# U-NetPlus: A Modified Encoder-Decoder U-Net Architecture for Semantic and Instance Segmentation of Surgical Instrument

S. M. Kamrul Hasan and Cristian A. Linte, *Senior Member, IEEE*

*Abstract*—Conventional therapy approaches limit surgeons' dexterity control due to limited field-of-view. With the advent of robot-assisted surgery, there has been a paradigm shift in medical technology for minimally invasive surgery. However, it is very challenging to track the position of the surgical instruments in a surgical scene, and accurate detection & identification of surgical tools is paramount. Deep learning-based semantic segmentation in frames of surgery videos has the potential to facilitate this task. In this work, we modify the U-Net architecture named U-NetPlus, by introducing a pre-trained encoder and re-design the decoder part, by replacing the transposed convolution operation with an upsampling operation based on nearest-neighbor (NN) interpolation. To further improve performance, we also employ a very fast and flexible data augmentation technique. We trained the framework on $8 \times 225$ frame sequences of robotic surgical videos, available through the MICCAI 2017 EndoVis Challenge dataset and tested it on $8 \times 75$ frame and $2 \times 300$ frame videos. Using our U-NetPlus architecture, we report a 90.20% DICE for binary segmentation, 76.26% DICE for instrument part segmentation, and 46.07% for instrument type (i.e., all instruments) segmentation, outperforming the results of previous techniques implemented and tested on these data.

*Index Terms*—Minimally invasive laparoscopic surgery, surgical instrument detection and identification, deep learning-based segmentation, U-Net framework, pre-trained encoder, nearest-neighbor interpolation.

## I. INTRODUCTION

Minimally invasive surgery has addressed many of the challenges of traditional surgical approaches by significantly reducing risk of infections and shortening hospitalization times, achieving similar outcome to traditional open surgery. There is a new paradigm shift in this field thanks to robot assistance under laparoscopic visualization [1]. To facilitate the manipulation of the laparoscopic surgical instruments while visualizing the endoscopic scene, surgical instrument identification is critical. Nevertheless, this task is challenging, due to the surrounding effects like illumination changes, visual occlusions, and presence of non-class objects. Hence, it is important to devise segmentation techniques that are sufficiently accurate and robust to ensure accurate tracking of the surgical tools to facilitate therapy via accurate manipulation of the laparoscopic instruments.

Although in recent years semantic segmentation methods applied to city-scapes, street scenes, and even Landsat image datasets [2], [3] have achieved ground-breaking performance by the virtue of deep convolutional neural networks (CNNs), image segmentation in clinical settings still requires more accuracy and precision, with even minimal segmentation errors being unacceptable. In the context of deep learning, Long *et al.* [4] proposed the first fully convolutional network (FCN) for semantic image segmentation, exploiting the capability of Convolutional Neural Networks (CNNs). However, their adoption for use in the medical domain was initially challenging, due to the limited availability of medical imaging data. These challenges were later circumvented by patch-based training, data augmentation, and transfer learning techniques [5]. These newer deep architectures learn to decode low-resolution images produced by VGG16 network into pixel-wise predictions. This network has 13 convolutional layers, with 3 fully connected layers and their weights are typically pre-trained on the large ImageNet object classification dataset [6], whereas, the decoder network upsamples feature maps from the bottleneck layer.

However, semantic segmentation is not sufficiently accurate for handling multi-class objects, due to the close presence of objects of the same class in the surgical scene. Therefore, the proposed work is motivated by the need to improve multi-class object segmentation, by leveraging the power of the existing U-Net architecture and augmenting it with new capabilities.

With the advent of U-Net architectures, a wide range of medical imaging tasks have been implemented and produced state-of-the-art results since 2015 [7]. Recently, Chen *et al.* modified U-Net architecture by introducing sub-pixel layers to improve low-light imaging [8]. Following the end-to-end training of a fully-convolutional network, they obtained promising results, with high signal-to-noise-ratio (SNR) and perfect color transformation on their own SID dataset. The authors in [9], [10] used nearest-neighbor interpolation for image reconstruction and super-resolution. The authors in [11] investigated the problem of transposed convolution and provided a solution by nearest-neighbor interpolation. However, the importance of integrating it into the deep CNN as part of the image upsampling operation was not fully explored so far. There have been a few papers tackling the segmentation and identification of surgical instruments from endoscopic video image, and, even fewer than half a dozen papers tackling this challenge using deep learning. One notable research contribution has been the use of a modified version of FCN-8, yet with no attempts for multi-class segmentation [12].

Multi-class (both instrument part and type) tool segmentation was first proposed by Shvets *et al.* [13], and Pakhomov *et al.* [14] and achieved promising results. They modified the classic U-Net model [7] that relies on the transposed convolution or deconvolution, in a similar, yet opposite fashion to the convolutional layers. As an example, instead of mapping from $4 \times 4$ input pixels to 1 output pixel, they map from 1 input pixel to $4 \times 4$ output pixels. However, its performance is much slower as the filters need additional weights and parameters that also require training. These large number of trainable parameters make the network hard to train end-to-end on a relevant task. Additionally, transposed convolution can easily lead to "uneven overlap", characterized by checkerboard-like patterns resulting in artifacts on a variety of scales and colors [11]. Redford *et al.* [15] and Salimans *et al.* [16] mentioned the problem of those artifacts and checkerboard patterns generated by the transposed convolution. While it is difficult to entirely remove these limitations and their resulting artifacts, our goal is to, at first, minimize their occurrence.

Hence, in the efforts to mitigate these challenges associated with the classic U-Net architecture, in this paper, we present the U-NetPlus model by introducing both VGG-11 and VGG-16 as an encoder with batch-normalized pre-trained weights and nearest-neighbor interpolation as the replacement of the transposed convolution in the decoder layer. This pre-trained encoder [17] speeds up convergence and leads to improved results by circumventing the optimization challenges associated with the target data [18]. Moreover, the nearest-neighbor interpolation used in the decoder section removes the artifacts generated by the transposed convolution.

To test the proposed U-NetPlus network, we implemented some of the recent state-of-the-art architectures for surgical tool segmentation and compared their results to those of the U-NetPlus architecture. From the above mentioned papers, only one seems to have achieved results comparable to ours [17], but it still suffers from several artifacts, which we have been able to further mitigate some of these artifacts using our proposed method. As such, while this paper leverages some of the existing infrastructure of fully convolutional network, it focuses on demonstrating the adaptation of existing infrastructure to refine its performance for a given task — in this case the segmentation and identification of surgical instruments from endoscopic images — rather than proposing a new fully convolutional framework. We believe there exist sufficient tools for segmentation and identification, however the integration and adaptation of these tools for improved performance is key to further improving the outcome of such tools. We demonstrate that the potential use of nearest-neighbor interpolation in the decoder removes artifacts and reduces the number of parameters.

## II. METHODOLOGY

### A. Overview of Proposed Method

U-NetPlus has an encoder network and a corresponding decoder network, followed by a final pixel-wise segmentation
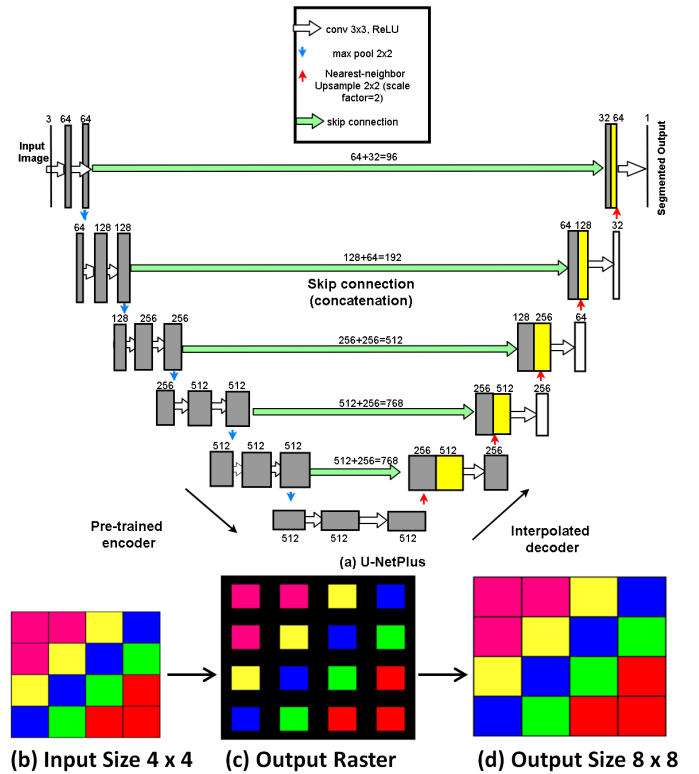


Fig. 1. (a) Modified U-Net with batch-normalized VGG11 as an encoder. Each box corresponds to a multi-channel featuring a map passing through a series of transformations. It consists of both an upsampling and a downsampling path and the height of the box represents the feature map resolution, while the width represents number of channels. Cyan arrows represent the max-pooling operation, whereas light-green arrows represent skip connections that transfer information from the encoder to the decoder. Red upward arrows represent the decoder which consists of nearest-neighbor upsampling with a scale factor of 2 followed by 2 convolution layers and a ReLU activation function; (b)-(d) working principle of nearest-neighbor interpolation where the low-resolution image is resized back to the original image.

layer, as illustrated in the architecture in Fig. 1. Similar to U-Net, our proposed U-NetPlus works like an auto-encoder with both a downsampling and an upsampling path. To maintain the number of channels the same as in the encoder part, skip connections are added between blocks of the same size in the downsampling and upsampling paths. This allows very precise alignment of the mask to the original image, which is particularly important in medical imaging. Furthermore, skip connections are known to ease network optimization by introducing multiple paths for backpropagation of the gradients, hence, mitigating the vanishing gradient problem.

Generally, weights are initialized randomly to train a network. However, limited training data can introduce overfitting problems, which become very "expensive" as far as manually altering the segmentation mask. Therefore, transfer learning can be used to initialize the network weights. But since a surgical instrument is not a class of ImageNet, one way to use transfer learning for a new task is to partially reuse ImageNet feature extractor — VGG11 or VGG16 as encoder — and then add a decoder. An improvement has been introduced for the encoder part, where we initiated a pre-trained VGG-11

and VGG-16 architecture with batch-normalization layers that has 11 and 16 sequential layers, respectively. Following this modification, it has been shown the pre-trained model is able to train the network within a very short time and with greater accuracy [19].

The feature map of VGG11 consists of seven convolutional layers of $3 \times 3$ kernel size followed by a ReLU activation function. For the reduction of the feature map size, max polling with stride 1 was used. The network starts by producing 64 channels in its first layer, and then continues by doubling the number of channels after each pooling operation until 512. Weights are copied from the original pre-trained VGG-11 on Imagenet.

The key effect of batch normalization has been investigated in a recent paper [20]. According to this work, batch normalization not only reduces the internal co-variate shift, but also re-parameterizes the underlying gradient optimization problem that makes the training more predictive at a faster convergence. After analyzing the impact of inserting BatchNorm layer, we applied BatchNorm layer after each convolutional layer.

The downsampling path decreases the feature size while increasing the number of feature maps, whereas the upsampling path increases the feature size while decreasing the number of feature maps, eventually leading to a pixel-wise mask. For the upsampling operation, we modified the existing architecture to reconstruct the high-resolution feature maps. Rather than using transposed convolution, we used the nearest-neighbor upsampling layer with a carefully selected stride and kernel size at the beginning of each block followed by two convolution layers and a ReLU function that would increase the spatial dimension in each block by a factor of 2.

Nearest-neighbor interpolation upsamples the input feature map by superimposing a regular grid onto it. Given $I_i$ be the input grid which is to be sampled, the output grid is produced by a linear transformation $\tau_\theta(I_i)$. Therefore, for an upsampling operation, $\tau_\theta$ can be defined as:

$$\begin{pmatrix} p_i^o \\ q_i^o \end{pmatrix} = \tau_\theta(I_i) = \begin{bmatrix} \theta & 0 \\ 0 & \theta \end{bmatrix} \begin{pmatrix} p_i^t \\ q_i^t \end{pmatrix}, \theta \geq 1, \qquad (1)$$

where $(p_i^o, q_i^o) \, \epsilon \, I_i$ are the original sampling input coordinates, $(p_i^t, q_i^t)$ are the target coordinates, and $\theta$ upsampling factor. The principle of how nearest-neighbor (NN) interpolation works to enlarge the image size, is shown in the Fig. 1. After locating the center pixel of the cell of the output raster dataset on the input raster, the location of the nearest center of the cell on the input raster will be determined and the value of that cell on the output raster will be assigned afterwards. As an example, we demonstrate the upsampling of a $4 \times 4$ image using this approach. The cell centers of the output raster are equidistant. A value is needed to be derived from the input raster for each output cell. Nearest-neighbor interpolation would select those cells centers from the input raster that are closest to that of output raster. The black areas of the middle image can be filled with the copies of the center pixel. Therefore, this fixed interpolation weights requires no learning for upsampling operation compared to strided or
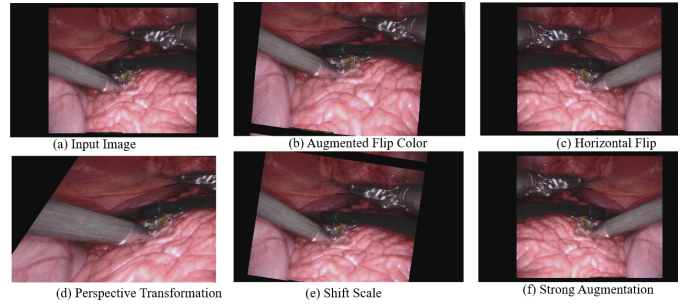


(a) Input Image    (b) Augmented Flip Color    (c) Horizontal Flip
(d) Perspective Transformation    (e) Shift Scale    (f) Strong Augmentation

Fig. 2. Example images of applying both affine and elastic transformation in albumentations library for data augmentation.

transposed convolution leading to create a more memory efficient upsampling operation. The algorithm is similar to the one proposed and used by the authors of [21] in their work.

### B. Dataset

For both training and validation, we used the Robotic instruments dataset from the sub-challenge of **MICCAI 2017 Endoscopic Vision Challenge** [22]. The training dataset has $8 \times 225$ frame sequences with 2 Hz frame rate of high resolution stereo camera images collected from a da Vinci Xi surgical system during laparoscopic cholecystectomy procedures. The frames were re-sampled from 30 Hz video to 2 Hz to avoid any redundancy issues. A stereo camera was used to capture the video sequences that consists of the left and right eye views with resolution of $1920 \times 1080$ in RGB format. In each frame, the articulated parts of the robotic surgical instrument consisting of a rigid shaft, an articulated wrist, and claspers, have been manually labeled by expert clinicians. The test set has $8 \times 75$ frame sequences and $2 \times 300$ frame videos. The challenge is to segment 7 classes such as prograsp forceps, needle driver, vessel sealer, grasping retractor etc.

### C. Data Augmentation

We augmented the MICCAI 2017 EndoVis Challenge data using the albumentations library that was reported as a fast and flexible implementation for data augmentation in [23]. These libraries include both affine and elastic transformations, and their effects on the image data during augmentation is illustrated in Fig. 2.

In short, the affine transformation includes scaling, translation, horizontal flip, vertical flip, random brightness and noise addition etc. For the elastic transformation (non-affine), first a random displacement field, $F(R)$ is generated for the horizontal and vertical directions, $\delta x$, and $\delta y$ respectively where, $[\delta x, \delta y] = [-1 \leq \delta x, \delta y \leq +1]$.

These random fields are then convolved with an intermediate value of $\sigma$ (in pixels) and the fields are multiplied by a scaling factor $\alpha$ that controls intensity. Thus, we obtain the elastically transformed image in which the global shape of the interest is undisturbed, unlike in the affine-transformed image.

### D. Implementation Details

We implemented our methodology using PyTorch[1]. During the pre-processing step, we cropped the un-wanted black border from each video frame. Images were normalized by subtracting their mean and dividing by their standard deviation (i.e., according to their z-scores). Batch normalization was used before each weighted layer, as it re-parameterizes the underlying gradient optimization problem that helps the training to converge faster [20]. For training, we use the Adam optimizer with a learning rate of 0.00001. We didn't use dropout as it degraded validation performance in our case. All models were trained for 100 epochs. The training set was shuffled before each epoch and the batch size was 4 in our case. All experiments were run on a machine equipped with a NVIDIA GTX 1080 Ti GPU (11GBs of memory).

### E. Performance Metrics

In this work, we used the common Jaccard index — also referred to as the intersection-over-union (IoU) — to evaluate segmentation results. It is an overlap index that quantifies the agreement between two segmented image regions: a ground truth segmentation and the predicted segmentation method. Given a vector of ground truth labels $T_1$ and a vector of predicted labels $P_1$, IoU can be defined as

$$J(T_1, P_1) = \frac{|T_1 \cap P_1|}{|T_1 \cup P_1|} = \frac{|T_1 \cap P_1|}{|T_1| + |P_1| - |T_1 \cap P_1|}, \quad (2)$$

Eqn. 2 can further be clarified. Given a pixel $j$, the label of the pixel $z_j$, and the probability of the same pixel for the predicted class $\hat{z}_j$, Eqn. 1 for $k$ number of dataset

$$J = \frac{1}{k} \sum_{j=1}^{k} \left( \frac{z_j \hat{z}_j}{z_j + \hat{z}_j - z_j \hat{z}_j} \right), \quad (3)$$

We can represent the loss function in a common ground of $log$ scale as this task is a pixel classification problem. So, for a given pixel $j$, the common loss can be defined as the function $H$ for $k$ number of dataset

$$H = -\frac{1}{k} \sum_{j=1}^{k} (z_j \log \hat{z}_j + (1 - z_j) \log(1 - \hat{z}_j)), \quad (4)$$

From both the Eqn. 1 and Eqn. 2, we can combine and can get a generalized loss

$$L = H - \log J \quad (5)$$

Our aim is to minimize the loss function, and, to do so, we can maximize the probabilities for correct pixels to be predicted and maximize the intersection, $J$ between masks and corresponding predictions.

Another commonly used performance metric is the DICE coefficient. Given the set of all pixels in the image, set of foreground pixels by automated segmentation $S_1^a$, and the set of pixels for ground truth $S_1^g$, DICE score can be compared

with $[S_1^a, S_1^g] \subseteq \Omega$, when a vector of ground truth labels $T_1$ and a vector of predicted labels $P_1$,

$$D(T_1, P_1) = \frac{2|T_1 \cap P_1|}{|T_1| + |P_1|} \quad (6)$$

DICE score will measure the similarity between two sets, $T_1$ and $P_1$ and $|T_1|$ denotes the cardinality of the set $T_1$ with the range of $D(T_1, P_1) \epsilon$ [0,1].

## III. RESULTS

### A. Quantitative Results

To illustrate the potential improvement in segmentation performance by using the nearest-neighbor interpolation (i.e., fixed upsampling) in the decoder, we conducted a pair comparison between the segmentation results obtained using the classical U-Net architecture, U-Net + NN, TernausNet, and U-NetPlus (our proposed method).

Training accuracy for binary segmentation is shown in Fig. 3 for 100 epochs. We compare our proposed architecture with three other models: U-Net, U-Net+NN, TernausNet. We can observe from this figure that after adding nearest-neighbor (NN) in the decoder of U-Net, the training accuracy of the classical U-Net framework (shown in blue) featuring the transposed convolution in the decodes, improves. Furthermore, the training of our proposed method (U-NetPlus) also converges faster and yields better training accuracy compared to TernausNet (shown in cyan). Hence, this graph alone illustrates the benefit of the nearest-neighbor interpolation on the segmentation performance.

The model was tested on the MICCAI 2017 EndoVis dataset. Table I summarizes the performance of our proposed U-NetPlus framework in the context of several state-of-the art multi-task segmentation techniques. The table clearly indicates the improvement in segmentation following the addition of nearest-neighbor interpolation in the decoder step across all frameworks — U-Net and TernausNet. Moreover, our model had been compared with four different structures other than U-Net and TernausNet — ToolNetH, ToolNetMS, FCN-8s, and CSL. The last one (CSL) was the first approach to multi-class surgical instrument segmentation. But, they used only two instrument classes (shaft and claspers) and omit wrist class which we introduced in our approach and the overall accuracy that we obtained was significantly higher than CSL approach.

We conducted a paired statistical test to compare the segmentation performance of each of these methods (U-Net, U-Net+NN, TernausNet, U-NetPlus) in terms of the IoU and DICE metric. As illustrated, our proposed U-NetPlus architecture yielded a statistically significant[2] 11.01% improvement ($p < 0.05$) in IoU and 6.91% DICE ($p < 0.05$) over the classical U-Net framework; a statistically significant 8.0% improvement ($p < 0.05$) in IoU and 5.79% DICE ($p < 0.05$) over the U-Net + NN framework; a statistically significant

---

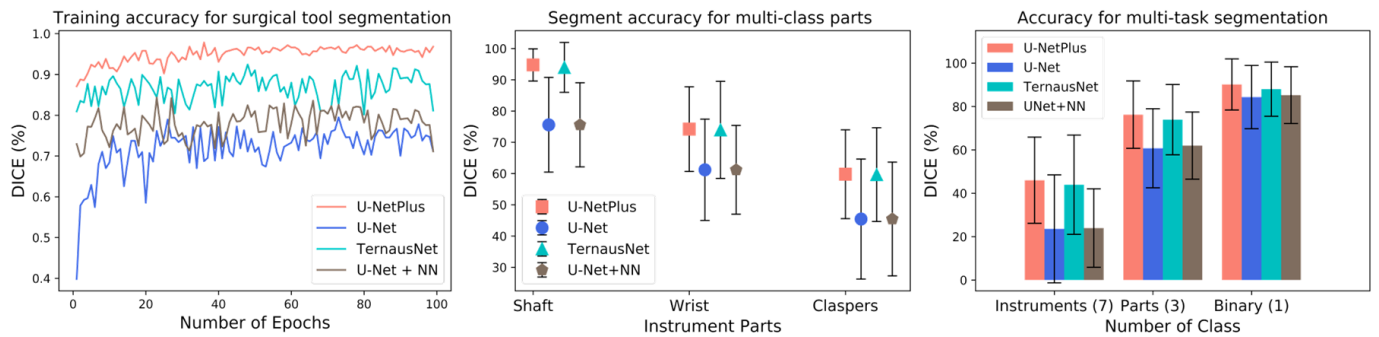[1] https://github.com/pytorch/pytorch

Fig. 3. Quantitative comparison of (a) training accuracy (left), (b) multi-class (class=3) instrument parts (middle) (c) multi-task segmentation accuracy (right).

TABLE I

QUANTITATIVE EVALUATION OF THE SEGMENTATION RESULTS. MEAN AND (STANDARD DEVIATION) VALUES ARE REPORTED FOR IOU(%) AND DICE COEFFICIENT(%) FROM ALL NETWORKS AGAINST THE PROVIDED REFERENCE SEGMENTATION. THE STATISTICAL SIGNIFICANCE OF THE RESULTS FOR U-NET + NN AND U-NETPLUS MODEL COMPARED AGAINST THE BASELINE MODEL (U-NET AND TERNASUNET) ARE REPRESENTED BY ∗ AND ∗∗ FOR P-VALUES 0.1 AND 0.05, RESPECTIVELY. U-NET HAS BEEN COMPARED WITH U-NET+NN, TERNAUSNET HAS BEEN COMPARED WITH U-NETPLUS. THE BEST PERFORMANCE METRIC (IOU AND DICE) IN EACH CATEGORY (BINARY, INSTRUMENT PART AND INSTRUMENT TYPE SEGMENTATION) IS INDICATED IN **BOLD** TEXT.

| Network | Binary Segmentation | | Instrument Part | | Instrument Type | |
|---|---|---|---|---|---|---|
| Metric | IoU | DICE | IoU | DICE | IoU | DICE |
| ToolNetH [12] | 74.4 | 82.2 | - | - | - | - |
| ToolNetMS [12] | 72.5 | 80.4 | - | - | - | - |
| FCN-8s [12] | 70.9 | 78.8 | - | - | - | - |
| CSL [24] | - | 88.9 | - | 87.70 (Shaft) | - | - |
| U-Net [7] | 75.44 | 84.37 | 48.41 | 60.75 | 15.80 | 23.59 |
| | (18.18) | (14.58) | (17.59) | (18.21) | (15.06) | (19.87) |
| **U-Net + NN** | **77.05**∗∗ | **85.26**∗ | **49.39**∗ | **61.98**∗ | **16.72**∗ | **23.97** |
| | (15.71) | (13.08) | (15.18) | (15.47) | (13.45) | (18.08) |
| TernausNet [13] | 83.60 | 90.01 | 65.50 | 75.97 | 33.78 | 44.95 |
| | (15.83) | (12.50) | (17.22) | (16.21) | (19.16) | (22.89) |
| **U-NetPlus-VGG-11** | 81.32 | 88.27 | 62.51 | 74.57 | **34.84**∗ | **46.07**∗∗ |
| | (16.76) | (13.52) | (18.87) | (16.51) | (14.26) | (16.16) |
| **U-NetPlus-VGG-16** | **83.75** | **90.20**∗ | **65.75** | **76.26**∗ | 34.19 | 45.32 |
| | (13.36) | (11.77) | (14.74) | (13.54) | (15.06) | (17.86) |
| | | | | 94.75(Shaft) | | |

0.18% improvement in IoU and 0.21% DICE (p < 0.1) over the state-of-the-art TernausNet framework [13].

Multi-class instrument segmentation was performed by labeling each instrument pixel with the corresponding index given in the training set. This application consisted of three classes: shaft, wrist, and claspers. The multi-class segmentation using our proposed U-NetPlus framework yielded a mean 65.75% IoU and 76.26% DICE. The accuracy and precision of the U-NetPlus architecture relative to the other three frameworks is illustrated in Fig. 3. As shown, the U-NetPlus framework outperforms the currently deemed best-in-class TernausNet framework.

The instrument type was segmented by labelling each instrument pixel with the corresponding instrument type, according to the training set, and all background pixels were labeled as 0. In the case of instrument type segmentation (for class = 7), U-NetPlus-VGG-11 encoder worked better than the U-NetPlus-VGG-16. Our results for instrument type segmentation can be further refined.

### B. Qualitative Results

The qualitative performance of our model both for a binary and multi-class segmentation is shown in Fig. 4. The second row of the figure shows that for the binary segmentation, the classical U-Net shows a portion of instrument which was not present in the binary mask of our ground truth data (second row and second column). U-netPlus shows the best performance for binary segmentation (i.e. it can clearly segment out the instruments from the background), whereas TernausNet still shows un-wanted regions in the segmentation output. For the instrument parts segmentation, U-Net still segments the un-wanted instrument (blue), whereas U-NetPlus can segment the 3 classes (blue:shaft, green:wrist, yellow:claspers) near perfectly compared to TernausNet. For the instrument type segmentation, we can clearly observe that U-Net can not differentiate between the blue and the green class, whereas U-NetPlus can differentiate these classes more accurately than TernausNet. Both the binary and multi-class segmented output have been overlayed onto the original image (sixth, seventh, eighth, and ninth column). The figure has a clear indication of
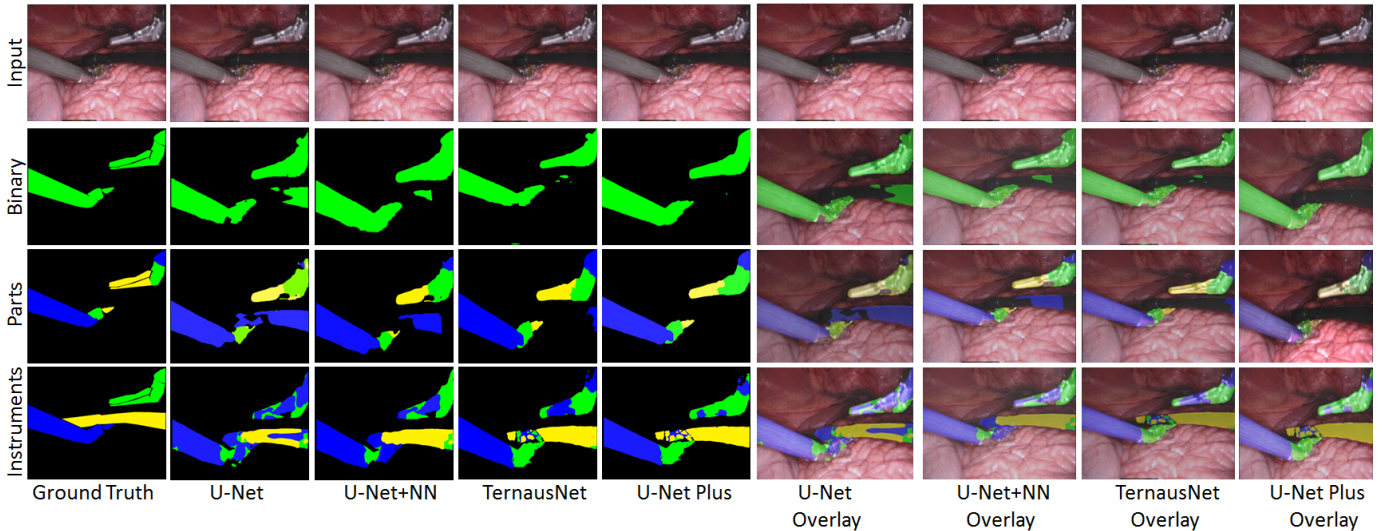
Fig. 4. Qualitative comparison of binary segmentation, instrument part and instrument type segmentation result and their overlay onto the native endoscopic images of the MICCAI 2017 EndoVis video dataset yielded by four different frameworks: U-Net, U-Net+NN, TernausNet, and U-NetPlus.
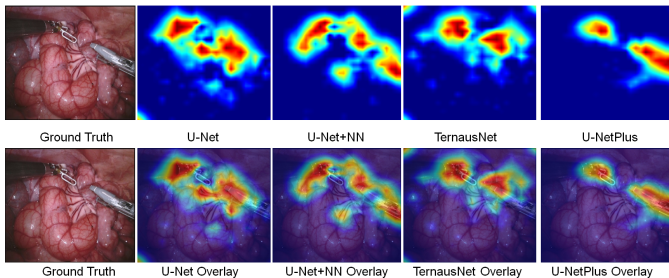


Fig. 5. Attention results: U-NetPlus "looks" at a focused target region, whereas U-Net, U-Net+NN and TernausNet appear less "focused", leading to less accurate segmentation.

qualitative improvement of U-NetPlus over U-Net, U-Net+NN and TernausNet.

### C. Attention Study

To further analyze the improvement in segmentation performance, we performed an attention analysis to visualize where our proposed algorithm "looks" in an image by using a novel image saliency technique [25] that learns the mask of an image by suppressing the softmax probability of its target class. Fig. 5 shows that using this class activation mapping, our approach (U-NetPlus) localizes the wrist and claspers of the bipolar forceps near perfectly compared to the classical U-Net, U-Net+NN and TernausNet frameworks. U-Net does not perform well compared to U-Net+NN where nearest-neighbor upsampling is added to the decoder path. TernausNet performs well than U-Net+NN due to the use of pre-trained VGG network in encoder. Pre-trained network performs well in this segmentation due to the limited dataset. Fig. 5 shows the heatmap image which has been overlayed onto it's original image. Therefore, the skillful integration and combination of

pre-trained encoder and nearest-neighbor interpolation as a fixed upsampling technique yields higher overall performance.

## IV. DISCUSSION AND CONCLUSION

In this paper, we proposed a modified U-Net model for surgical tool segmentation. To improve robustness beyond that of the U-Net framework, we used the pre-trained model as the encoder with batch-normalization, which converges much faster in comparison to the non-pre-trained network. In the decoder part, we substituted the deconvolution layer with an upsampling layer that uses nearest-neighbor interpolation followed by two convolution layers. Moreover, we used a fast and effective data augmentation technique to avoid the overfitting problem. We evaluated its performance on the MICCAI 2017 EndoVis Challenge dataset. We also visualized the output of our proposed method both as stand-alone surgical instrument segmentation, as well as overlays onto the native endoscopic images. Apart from that, we also conducted an "attention" study to determine where our proposed algorithm "looks" in an image.

Our proposed model with batch-normalized U-NetPlus-VGG-16 outperforms existing methods in terms of both Jaccard and DICE, achieving 90.20% DICE for binary class segmentation and 76.26% for parts segmentation, both of which showed at least 0.21% improvement over the current methods and more than 6% improvement over the traditional U-Net architecture. Nevertheless, U-NetPlus-VGG-16's performance with regards to identifying the instrument type was inferior to that of U-NetPlus-VGG-11, which was slightly superior to the other disseminated techniques rendering this paper as a first demonstration of a modified version of U-Net Decoder via nearest-neighbor interpolation to remove artifacts induced by the transposed convolution being used for surgical instrument segmentation application and showing improved performance over the TernausNet framework.

## REFERENCES

[1] P. P. Rao, "Robotic surgery: new robots and finally some real competition!" *World Journal of Urology*, vol. 36, no. 4, pp. 537–541, 2018.

[2] M. Yang, K. Yu, C. Zhang, Z. Li, and K. Yang, "Denseaspp for semantic segmentation in street scenes," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3684–3692.

[3] K. He, G. Gkioxari, P. Dollár, and R. Girshick, "Mask r-cnn," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2980–2988.

[4] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3431–3440.

[5] D. Shen, G. Wu, and H.-I. Suk, "Deep learning in medical image analysis," *Annual review of biomedical engineering*, vol. 19, pp. 221–248, 2017.

[6] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, "ImageNet: A Large-Scale Hierarchical Image Database," in *CVPR09*, 2009.

[7] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2015, pp. 234–241.

[8] C. Chen, Q. Chen, J. Xu, and V. Koltun, "Learning to see in the dark," *arXiv preprint arXiv:1805.01934*, 2018.

[9] N. Jiang and L. Wang, "Quantum image scaling using nearest neighbor interpolation," *Quantum Information Processing*, vol. 14, no. 5, pp. 1559–1571, 2015.

[10] X. Jia, H. Chang, and T. Tuytelaars, "Super-resolution with deep adaptive image resampling," *arXiv preprint arXiv:1712.06463*, 2017.

[11] A. Odena, V. Dumoulin, and C. Olah, "Deconvolution and checkerboard artifacts," *Distill*, vol. 1, no. 10, p. e3, 2016.

[12] L. C. García-Peraza-Herrera, W. Li, L. Fidon, C. Gruijthuijsen, A. Devreker, G. Attilakos, J. Deprest, E. Vander Poorten, D. Stoyanov, T. Vercauteren *et al.*, "Toolnet: holistically-nested real-time segmentation of robotic surgical tools," in *Intelligent Robots and Systems (IROS), 2017 IEEE/RSJ International Conference on*. IEEE, 2017, pp. 5717–5722.

[13] A. Shvets, A. Rakhlin, A. A. Kalinin, and V. Iglovikov, "Automatic instrument segmentation in robot-assisted surgery using deep learning," *arXiv preprint arXiv:1803.01207*, 2018.

[14] D. Pakhomov, V. Premachandran, M. Allan, M. Azizian, and N. Navab, "Deep residual learning for instrument segmentation in robotic surgery," *arXiv preprint arXiv:1703.08580*, 2017.

[15] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.

[16] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," in *Advances in Neural Information Processing Systems*, 2016, pp. 2234–2242.

[17] V. Iglovikov and A. Shvets, "Ternausnet: U-net with vgg11 encoder pre-trained on imagenet for image segmentation," *arXiv e-prints arXiv:1801.05746*, 2018.

[18] K. He, R. Girshick, and P. Dollár, "Rethinking imagenet pre-training," *arXiv preprint arXiv:1811.08883*, 2018.

[19] S. Kornblith, J. Shlens, and Q. V. Le, "Do better imagenet models transfer better?" *arXiv preprint arXiv:1805.08974*, 2018.

[20] S. Santurkar, D. Tsipras, A. Ilyas, and A. Madry, "How does batch normalization help optimization?(no, it is not about internal covariate shift)," *arXiv preprint arXiv:1805.11604*, 2018.

[21] C. Dong, C. C. Loy, and X. Tang, "Accelerating the super-resolution convolutional neural network," in *European Conference on Computer Vision*. Springer, 2016, pp. 391–407.

[22] *MICCAI 2017 Endoscopic Vision Challenge: Robotic Instrument Segmentation Sub-Challenge*, 2017, https://endovissub2017-roboticinstrumentsegmentation.grand-challenge.org/Data/.

[23] E. K. V. I. I. A. Buslaev, A. Parinov and A. A. Kalinin, "Albumentations: fast and flexible image augmentations," *arXiv e-prints arXiv:1809.06839*, 2018.

[24] I. Laina, N. Rieke, C. Rupprecht, J. P. Vizcaíno, A. Eslami, F. Tombari, and N. Navab, "Concurrent segmentation and localization for tracking of surgical instruments," in *International conference on medical image computing and computer-assisted intervention*. Springer, 2017, pp. 664–672.

[25] R. C. Fong and A. Vedaldi, "Interpretable explanations of black boxes by meaningful perturbation," *arXiv preprint arXiv:1704.03296*, 2017.