

# Supplementary information to the manuscript: A Benchmark for Virus Infection Reporter Virtual Staining in Fluorescence and Brightfield Microscopy

**Authors:** Maria Wyrzykowska\*, Gabriel della Maggiora\*, Nikita Deshpande, Ashkan Mokarian, Artur Yakimovich

\*these authors contributed equally to this work

*correspondence: a.yakimovich@hzdr.de*

## Standard Deviation Across Dataset Splits

To evaluate the impact of data splitting on model performance, we trained and tested the U-Net model for 100000 steps on datasets' splits generated using three different random seeds (43, 44, and 45) for each virus. For viruses represented as time-lapse sequences, we ensured that the same sample across different timepoints was included in only one of the data splits to prevent data leakage.

We then computed the mean and standard deviation of the performance metrics across all models and datasets to assess variability in model performance. The results are summarised in Table S1. The analysis indicates that while certain datasets, notably VACV and HAdV, exhibit higher variability across different data splits, particularly concerning the PSNR metric. Such fluctuations are generally limited. This consistency across most metrics and datasets reinforces the robustness and reliability of the study's conclusions.

**Table S1.** Performance of the U-Net model trained on HAdV, HAdV (2ch), VACV, HSV, IAV, and RV datasets. Reported metrics include the mean and standard deviation across models trained with three different data split seeds. Metrics shown are MSE, PSNR, and SSIM, along with their corresponding foreground (FG) and background (BG) versions.

Dataset	MSE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )	BG MSE ( $\downarrow$ )	BG PSNR ( $\uparrow$ )	BG SSIM ( $\uparrow$ )	FG MSE ( $\downarrow$ )	FG PSNR ( $\uparrow$ )	FG SSIM ( $\uparrow$ )
HAdV	$0.029 \pm 0.002$	$25.282 \pm 0.343$	$0.897 \pm 0.007$	$0.006 \pm 0.002$	$29.240 \pm 1.262$	$0.952 \pm 0.003$	$0.102 \pm 0.006$	$16.906 \pm 0.071$	$0.724 \pm 0.010$
HAdV (2ch)	$0.032 \pm 0.004$	$24.865 \pm 1.078$	$0.892 \pm 0.013$	$0.008 \pm 0.005$	$28.210 \pm 2.765$	$0.948 \pm 0.008$	$0.105 \pm 0.002$	$16.635 \pm 0.260$	$0.714 \pm 0.022$
VACV	$0.046 \pm 0.005$	$23.963 \pm 1.539$	$0.861 \pm 0.013$	$0.012 \pm 0.005$	$27.323 \pm 1.863$	$0.970 \pm 0.010$	$0.156 \pm 0.022$	$14.379 \pm 0.685$	$0.492 \pm 0.044$
HSV	$0.011 \pm 0.000$	$29.465 \pm 0.350$	$0.935 \pm 0.002$	$0.001 \pm 0.000$	$35.885 \pm 0.114$	$0.974 \pm 0.001$	$0.096 \pm 0.001$	$16.555 \pm 0.017$	$0.540 \pm 0.004$
IAV	$0.040 \pm 0.002$	$24.609 \pm 0.160$	$0.907 \pm 0.002$	$0.001 \pm 0.000$	$35.950 \pm 0.713$	$0.974 \pm 0.000$	$0.521 \pm 0.025$	$9.267 \pm 0.178$	$0.216 \pm 0.011$
RV	$0.025 \pm 0.001$	$24.681 \pm 0.163$	$0.913 \pm 0.002$	$0.003 \pm 0.001$	$31.946 \pm 0.989$	$0.958 \pm 0.002$	$0.419 \pm 0.014$	$9.862 \pm 0.157$	$0.205 \pm 0.004$

## Background-foreground metrics

To gain a more comprehensive understanding of the models' behavior, we employed additional sets of metrics. We calculated traditional metrics on the data segmented into "foreground" and "background" using a threshold-based approach. Specifically, regions of an image with values higher than a certain manually chosen threshold in the ground truth viral signal channel were classified as "foreground," while the rest were classified as "background." The threshold was set to -0.8 for all viruses except HAdV, for which we used -0.9. Table S2 contains the results.

The results suggest that predicting the background is easier, as shown by lower MSE and higher PSNR and SSIM scores across all models. VACV predictions had notably higher MSE, likely due to overestimating infection extent, which was confirmed by a manual review. Foreground metrics were worse overall, with HSV and HAdV performing best, while IAV and RV performed the worst, likely due to models hallucinating the signal. U-Net generally outperformed pix2pix, except for IAV, where pix2pix produced clearer, more accurate results despite U-Net's blurrier predictions.

## Cell classification metrics

Table S3 reports IoU, F1, accuracy, precision and recall metrics, but based on whole-cell masks, as opposed to previously used nuclei masks, which were available only for the HAdV and HAdV (2ch) datasets. The conclusions remain consistent: pix2pix generally outperforms U-Net, except in accuracy for both datasets and precision for the HAdV (2ch) data. Pix2pix achieves better IoU, F1, and recall scores, with precision being either comparable or slightly lower for HAdV. The accuracy, however, is lower when compared to metrics calculated from the nuclei masks. This can be attributed to the larger size of cells, which makes it easier to achieve higher IoU and recall scores. Because cells occupy a larger portion of the image, achieving high accuracy, which is heavily influenced by the background, becomes more challenging.

**Table S2.** Performance of U-Net and pix2pix models, trained on HAdV, HAdV (2ch), VACV, HSV, IAV and RV data. Metrics reported are MSE, PSNR and SSIM, calculated separately for background (BG) or foreground (FG). For pix2pix, we report the mean and standard deviation of each metric based on the 3 trials.

Dataset	Model	BG MSE ( $\downarrow$ )	BG PSNR ( $\uparrow$ )	BG SSIM ( $\uparrow$ )	FG MSE ( $\downarrow$ )	FG PSNR ( $\uparrow$ )	FG SSIM ( $\uparrow$ )
HAdV	pix2pix	$0.012 \pm 4.0\text{e-}05$	$26.272 \pm 3.8\text{e-}02$	$0.923 \pm 3.0\text{e-}05$	$0.092 \pm 2.0\text{e-}04$	$17.044 \pm 2.0\text{e-}02$	$0.686 \pm 1.0\text{e-}04$
	U-Net	<b>0.007</b>	<b>28.435</b>	<b>0.950</b>	<b>0.087</b>	<b>17.500</b>	<b>0.738</b>
HAdV (2ch)	pix2pix	$0.010 \pm 3.0\text{e-}05$	<b><math>27.597 \pm 3.3\text{e-}02</math></b>	$0.944 \pm 6.0\text{e-}05$	$0.097 \pm 8.0\text{e-}04$	$16.882 \pm 1.7\text{e-}02$	$0.716 \pm 1.0\text{e-}03$
	U-Net	<b>0.009</b>	27.263	<b>0.947</b>	<b>0.094</b>	<b>16.989</b>	<b>0.726</b>
VACV	pix2pix	$0.057 \pm 2.0\text{e-}05$	$21.144 \pm 5.0\text{e-}04$	$0.919 \pm 5.0\text{e-}05$	$0.201 \pm 6.0\text{e-}04$	$13.304 \pm 1.4\text{e-}02$	$0.426 \pm 1.6\text{e-}03$
	U-Net	<b>0.031</b>	<b>26.673</b>	<b>0.949</b>	<b>0.163</b>	<b>14.054</b>	<b>0.520</b>
HSV	pix2pix	$0.003 \pm 2.0\text{e-}06$	$32.252 \pm 7.0\text{e-}03$	$0.961 \pm 1.0\text{e-}05$	$0.143 \pm 8.0\text{e-}06$	$14.607 \pm 5.0\text{e-}04$	$0.438 \pm 6.0\text{e-}05$
	U-Net	<b>0.002</b>	<b>35.137</b>	<b>0.974</b>	<b>0.089</b>	<b>16.823</b>	<b>0.541</b>
IAV	pix2pix	$0.019 \pm 2.0\text{e-}05$	$23.403 \pm 6.0\text{e-}03$	$0.942 \pm 3.0\text{e-}05$	<b><math>0.498 \pm 2.0\text{e-}04</math></b>	<b><math>9.429 \pm 2.0\text{e-}03</math></b>	$0.174 \pm 8.0\text{e-}06$
	U-Net	<b>0.001</b>	<b>35.856</b>	<b>0.974</b>	0.544	9.075	<b>0.208</b>
RV	pix2pix	$0.022 \pm 9.0\text{e-}06$	$22.622 \pm 2.0\text{e-}03$	$0.903 \pm 2.0\text{e-}05$	$0.465 \pm 1.0\text{e-}04$	$9.377 \pm 9.6\text{e-}04$	$0.160 \pm 1.4\text{e-}04$
	U-Net	<b>0.002</b>	<b>33.174</b>	<b>0.960</b>	<b>0.438</b>	<b>9.691</b>	<b>0.205</b>

**Table S3.** Performance of U-Net and pix2pix models, trained on HAdV and HAdV (2ch) data. Metrics reported are IoU, F1, accuracy, precision, and recall, calculated based on cells segmentation masks. For pix2pix, we report the mean and standard deviation of each metric based on 3 trials.

Dataset	Model	IoU ( $\uparrow$ )	F1 ( $\uparrow$ )	Accuracy ( $\uparrow$ )	Precision ( $\uparrow$ )	Recall ( $\uparrow$ )
HAdV	pix2pix	<b><math>0.326 \pm 6.0\text{e-}04</math></b>	<b><math>0.462 \pm 6.0\text{e-}04</math></b>	$0.856 \pm 2.0\text{e-}04$	$0.500 \pm 7.0\text{e-}04$	<b><math>0.548 \pm 6.0\text{e-}04</math></b>
	U-Net	0.303	0.440	<b>0.871</b>	<b>0.528</b>	0.447
HAdV (2ch)	pix2pix	<b><math>0.294 \pm 2.0\text{e-}03</math></b>	<b><math>0.428 \pm 2.0\text{e-}03</math></b>	$0.868 \pm 4.0\text{e-}04$	<b><math>0.529 \pm 3.0\text{e-}03</math></b>	<b><math>0.438 \pm 2.0\text{e-}03</math></b>
	U-Net	0.282	0.413	<b>0.871</b>	0.516	0.408

## Standard Deviation Across Samples

We calculated the standard deviations of performance metrics across all models and datasets to assess the variation in model performance. The results are detailed in Tables S4 and S5.

The analysis reveals that VACV exhibits the highest variance in both general and background metrics. This suggests that the models made substantial errors on a few samples, likely exacerbated by the smaller size of the VACV test set. In contrast, the standard deviations for background metrics across other viruses were relatively low. However, foreground metrics showed higher variability, with IAV demonstrating the greatest standard deviation in MSE.

**Table S4.** Performance of U-Net and pix2pix models, trained on HAdV, HAdV (2ch), VACV, HSV, IAV, and RV data. Metrics reported are the mean and standard deviation between results for each datapoint for MSE, PSNR, and SSIM, along with their foreground (FG) and background (BG) versions.

Dataset	Model	MSE ( $\downarrow$ )	PSNR ( $\uparrow$ )	SSIM ( $\uparrow$ )	BG MSE ( $\downarrow$ )	BG PSNR ( $\uparrow$ )	BG SSIM ( $\uparrow$ )	FG MSE ( $\downarrow$ )	FG PSNR ( $\uparrow$ )	FG SSIM ( $\uparrow$ )
HAdV	pix2pix	$0.034 \pm 0.032$	$24.214 \pm 4.218$	$0.870 \pm 0.060$	$0.012 \pm 0.009$	$26.272 \pm 3.267$	$0.923 \pm 0.016$	$0.092 \pm 0.055$	$17.044 \pm 2.433$	$0.686 \pm 0.084$
	U-Net	<b><math>0.031 \pm 0.038</math></b>	<b><math>25.331 \pm 4.753</math></b>	<b><math>0.899 \pm 0.067</math></b>	<b><math>0.007 \pm 0.006</math></b>	<b><math>28.435 \pm 3.112</math></b>	<b><math>0.950 \pm 0.011</math></b>	<b><math>0.087 \pm 0.060</math></b>	<b><math>17.500 \pm 2.796</math></b>	<b><math>0.738 \pm 0.076</math></b>
HAdV (2ch)	pix2pix	<b><math>0.032 \pm 0.032</math></b>	<b><math>24.824 \pm 4.776</math></b>	$0.892 \pm 0.065$	$0.010 \pm 0.007$	<b><math>27.597 \pm 3.954</math></b>	$0.944 \pm 0.012$	$0.097 \pm 0.057$	$16.882 \pm 2.698$	$0.716 \pm 0.083$
	U-Net	<b><math>0.032 \pm 0.034</math></b>	$24.767 \pm 4.485$	<b><math>0.895 \pm 0.067</math></b>	<b><math>0.009 \pm 0.006</math></b>	$27.263 \pm 3.056$	<b><math>0.947 \pm 0.009</math></b>	<b><math>0.094 \pm 0.057</math></b>	<b><math>16.989 \pm 2.542</math></b>	<b><math>0.726 \pm 0.078</math></b>
VACV	pix2pix	$0.083 \pm 0.070$	$21.087 \pm 3.593$	$0.825 \pm 0.096$	$0.057 \pm 0.069$	$21.144 \pm 4.480$	$0.919 \pm 0.073$	$0.201 \pm 0.078$	$13.304 \pm 1.621$	$0.426 \pm 0.049$
	U-Net	<b><math>0.057 \pm 0.047</math></b>	<b><math>24.563 \pm 7.304</math></b>	<b><math>0.860 \pm 0.090</math></b>	<b><math>0.031 \pm 0.048</math></b>	<b><math>26.673 \pm 6.911</math></b>	<b><math>0.949 \pm 0.059</math></b>	<b><math>0.163 \pm 0.044</math></b>	<b><math>14.054 \pm 1.195</math></b>	<b><math>0.520 \pm 0.068</math></b>
HSV	pix2pix	$0.016 \pm 0.007$	$27.514 \pm 2.647$	$0.917 \pm 0.030$	$0.003 \pm 0.002$	$32.252 \pm 2.108$	$0.961 \pm 0.015$	$0.143 \pm 0.037$	$14.607 \pm 1.055$	$0.438 \pm 0.062$
	U-Net	<b><math>0.010 \pm 0.005</math></b>	<b><math>29.673 \pm 2.714</math></b>	<b><math>0.938 \pm 0.023</math></b>	<b><math>0.002 \pm 0.002</math></b>	<b><math>35.137 \pm 2.704</math></b>	<b><math>0.974 \pm 0.010</math></b>	<b><math>0.089 \pm 0.034</math></b>	<b><math>16.823 \pm 1.503</math></b>	<b><math>0.541 \pm 0.073</math></b>
IAV	pix2pix	$0.054 \pm 0.014$	$22.266 \pm 1.103$	$0.875 \pm 0.023$	$0.019 \pm 0.004$	$23.403 \pm 0.855$	$0.942 \pm 0.007$	<b><math>0.498 \pm 0.211</math></b>	<b><math>9.429 \pm 1.818</math></b>	$0.174 \pm 0.088$
	U-Net	<b><math>0.041 \pm 0.016</math></b>	<b><math>24.374 \pm 1.861</math></b>	<b><math>0.905 \pm 0.023</math></b>	<b><math>0.001 \pm 0.000</math></b>	<b><math>35.856 \pm 1.587</math></b>	<b><math>0.974 \pm 0.003</math></b>	$0.544 \pm 0.235$	$9.075 \pm 1.887$	<b><math>0.208 \pm 0.109</math></b>
RV	pix2pix	$0.045 \pm 0.009$	$21.668 \pm 0.927$	$0.863 \pm 0.019$	$0.022 \pm 0.004$	$22.622 \pm 0.818$	$0.903 \pm 0.011$	$0.465 \pm 0.056$	$9.377 \pm 0.546$	$0.160 \pm 0.028$
	U-Net	<b><math>0.024 \pm 0.007</math></b>	<b><math>24.738 \pm 1.400</math></b>	<b><math>0.917 \pm 0.018</math></b>	<b><math>0.002 \pm 0.001</math></b>	<b><math>33.174 \pm 1.650</math></b>	<b><math>0.960 \pm 0.007</math></b>	<b><math>0.438 \pm 0.088</math></b>	<b><math>9.691 \pm 0.888</math></b>	<b><math>0.205 \pm 0.053</math></b>

**Table S5.** Performance of U-Net and pix2pix models, trained on HAdV, HAdV (2ch), HSV, IAV, and RV data. Metrics reported are the mean and standard deviation between results for each datapoint for IoU, F1, accuracy, precision, and recall calculated based on the nuclei masks.

Dataset	Model	IoU ( $\uparrow$ )	F1 ( $\uparrow$ )	Accuracy ( $\uparrow$ )	Precision ( $\uparrow$ )	Recall ( $\uparrow$ )
HAdV	pix2pix	<b><math>0.300 \pm 0.163</math></b>	<b><math>0.437 \pm 0.199</math></b>	$0.959 \pm 0.022$	$0.503 \pm 0.302$	<b><math>0.473 \pm 0.141</math></b>
	U-Net	$0.268 \pm 0.137$	$0.404 \pm 0.174$	$0.963 \pm 0.025$	<b><math>0.557 \pm 0.309</math></b>	$0.357 \pm 0.124$
HAdV (2ch)	pix2pix	<b><math>0.239 \pm 0.138</math></b>	<b><math>0.364 \pm 0.182</math></b>	$0.960 \pm 0.025$	<b><math>0.524 \pm 0.320</math></b>	<b><math>0.322 \pm 0.141</math></b>
	U-Net	$0.235 \pm 0.140$	$0.360 \pm 0.183$	<b><math>0.961 \pm 0.026</math></b>	$0.522 \pm 0.321$	$0.306 \pm 0.129$
HSV	pix2pix	$0.660 \pm 0.054$	$0.794 \pm 0.040$	$0.985 \pm 0.005$	<b><math>0.862 \pm 0.052</math></b>	$0.741 \pm 0.063$
	U-Net	<b><math>0.727 \pm 0.048</math></b>	<b><math>0.841 \pm 0.034</math></b>	<b><math>0.988 \pm 0.005</math></b>	$0.858 \pm 0.050$	<b><math>0.830 \pm 0.056</math></b>
IAV	pix2pix	<b><math>0.249 \pm 0.082</math></b>	<b><math>0.391 \pm 0.107</math></b>	$0.967 \pm 0.008$	$0.485 \pm 0.154$	<b><math>0.341 \pm 0.095</math></b>
	U-Net	$0.174 \pm 0.067$	$0.290 \pm 0.098$	<b><math>0.971 \pm 0.010</math></b>	<b><math>0.655 \pm 0.186</math></b>	$0.190 \pm 0.072$
RV	pix2pix	$0.219 \pm 0.047$	$0.357 \pm 0.062$	$0.963 \pm 0.015$	$0.347 \pm 0.088$	<b><math>0.379 \pm 0.056</math></b>
	U-Net	<b><math>0.294 \pm 0.085</math></b>	<b><math>0.448 \pm 0.104</math></b>	<b><math>0.978 \pm 0.008</math></b>	<b><math>0.614 \pm 0.083</math></b>	$0.358 \pm 0.102$