906 <sup>*</sup> <sub>.047</sub>	<b>.026</b> <sub>.013</sub>
<b>SGPMIL</b>	<b>.900</b> <sub>.065</sub>	<b>.955</b> <sub>.037</sub>	<u>.028</u> <sub>.022</sub>

Table 5. Slide-level performance on PANDA. \* indicates significance based on a one-sided paired t-test across folds ( $p < .05$ ).

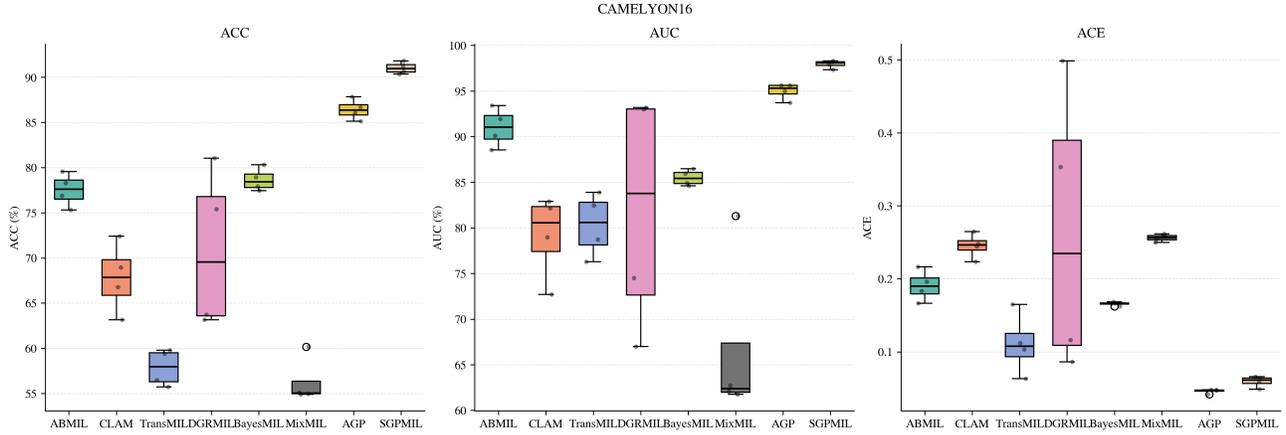


Figure 1. Instance-level performance across folds on CAMELYON16. Boxplots show the distribution of accuracy (ACC), area under the ROC curve (AUC), and adaptive expected calibration error (ACE) for each model. \* denotes significance according to one-sided paired t-tests ( $p < .05$ ).

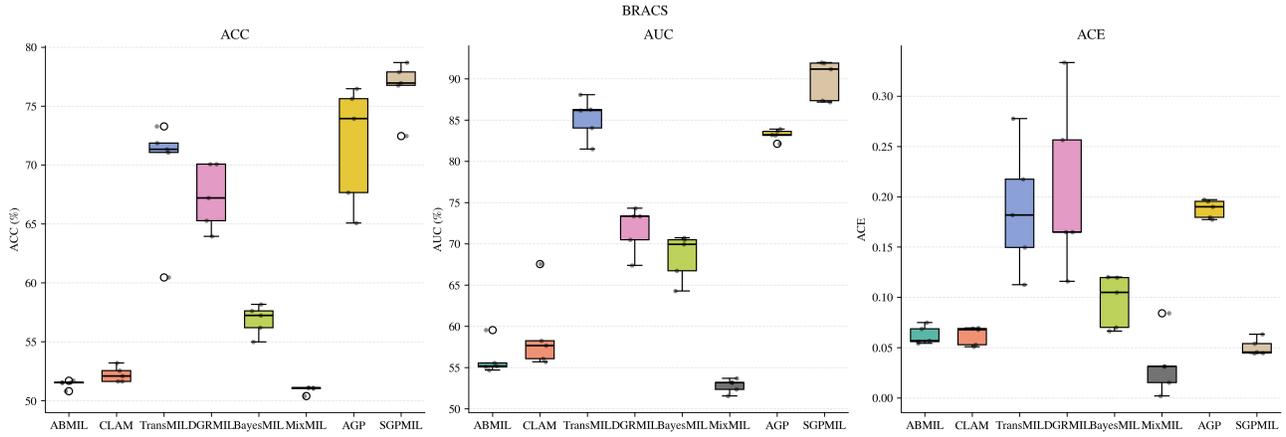


Figure 2. Instance-level performance across folds on BRACS. Boxplots show the distribution of accuracy (ACC), area under the ROC curve (AUC), and adaptive expected calibration error (ACE) for each model. \* denotes significance according to one-sided paired t-tests ( $p < .05$ ).

	ACC ( $p$ )	$\kappa$ ( $p$ )	ACE ( $p$ )
ABMIL	<b>1.86e-3</b>	<b>8.04e-4</b>	<b>1.01e-3</b>
CLAM	<b>7.51e-3</b>	<b>2.03e-3</b>	<b>4.44e-3</b>
TransMIL	<b>4.50e-3</b>	<b>1.53e-3</b>	<b>1.56e-3</b>
DGRMIL	<b>1.51e-2</b>	<b>3.50e-2</b>	<b>2.69e-2</b>
BayesMIL	<b>1.60e-3</b>	<b>2.92e-3</b>	<b>6.92e-3</b>
MixMIL	<b>2.23e-3</b>	<b>2.95e-3</b>	<b>8.19e-5</b>
AGP	<b>1.57e-2</b>	<b>1.71e-2</b>	2.31e-1

Table 6.  $p$ -values from one-sided paired t-tests comparing SGP-MIL with baseline models across folds on PANDA (slide level). Bold values indicate  $p < .05$ .

	ACC	AUC	ACE
ABMIL	.694 <sup>*,010</sup>	.852 <sup>*,025</sup>	.175 <sup>*,007</sup>
CLAM	.699 <sup>*,034</sup>	.850 <sup>*,021</sup>	.183 <sup>*,011</sup>
TransMIL	.676 <sup>*,005</sup>	.826 <sup>*,032</sup>	.186 <sup>*,012</sup>
DGRMIL	.703 <sup>*,033</sup>	.818 <sup>*,035</sup>	.186 <sup>*,023</sup>
BayesMIL	.648 <sup>*,058</sup>	.829 <sup>*,022</sup>	.183 <sup>*,028</sup>
MixMIL	.662 <sup>*,031</sup>	.855 <sup>*,010</sup>	.254 <sup>*,002</sup>
AGP	.634 <sup>*,030</sup>	.830 <sup>*,010</sup>	.134 <sup>*,014</sup>
<b>SGPMIL</b>	<b>.736<sup>*,029</sup></b>	<b>.870<sup>*,026</sup></b>	<b>.142<sup>*,032</sup></b>

Table 7. Slide-level performance on BRACS. \* indicates significance based on a one-sided paired t-test across folds ( $p < .05$ ).

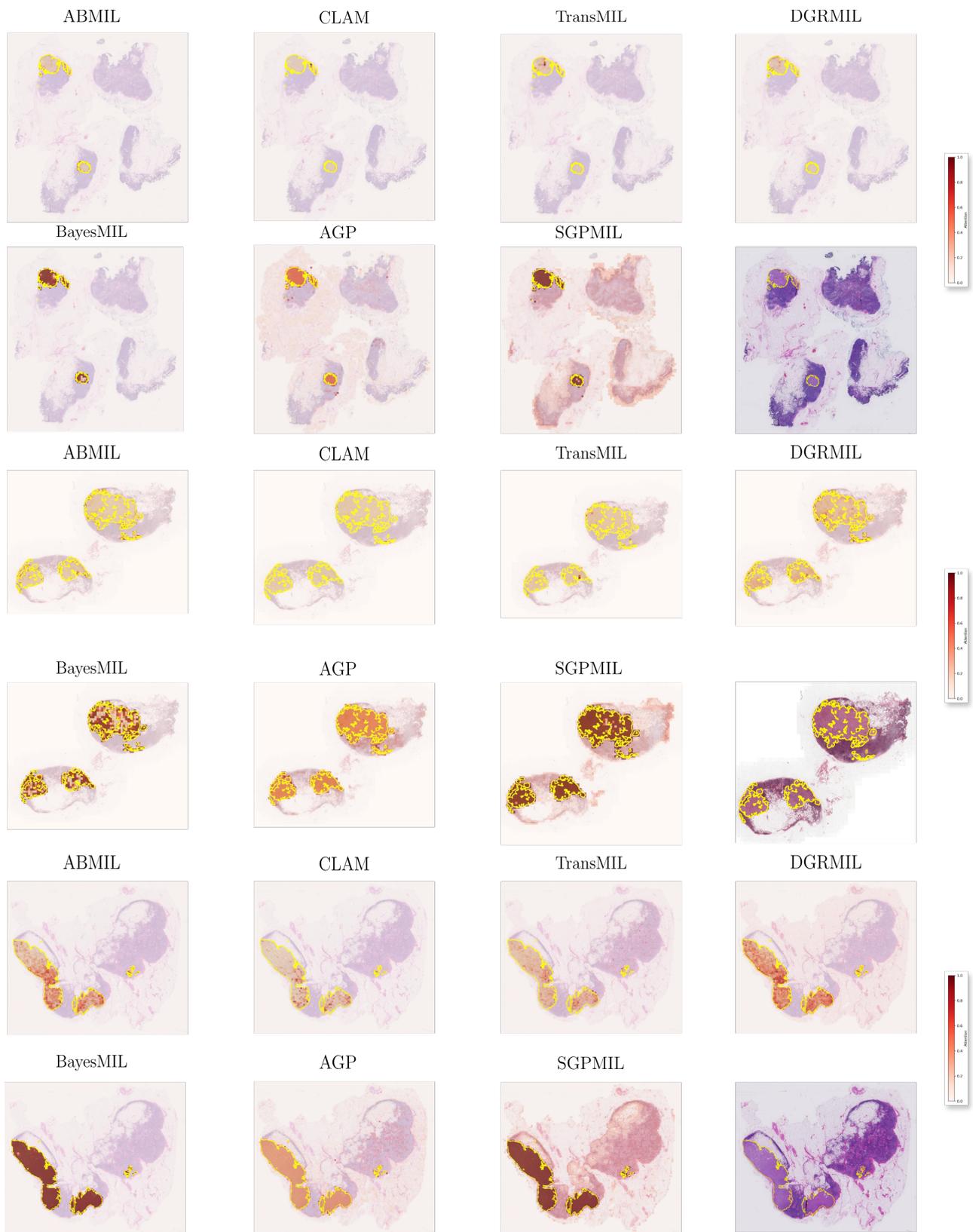


Figure 3. Additional attention heatmaps on CAMELYON16. Ground-truth tumor annotations are overlaid as yellow contours.

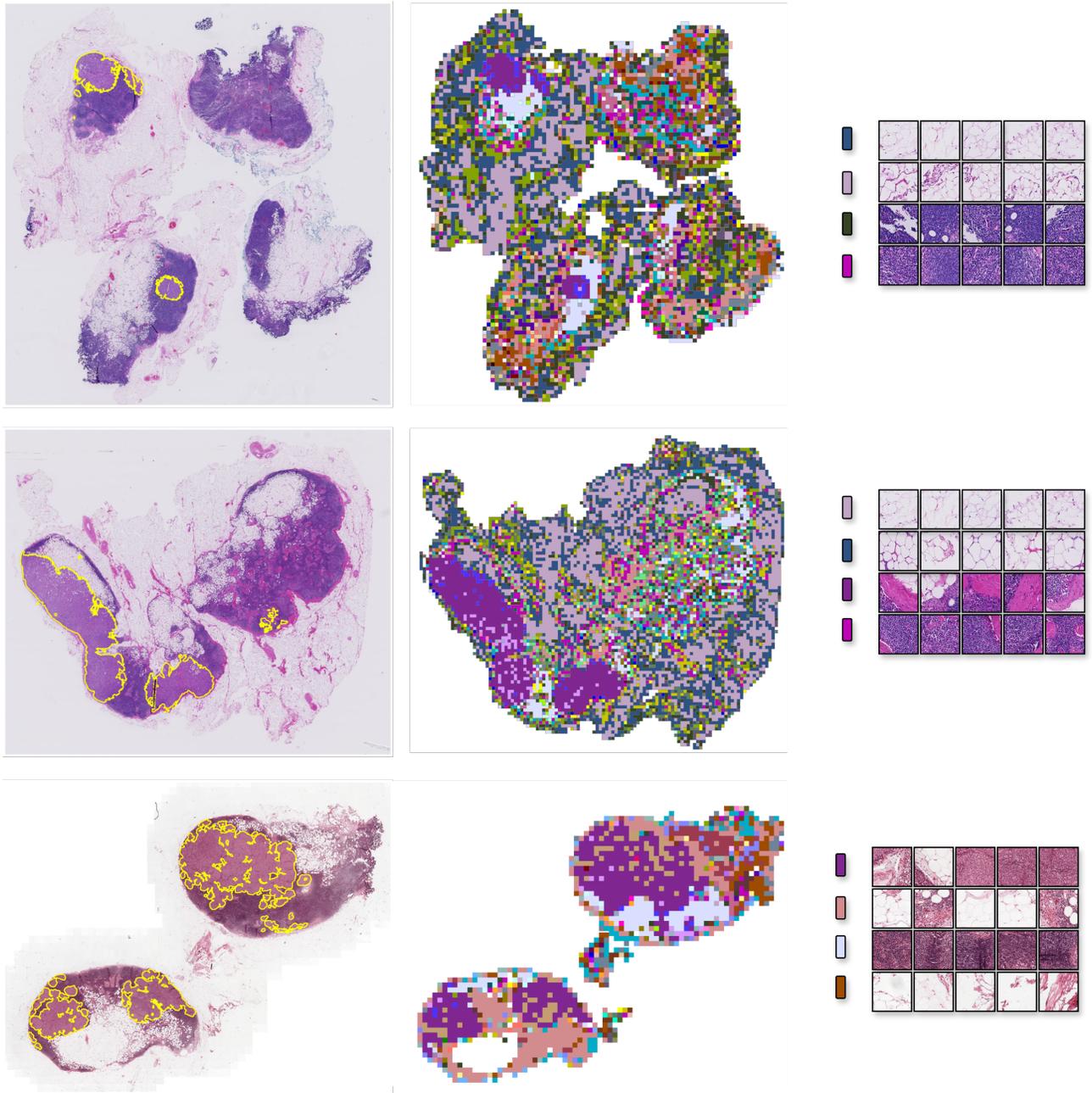


Figure 4. Additional inducing point label maps. Left: CAMELYON16 test set WSIs with ground-truth tumor annotations (yellow contours). Middle: label maps where each patch embedding is assigned to its most similar inducing point based on cosine similarity. Right: top-5 most similar patches for selected inducing points. These visualizations highlight how inducing points specialize in different tissue morphologies, including tumor, stroma, and interface regions.

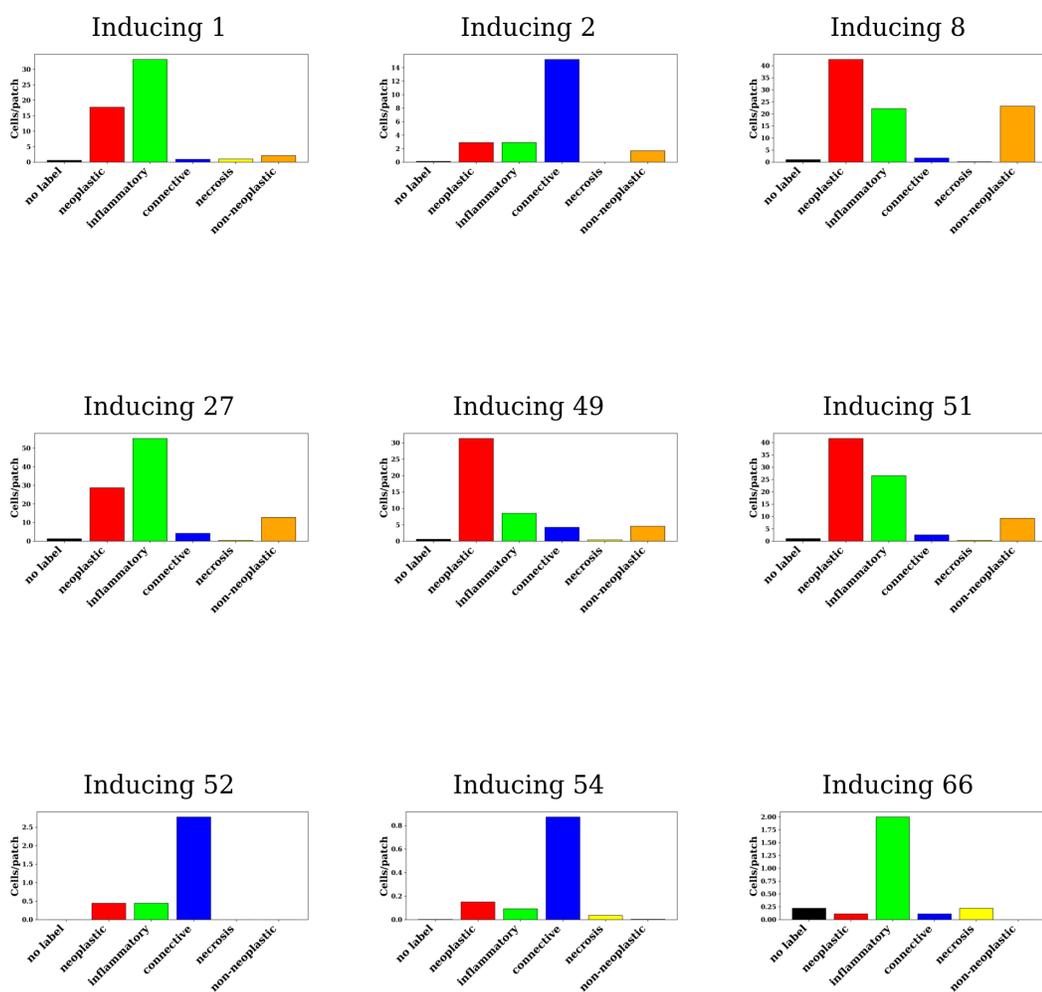
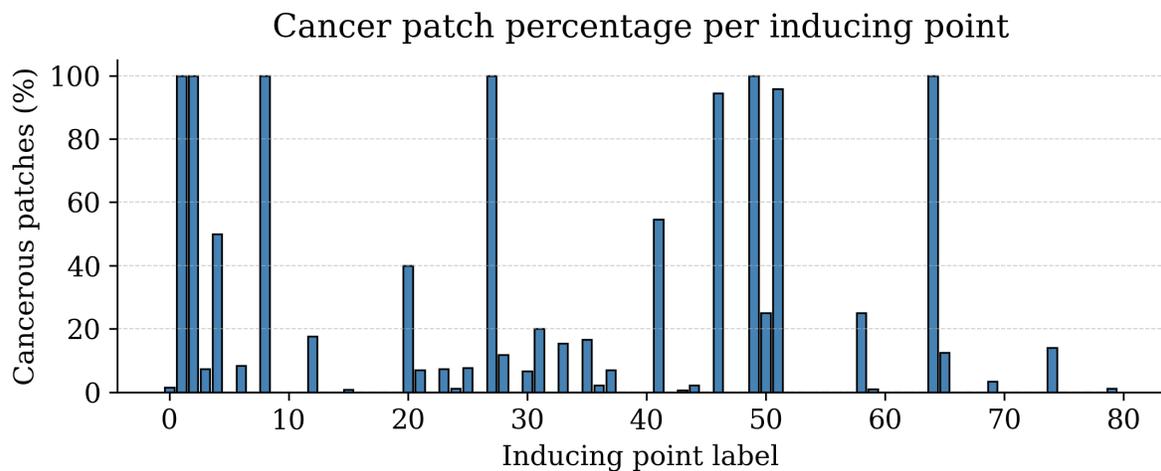


Figure 5. Cancer patch percentage per inducing point (top) and nucleus type distributions (bottom) for selected inducing points in CAMELYON16. Percentages were computed as the fraction of cancerous to total patches assigned to each inducing point in the slide. Inducing points 1, 2, 8, 27, 49, and 51 are highly enriched in cancerous regions and show a predominance of neoplastic nuclei, with inflammatory cells often the next most abundant type. By contrast, inducing points not associated with cancerous patches (52, 54, 66) mainly reflect connective or non-tumor tissue, with sparse isolated tumor cells, consistent with the lack of pixel-level annotations in CAMELYON16.

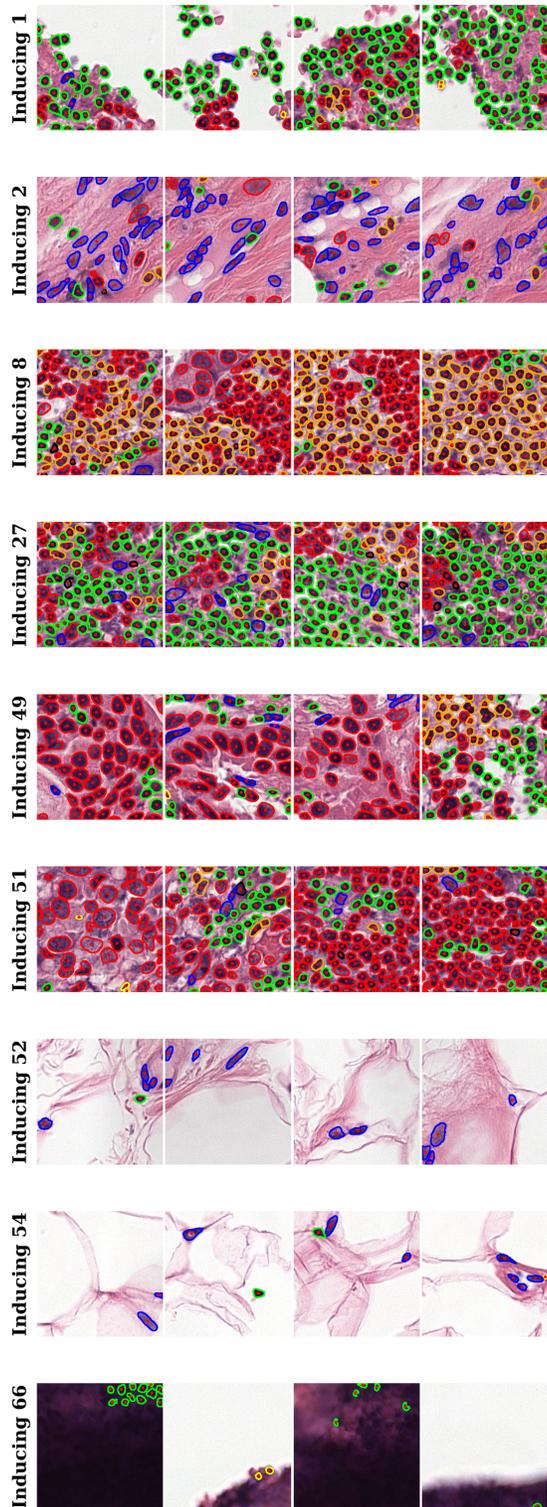


Figure 6. Representative image patches most similar to inducing points in CAMELYON16 slide 75. Each row corresponds to one inducing point (labels shown on the left), with four example patches illustrating the tissue context and associated nuclear types as segmented by HoverNet. Inducing points 1, 2, 8, 27, 49, and 51 predominantly align with tumor-rich regions, while others (e.g., 52, 54, 66) capture connective or non-tumor tissue, occasionally containing isolated tumor cells.

	ACC ( $p$ )	AUC ( $p$ )	ACE ( $p$ )
ABMIL	<b>1.81e-3</b>	<b>3.75e-2</b>	<b>1.53e-2</b>
CLAM	<b>3.11e-2</b>	<b>4.64e-2</b>	<b>7.92e-3</b>
TransMIL	<b>7.29e-4</b>	<b>1.30e-2</b>	<b>8.27e-3</b>
DGRMIL	<b>3.07e-2</b>	<b>8.13e-3</b>	<b>6.68e-3</b>
BayesMIL	<b>3.69e-3</b>	<b>8.01e-3</b>	<b>1.03e-2</b>
MixMIL	<b>8.88e-4</b>	<b>2.94e-2</b>	<b>1.49e-4</b>
AGP	<b>4.81e-4</b>	<b>3.85e-3</b>	1.80e-1

Table 8.  $p$ -values from one-sided paired t-tests comparing SGP-MIL with baseline models across folds on BRACS (slide level). Bold values indicate  $p < .05$ .

	ACC ( $p$ )	AUC ( $p$ )	ACE ( $p$ )
ABMIL	<b>1.27e-4</b>	<b>3.84e-3</b>	<b>6.41e-4</b>
CLAM	<b>8.92e-4</b>	<b>2.29e-3</b>	<b>1.60e-4</b>
TransMIL	<b>4.53e-5</b>	<b>1.14e-3</b>	<b>4.08e-2</b>
DGRMIL	<b>8.89e-3</b>	<b>4.82e-2</b>	<b>6.19e-4</b>
BayesMIL	<b>2.31e-4</b>	<b>1.12e-4</b>	<b>1.51e-5</b>
MixMIL	<b>3.03e-5</b>	<b>3.49e-3</b>	<b>3.96e-7</b>
AGP	<b>4.15e-3</b>	<b>6.15e-3</b>	9.75e-1

Table 9.  $p$ -values from one-sided paired t-tests comparing SGP-MIL with baseline models across folds on CAMELYON16 (instance level). Bold values indicate  $p < .05$ .

	ACC ( $p$ )	AUC ( $p$ )	ACE ( $p$ )
ABMIL	<b>1.03e-5</b>	<b>2.26e-5</b>	8.75e-2
CLAM	<b>1.50e-5</b>	<b>2.17e-4</b>	7.37e-2
TransMIL	<b>3.75e-3</b>	<b>6.65e-3</b>	<b>4.17e-3</b>
DGRMIL	<b>1.79e-4</b>	<b>1.20e-4</b>	<b>8.47e-3</b>
BayesMIL	<b>9.11e-5</b>	<b>3.22e-4</b>	<b>5.40e-3</b>
MixMIL	<b>1.01e-5</b>	<b>2.81e-6</b>	9.11e-1
AGP	<b>3.17e-2</b>	<b>3.17e-3</b>	<b>8.10e-7</b>

Table 10.  $p$ -values from one-sided paired t-tests comparing SGP-MIL with baseline models across folds on BRACS (instance level). Bold values indicate  $p < 0.05$ .

$U$	WSI			INSTANCE		
	ACC	AUC	ACE	ACC	AUC	ACE
16	.971	.977	.045	<u>.893</u>	.963	.271
32	<u>.979</u>	<u>.985</u>	<u>.026</u>	.880	<u>.970</u>	<b>.046</b>
64	.964	.981	.041	.832	.910	.210
80	<b>.980</b>	<b>.986</b>	<b>.021</b>	<b>.914</b>	<b>.973</b>	<u>.051</u>

Table 11. Performance of SGPMIL on CAMELYON16 for varying inducing points.

Model	Training (s)	Inference (s)	Params (M)
ABMIL	5.5*	0.8*	0.66
CLAM	7.0	0.9	0.92
TransMIL	13.4	1.2	2.67
DGRMIL	16.7	1.5	4.34
BayesMIL	9.5	1.1	1.32
MixMIL	11	1.5	1.57
AGP	31.0	6.0	1.21
<b>SGPMIL</b>	<b>9.0**</b>	<b>1.0**</b>	1.21

Table 12. Training and inference times (in seconds) and model sizes (number of trainable parameters in millions, M). Training times are averaged over 30 epochs, while inference times correspond to processing the full test set of 129 slides. For MixMIL, we report the full variational posterior variant (non-mean field approximation). \* marks the overall fastest model (ABMIL) and \*\* highlights the fastest probabilistic model among AGP, BayesMIL, MixMIL and SGPMIL.

Kernel	WSI		
	ACC	AUC	ACE
RBF	.980 <sub>.007</sub>	.986 <sub>.005</sub>	.021 <sub>.005</sub>
Linear	.973 <sub>.010</sub>	.985 <sub>.005</sub>	.048 <sub>.008</sub>
Matern	.985 <sub>.006</sub>	.988 <sub>.003</sub>	.024 <sub>.006</sub>

Table 13. Performance of SGPMIL on CAMELYON16 for various kernels, across folds.

Kernel	INSTANCE		
	ACC	AUC	ACE
RBF	.914 <sub>.006</sub>	.973 <sub>.004</sub>	.051 <sub>.008</sub>
Linear	.851 <sub>.017</sub>	.950 <sub>.020</sub>	.113 <sub>.016</sub>
Matern	.850 <sub>.011</sub>	.960 <sub>.012</sub>	.062 <sub>.014</sub>

Table 14. Performance of SGPMIL on CAMELYON16 for various kernels, across folds.

	Bag-level		Instance-level			
	AUC	ACE	F1	FROC	AUC	ACE
ABMIL	0.986	0.030	0.958	0.816	0.987	0.033
CLAM	0.989	0.021	0.952	0.813	<b>0.999</b>	0.039
TransMIL	0.990	0.021	<b>0.980</b>	0.826	0.993	0.025
DGRMIL	<b>0.996</b>	<u>0.016</u>	<u>0.979</u>	<u>0.827</u>	0.996	<b>0.013</b>
BayesMIL	0.990	0.019	0.968	0.824	0.994	<u>0.019</u>
AGP	0.994	0.046	0.957	0.820	0.981	0.083
<b>SGPMIL</b>	<u>0.995</u>	<b>0.009</b>	0.978	<b>0.830</b>	<u>0.998</u>	0.050

Table 15. Bag-level and instance-level performance for the MNIST-bags dataset. Bags are comprised of 9 instances each. All baselines are trained end-to-end with a CNN as initial feature extractor.