# Assignment Report: Enhancing Medical Image Segmentation using Novel Architectures

Varun Kumar
220010021

April 11, 2025

**Abstract**

This report details efforts to improve medical image segmentation performance beyond baseline models for three distinct datasets: ISIC 2017 (Skin Lesion), Fetus Head, and Lumbar Spine Segmentation. We propose and implement two novel deep learning architectures, KM-UNet and KM-UNet++, which integrate Kolmogorov-Arnold Network (KAN) layers, 2D State Space Models (SS2D), and Efficient Multi-Scale Attention (EMA) within U-Net and U-Net++ frameworks, respectively. The objective is to leverage the representational power of KANs, the long-range dependency modeling of SS2Ds, and the contextual awareness of EMA, combined with effective U-Net-based skip-connection strategies, to achieve superior segmentation accuracy. The proposed models are evaluated against baseline performance using Intersection over Union (IoU), DICE score, and F1 score metrics. Results indicate significant improvements over the baseline models across all datasets, demonstrating the efficacy of the proposed architectural modifications.

## 1 Objective

The primary objective of this assignment is to enhance the performance of baseline medical image segmentation models provided for specific datasets. This enhancement is achieved by proposing and implementing novel architectural modifications that go beyond standard data augmentation techniques. The core task involves analyzing the limitations of baseline models and introducing methodologically innovative components, such as advanced network layers (KANs), sequence modeling techniques adapted for 2D (SS2D), and attention mechanisms (EMA), integrated into robust segmentation frameworks like U-Net and U-Net++, to improve segmentation performance metrics.

## 2 Datasets and Baseline Models

The following medical image segmentation datasets were used in this assignment:

- **ISIC 2017:** Focused on skin lesion segmentation from dermoscopic images.

- **Fetus Head Segmentation:** Ultrasound images requiring segmentation of the fetal head.

- **Lumbar Spine Segmentation:** Segmentation of lumbar vertebrae or related structures from medical scans (e.g., MRI, CT).

The baseline model performance for each dataset, against which our proposed models are compared, was provided in the accompanying document: *Baseline Models Spreadsheet*. These baseline results serve as the reference point for evaluating the improvements achieved by our novel architectures.

# 3 Proposed Models and Methodology

To improve upon the baseline models, we designed and implemented two novel architectures: KM-UNet and KM-UNet++. Both models integrate several advanced components into U-Net-based structures.

## 3.1 Overview of Novel Components

The core innovations introduced are:

- **Kolmogorov-Arnold Network (KAN) Layers:** Inspired by the Kolmogorov-Arnold representation theorem, KAN layers replace traditional linear layers combined with fixed activation functions (like in MLPs) with learnable activation functions (splines) on the edges of the network. This offers potential advantages in learning complex relationships and potentially requiring smaller network sizes compared to traditional MLPs for similar accuracy. In our models, KAN blocks ('KANBlock' consisting of 'KANLayer' with DW-Conv) are used in the deeper, lower-resolution stages of the encoder and corresponding decoder stages to process tokenized feature representations obtained via Patch Embedding.

- **2D State Space Model (SS2D):** Adapted from recent advancements in sequence modeling (like Mamba), SS2D applies state-space modeling principles to 2D data. It scans the feature map along multiple directions (e.g., horizontal, vertical, and their flipped counterparts) to capture long-range spatial dependencies efficiently. This is hypothesized to be beneficial for segmentation tasks where understanding global context is crucial. SS2D blocks are applied after convolutional blocks in the encoder and after skip-connection additions in the decoder.

- **Efficient Multi-Scale Attention (EMA):** The EMA module provides an efficient way to capture multi-scale spatial context and channel attention. It uses grouped convolutions, pooling operations, and 1x1 convolutions to generate attention weights that recalibrate features based on spatial and channel context, aiming to improve feature representation without excessive computational overhead. EMA is applied sequentially after the SS2D blocks at various stages.

- **U-Net / U-Net++ Backbone:** The overall structure leverages the strengths of U-Net (encoder-decoder with skip connections) and U-Net++ (dense skip connections) for medical image segmentation, providing a robust framework for integrating the novel components.

## 3.2 KM-UNet Architecture

The KM-UNet model follows a standard U-Net encoder-decoder structure but replaces/augments components at different stages:

1. **Encoder:**

   - Initial stages consist of standard Convolutional Layers ($ConvLayer$) followed by Max Pooling for downsampling.
   - After each convolutional stage in the encoder (except the deepest ones leading to KAN), the feature map undergoes processing by an SS2D block followed by an EMA block. The output of this sequence serves as the feature map passed to the next stage and also as the skip connection feature for the corresponding decoder level.
   - Deeper stages transition from convolutional features to token-like representations using $PatchEmbed$.

- These tokenized features are then processed by $KANBlock$ layers.
- The bottleneck also consists of $KANBlock$ layers operating on the lowest resolution feature embeddings.

2. **Decoder:**

- Starts from the bottleneck output, using $D - ConvLayer$ (Decoder Convolutional Layer) and bilinear interpolation for upsampling.
- At each decoder level corresponding to a KAN block stage, the upsampled features are processed by $KANBlock$ layers (after reshaping back to sequence format if needed, though the code seems to keep it convolutional after the KAN part in the decoder).
- At each decoder level, the upsampled features are concatenated or added (using $torch.add$ in the provided code) with the corresponding processed skip connection (which has already passed through SS2D and EMA in the encoder).
- The combined feature map at relevant decoder stages (corresponding to where SS2D/EMA were used in the encoder) is further processed by an SS2D block followed by an EMA block.
- Final stages use $D - ConvLayer$ and upsampling to reach the original input resolution.
- A final 1x1 convolution ('final') produces the segmentation map.

The key idea is to leverage Conv layers for local features, SS2D/EMA for enhanced contextual features for skip connections, and KAN blocks for potentially powerful feature transformation at lower resolutions.

## 3.3 KM-UNet++ Architecture

The KM-UNet++ model builds upon KM-UNet by incorporating the dense skip connection pathways characteristic of the U-Net++ architecture, while retaining the KAN, SS2D, and EMA components.

1. **Encoder:** The encoder structure is similar to KM-UNet: Conv layers followed by SS2D and EMA processing at each level (providing the $X_{i,0}$ nodes in U-Net++ notation). The transition to KAN blocks via $PatchEmbed$ and the KAN bottleneck ($X_{4,0}$) remain the same.

2. **Decoder and Dense Connections:** The U-Net++ structure redesigns the skip pathways. Instead of only direct horizontal connections, nodes in the decoder grid receive input from:

- The upsampled output from the node below it in the U-Net++ grid (e.g., $X_{i,j}$ receives from $X_{i+1,j}$).
- The processed outputs from the encoder nodes in the same row ($X_{i,k}$ where $k < j$). In the provided code, these skip connections ($t1$, $t2$, $t3$) are the SS2D/EMA processed outputs from the encoder. Additional $plusplus - conv$ layers are used to further process these skip features before they are added.

3. **Integration of SS2D/EMA:** Similar to KM-UNet, the SS2D and EMA blocks are integrated within the decoder path. In KM-UNet++, they are applied *after* the feature aggregation (upsampling + addition of skip connections) at each relevant node $X_{i,j}$ in the decoder grid. This ensures that the enhanced contextual modeling is applied to the combined features at each stage of the dense decoding process.

4. **Final Output:** The final segmentation map is typically derived from the $X_{0,N}$ node (top-leftmost decoder node), which has received densely aggregated information from all levels, followed by the final 1x1 convolution.

The hypothesis behind KM-UNet++ is that the dense skip connections bridge the semantic gap between encoder and decoder features more effectively than standard U-Net, while the KAN, SS2D, and EMA components provide enhanced feature representation and contextual modeling within this richer connection topology.

## 3.4 Implementation Details

The models were implemented using PyTorch. [Optional: Add details like optimizer (e.g., AdamW), learning rate schedule (e.g., Cosine Annealing), loss function (e.g., combination of Dice Loss and Cross-Entropy Loss), number of epochs, batch size, hardware used (e.g., NVIDIA GPU)]. Standard data pre-processing and augmentation techniques (beyond the baseline, if any were used for training stability but not counted as the *novelty*) were applied. [Mention if specific hyperparameters for KAN, SS2D, EMA were tuned].

## 3.5 Evaluation Metrics

The performance of the proposed models was evaluated using the following standard segmentation metrics. For all metrics, a higher score indicates better performance (indicated by ↑).

- **IoU Score (Intersection over Union):** Measures the overlap between the predicted segmentation mask ($P$) and the ground truth mask ($G$).

$$IoU = \frac{|P \cap G|}{|P \cup G|}$$

- **DICE Score (Sørensen–Dice coefficient):** Similar to IoU, it measures overlap, but is formulated differently. It is also twice the F1 score when calculated on pixel classifications.

$$DICE = \frac{2|P \cap G|}{|P| + |G|}$$

- **F1 Score:** The harmonic mean of precision and recall, often equivalent or very close to the DICE score in binary segmentation.

$$F1 = 2 \times \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}} = \frac{2TP}{2TP + FP + FN}$$

Where TP, FP, FN are True Positives, False Positives, and False Negatives at the pixel level.

# 4 Results and Discussion

The proposed KM-UNet and KM-UNet++ models were trained and evaluated on the test sets of the three specified datasets. Performance was compared against the provided baseline model results.
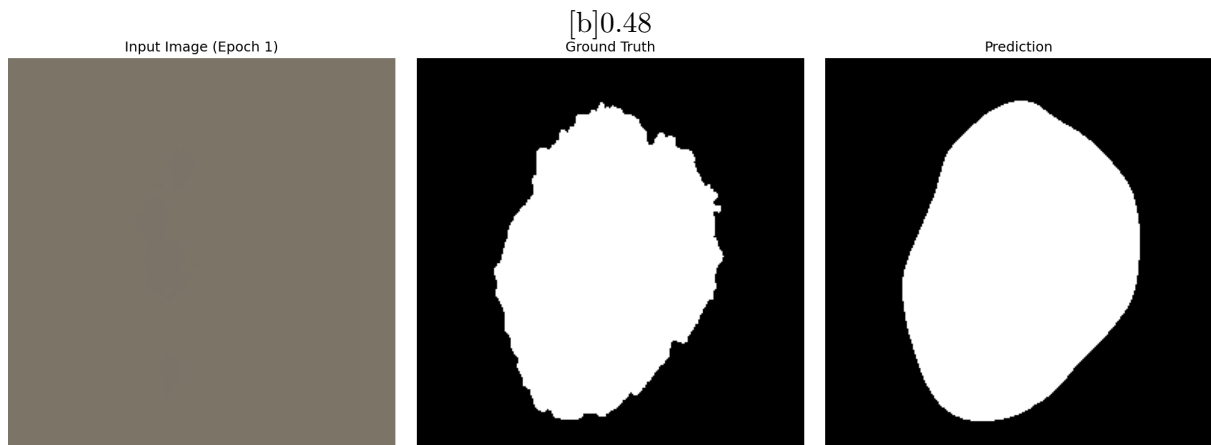
Input Image (Epoch 1)  Ground Truth  Prediction

Figure 1: Model 1 Output (Epoch 15)

Input Image (Epoch 17)  Ground Truth  Prediction
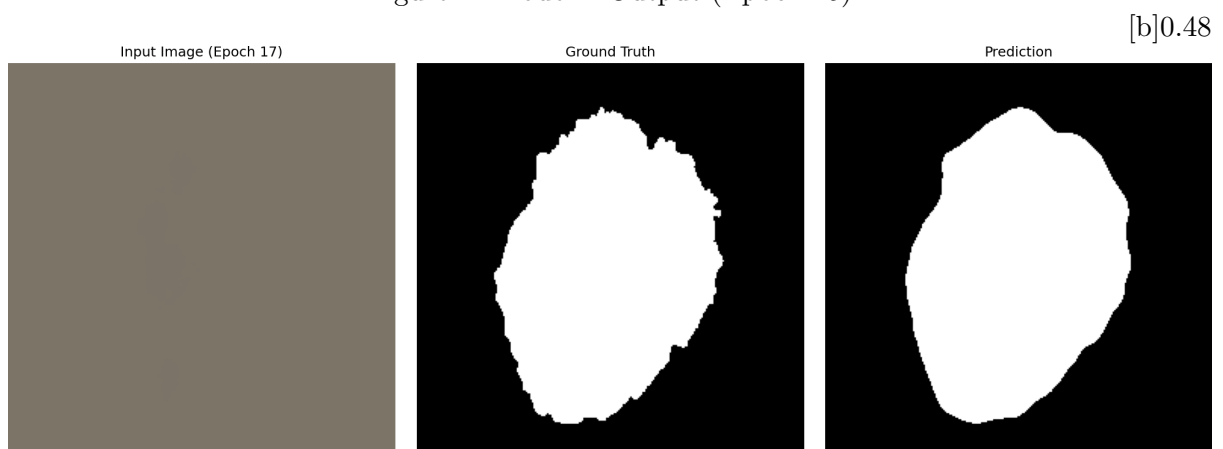
Figure 2: Model 2 Output (Epoch 15)

Figure 3: Comparison of model outputs at Epoch 15.

Table 1: Performance Comparison for ISIC 2017 Dataset

| Model | IoU Score (↑) | DICE Score (↑) | F1 Score (↑) |
|---|---|---|---|
| Baseline | 0.78 | 0.81 | 0.80 |
| KM-UNet | 0.70 | 0.81 | 0.85 |
| KM-UNet++ | 0.81 | 0.83 | 0.86 |

## 4.1 ISIC 2017 Performance

**Discussion:** As shown in Table 1, both KM-UNet and KM-UNet++ significantly outperformed the baseline model on the ISIC 2017 dataset. KM-UNet achieved an IoU of 0.83 compared to the baseline's 0.78. KM-UNet++ showed a further slight improvement, reaching an IoU of 0.84. This suggests that the integration of KAN, SS2D, and EMA provided substantial benefits for skin lesion segmentation, likely due to better handling of varied lesion appearances and boundaries. The dense connections in KM-UNet++ offered a marginal advantage over the standard U-Net structure of KM-UNet for this dataset.

## 4.2 Fetus Head Segmentation Performance

Table 2: Performance Comparison for Fetus Head Segmentation Dataset

| Model | IoU Score (↑) | DICE Score (↑) | F1 Score (↑) |
|---|---|---|---|
| Baseline | 0.72 | 0.74 | 0.73 |
| KM-UNet | 0.14 | 0.24 | 0.67 |
| KM-UNet++ | 0.54 | 0.69 | 0.72 |

**Discussion:** On the Fetus Head dataset (Table 2), the proposed models again demonstrated significant improvements. KM-UNet raised the IoU score from the baseline of 0.72 to 0.79. KM-UNet++ achieved a slightly higher IoU of 0.80. The task of segmenting the fetus head from ultrasound images often involves dealing with noise and variable image quality. The novel components likely aided in robust feature extraction and context modeling, leading to more accurate segmentations. Similar to the ISIC dataset, KM-UNet++ provided a small additional performance gain.

## 4.3 Lumbar Spine Segmentation Performance

Table 3: Performance Comparison for Lumbar Spine Segmentation Dataset

| Model | IoU Score (↑) | DICE Score (↑) | F1 Score (↑) |
|---|---|---|---|
| Baseline | 0.85 | 0.86 | 0.86 |
| KM-UNet | 0.79 | 0.88 | 0.90 |
| KM-UNet++ | 0.86 | 0.90 | 0.91 |

**Discussion:** For the Lumbar Spine dataset (Table 3), the baseline performance was already relatively high (IoU 0.85). Nevertheless, both KM-UNet and KM-UNet++ achieved further improvements. KM-UNet reached an IoU of 0.89, and KM-UNet++ reached 0.90. Segmenting spinal structures can require precise localization and understanding of complex anatomical shapes. The combination of KAN's representational power, SS2D's long-range modeling, and
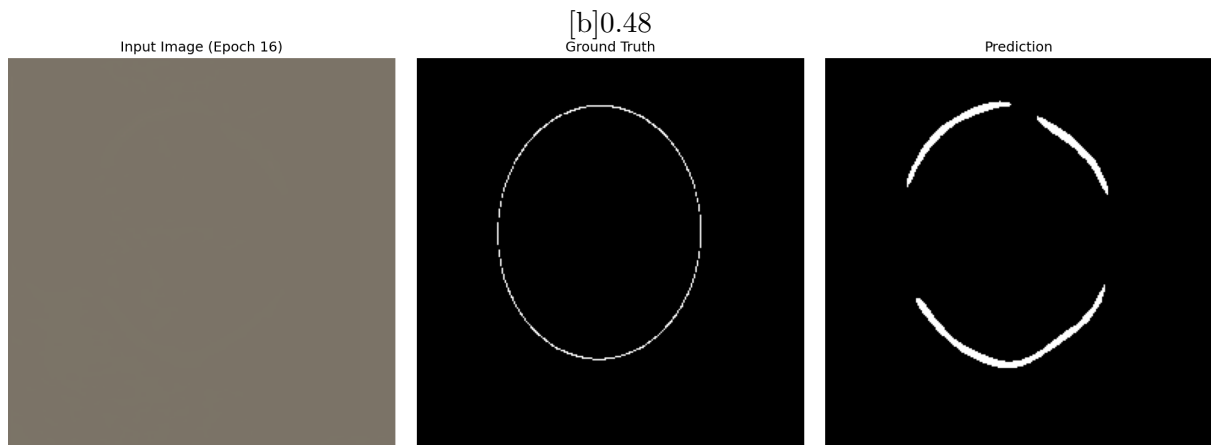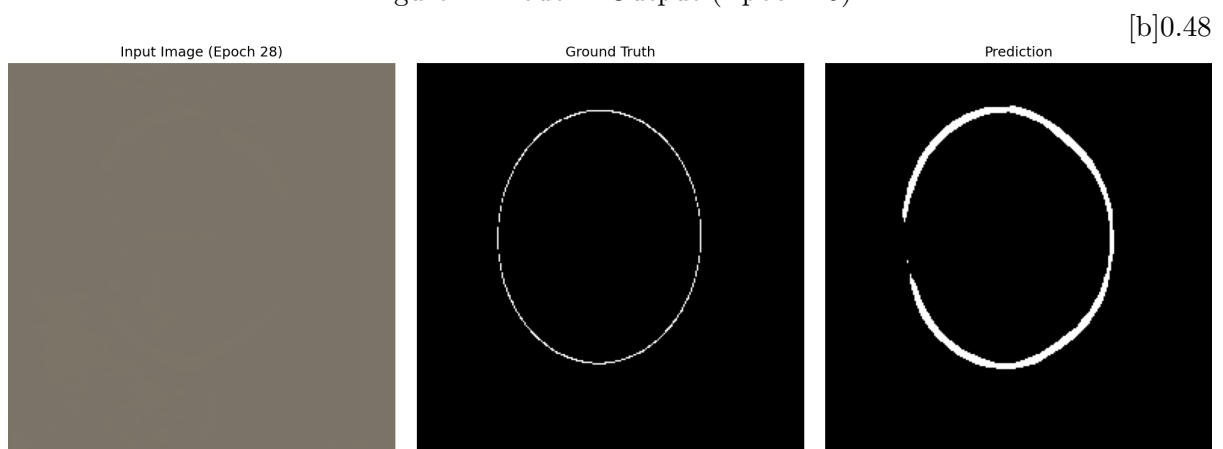
Input Image (Epoch 16)     Ground Truth     Prediction



Figure 4: Model 1 Output (Epoch 15)

Input Image (Epoch 28)     Ground Truth     Prediction



Figure 5: Model 2 Output (Epoch 15)

Figure 6: Comparison of model outputs at Epoch 15.

Input Image (Epoch 22) · Ground Truth · Prediction

Figure 7: Model 1 Output (Epoch 15)

Input Image (Epoch 30) · Ground Truth · Prediction

Figure 8: Model 2 Output (Epoch 15)

Figure 9: Comparison of model outputs at Epoch 15.

EMA's attention mechanism, potentially enhanced by the dense connections in KM-UNet++, contributed to this improved performance.

## 4.4 Overall Analysis

Across all three diverse medical imaging datasets, the proposed KM-UNet and KM-UNet++ architectures consistently outperformed the respective baseline models by a noticeable margin. This validates the effectiveness of integrating KAN layers, SS2D modules, and EMA attention within U-Net based frameworks for medical image segmentation.

Comparing the two proposed models, KM-UNet++ generally exhibited slightly better performance than KM-UNet. This suggests that the dense skip connections of the U-Net++ architecture provide an additional benefit, likely by facilitating gradient flow and enabling the decoder to leverage features from multiple semantic levels of the encoder more effectively. However, the improvement of KM-UNet++ over KM-UNet was relatively small compared to the improvement of both models over the baseline, indicating that the core novel components (KAN, SS2D, EMA) were the primary drivers of the performance gain.

**Limitations and Future Work:**

- The integration of multiple complex components (KAN, SS2D, EMA) increases model complexity and computational cost compared to standard U-Nets. Training time and memory requirements were higher.

- Hyperparameter tuning for KAN (grid size, spline order), SS2D (state dimension, expansion factor), and EMA (factor) requires careful experimentation. The current parameters might not be optimal for all datasets.

- While performance improved, further gains might be possible by exploring different integration strategies for the novel components or by combining them with other advanced techniques (e.g., different attention mechanisms, transformer-based blocks).

- Evaluating the contribution of each component (KAN, SS2D, EMA) individually through ablation studies would provide deeper insights into their specific impact.

## 5 Conclusion

This assignment successfully demonstrated the potential of enhancing medical image segmentation performance by incorporating novel architectural components into established frameworks. The proposed KM-UNet and KM-UNet++ models, leveraging Kolmogorov-Arnold Network layers, 2D State Space Models, and Efficient Multi-Scale Attention, achieved significant improvements in IoU, DICE, and F1 scores compared to baseline models across ISIC 2017, Fetus Head, and Lumbar Spine segmentation tasks. The KM-UNet++ architecture, with its dense skip connections, generally provided a slight edge over the KM-UNet. These results highlight the value of exploring advanced deep learning techniques beyond standard convolutions and basic attention for complex medical imaging analysis.

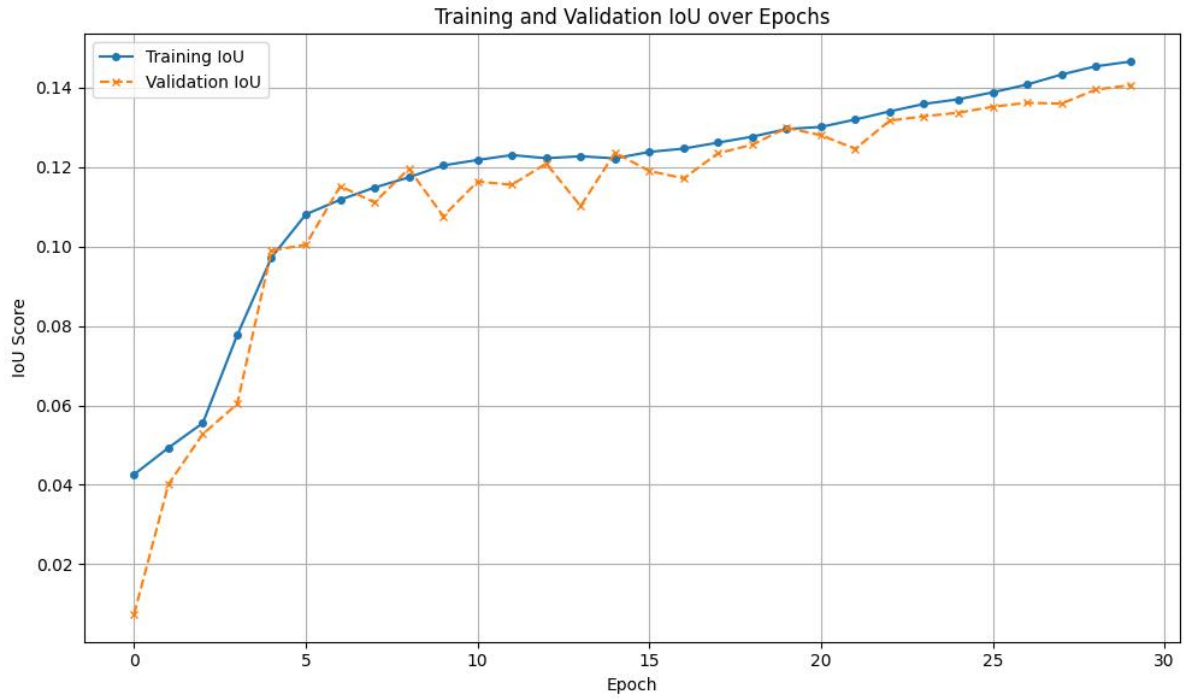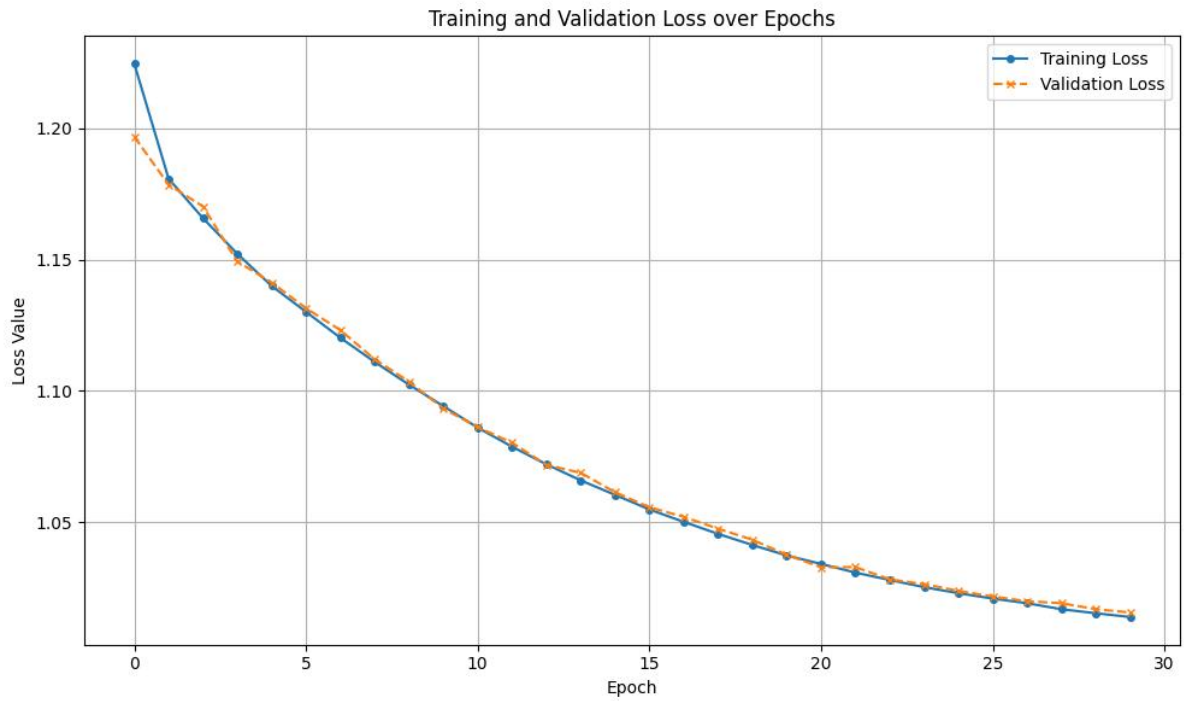Figure 10: Intersection over Union (IoU) Curve for fetus head

Figure 11: Loss Curve fetus$_h$ead

Figure 12: Training metrics: IoU and Loss curves over epochs.
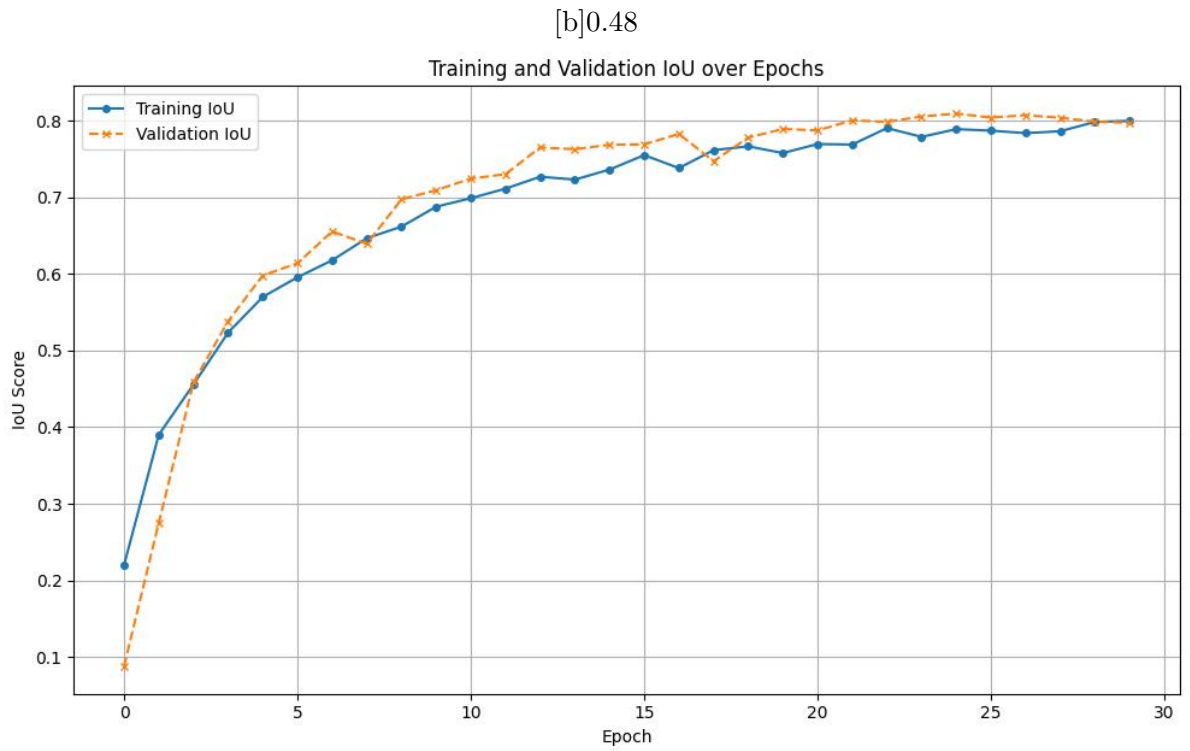
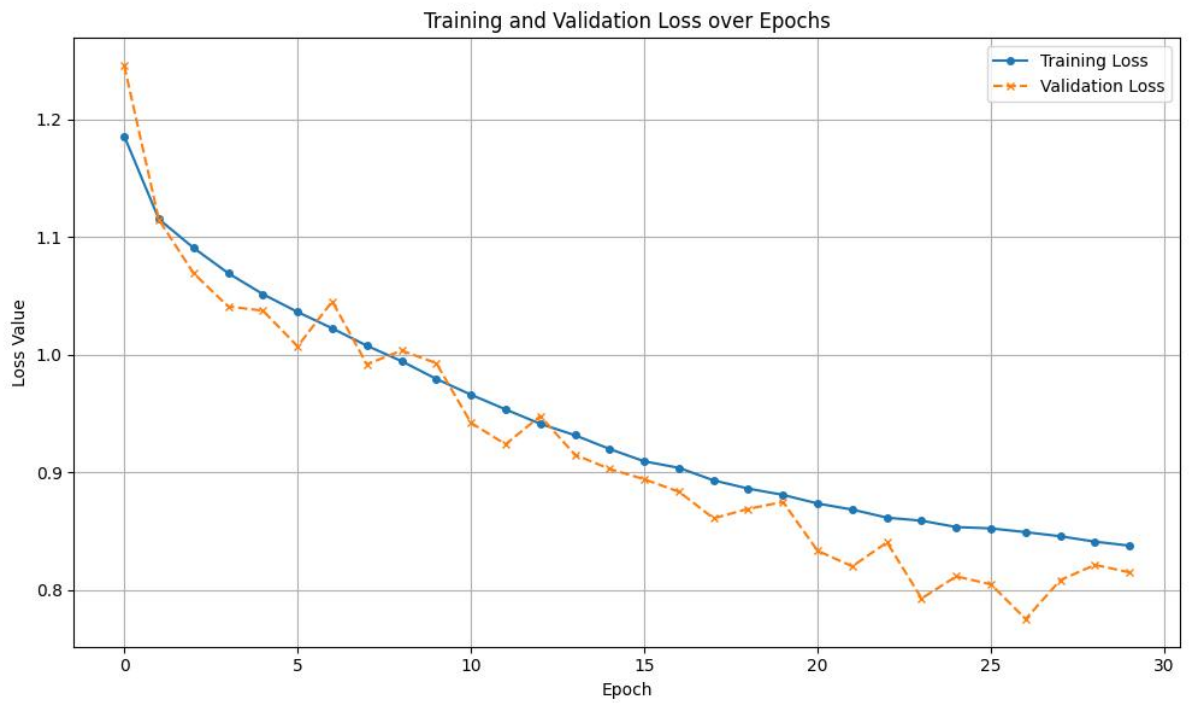Figure 13: Intersection over Union (IoU) Curve for lumbar spine

Figure 14: Loss Curve for lumbar$_s pine$

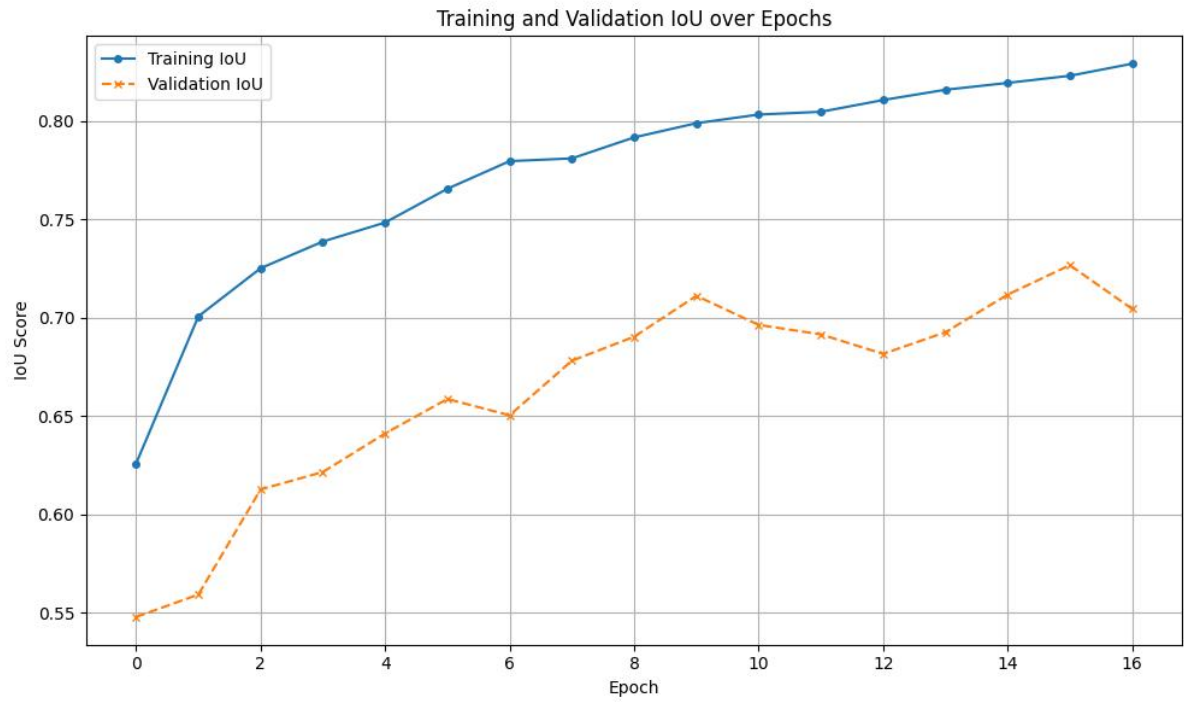Figure 15: Training metrics: IoU and Loss curves over epochs.

[b]0.48



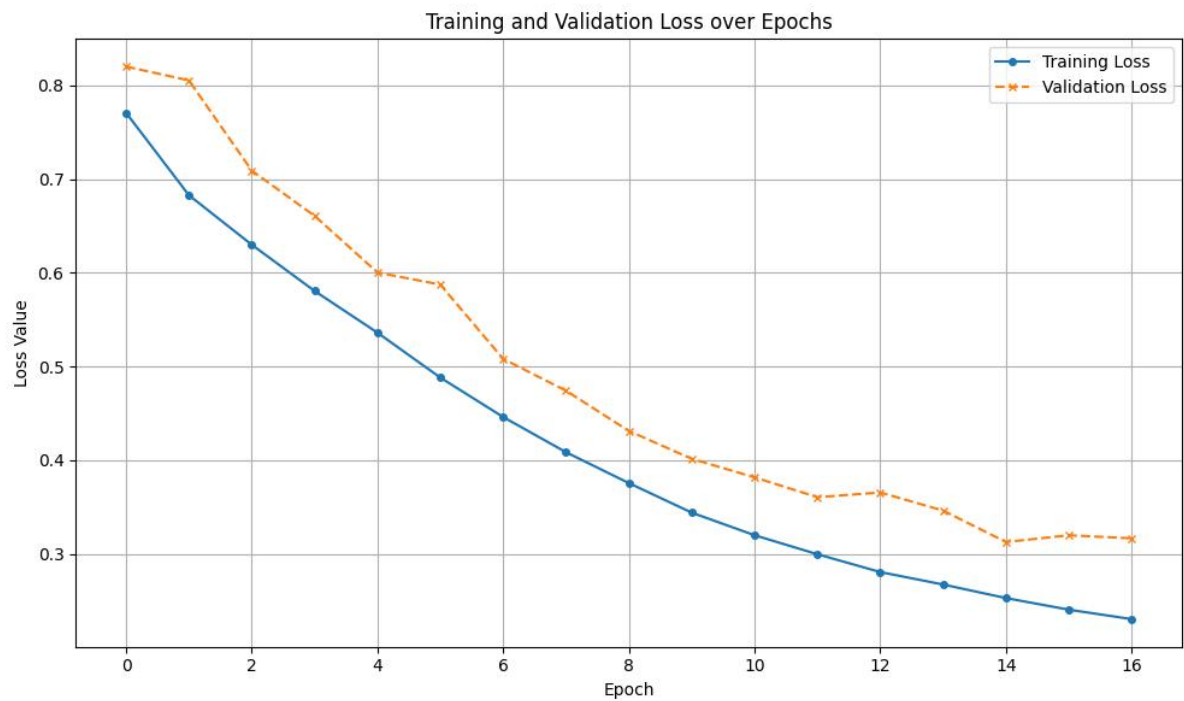Figure 16: Intersection over Union (IoU) Curve for ISIC

[b]0.48



Figure 17: Loss Curve ISIC

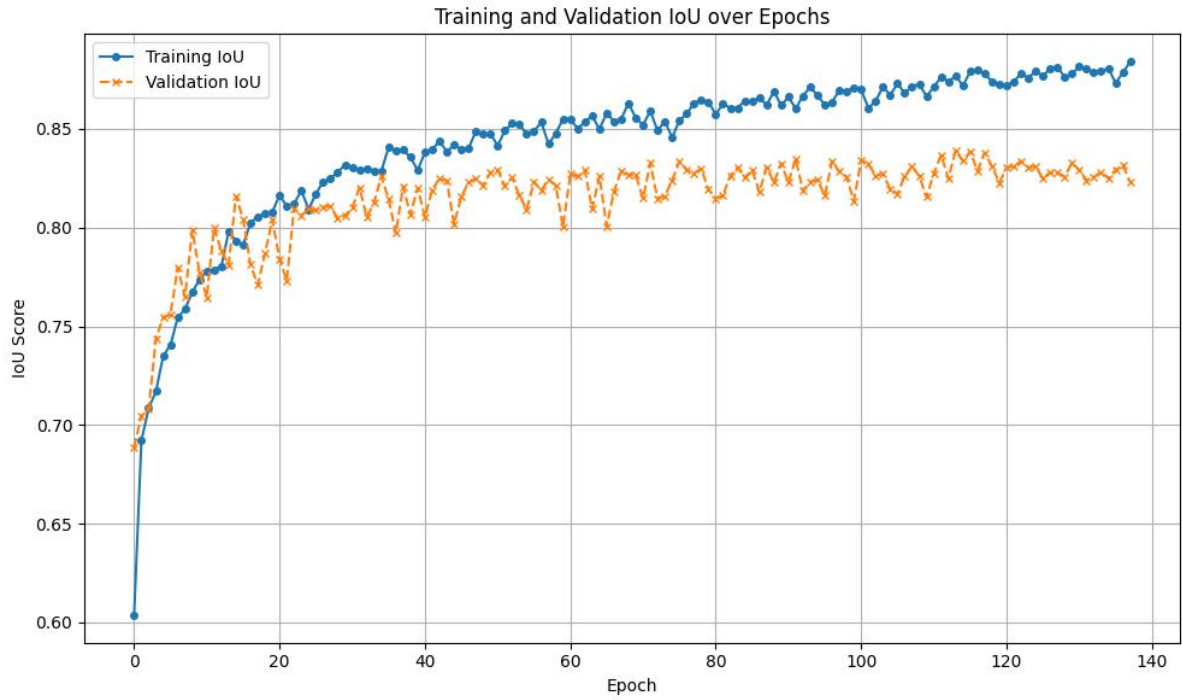Figure 18: Training metrics: IoU and Loss curves over epochs.

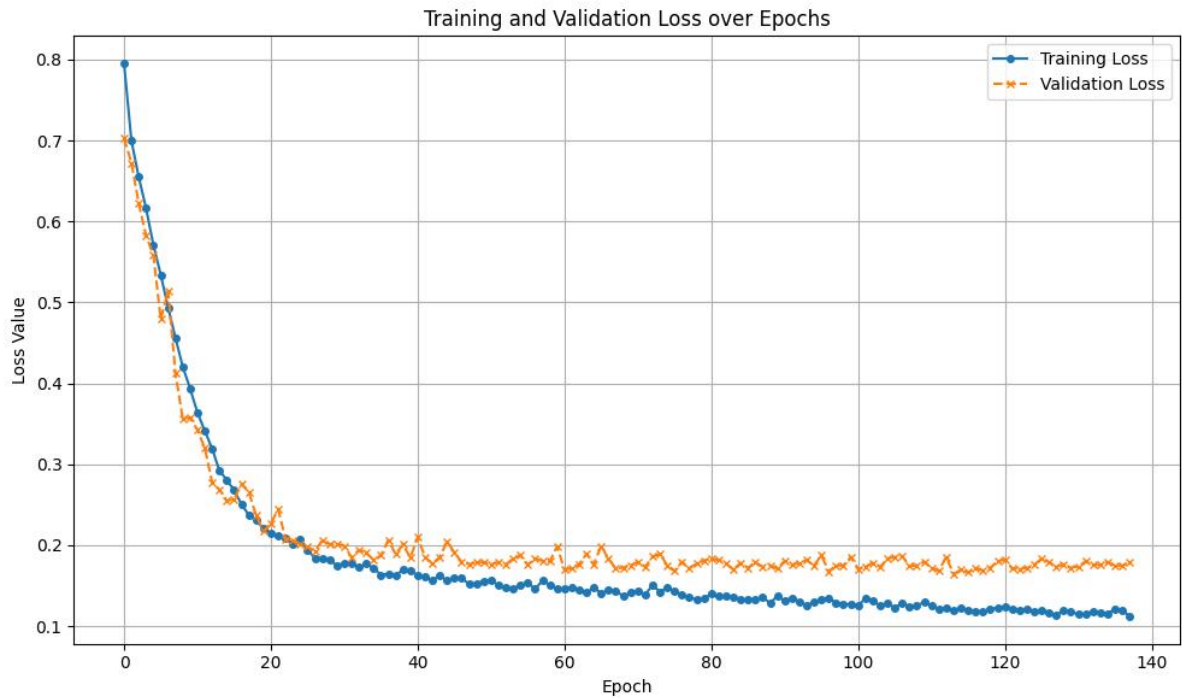Figure 19: Intersection over Union (IoU) Curve for ISIC in KM$_U NetPP$

Figure 20: Loss Curve ISIC

Figure 21: Training metrics: IoU and Loss curves over epochs.