

Enhancing Medical Image Segmentation with Self-Attention-Embedded LSTM Mechanism

Abstract—Medical image segmentation is crucial for clinical applications but remains challenging due to complex anatomical structures and the need for precise localization. This paper introduces LSTMSA, a novel deep learning model that combines the Self-Attention (SA) mechanism with a Long Short-Term Memory (LSTM) module to address key issues in medical image segmentation. Traditional methods often face premature convergence, leading to suboptimal outcomes. LSTMSA mitigates this by embedding SA within LSTM, ensuring a balanced training process that effectively utilizes LSTM’s temporal capabilities while focusing on local details. Additionally, SA can cause rank collapse, leading to information loss and blurred segmentations. LSTMSA addresses this by preserving high-dimensional feature representations, maintaining both details and context. Experimental results across four medical image datasets and comparisons with twenty baselines confirm LSTMSA’s superiority, achieving state-of-the-art accuracy and preserving essential anatomical features. The code is available at <https://github.com/yeshunlong/LSTMSA>.

Index Terms—medical image segmentation, deep learning, long short-term memory (LSTM), self-attention (SA), rank preservation, feature representation.

I. INTRODUCTION

THE nature of medical images, characterized by intricate anatomical structures and the need for precise localization, demands innovative approaches to achieve accurate and robust segmentation results [1], [2], [3]. Recent years, researchers have been dedicated to enhance model performance through various contextual modules. For example, they have explored the individual implementation of Long Short-Term Memory (LSTM) and Self-Attention (SA) mechanism to enhance the extraction of contextual information. LSTM is adept at capturing temporal and sequential dependencies within the data [4], while SA mechanism excels in modeling global contextual relationships [5]. Each approach independently contributes to the improvement of context-awareness in segmentation models, addressing the critical challenge of capturing long-range dependencies within medical images [6]. These strategies have yielded promising results in medical imaging applications, contributing to more precise segmentation and clinically relevant outcomes [7]. Furthermore, such advanced methodologies enhance the model’s ability to interpret complex image sequences and enables the dynamic adjustment of focus [8], thereby improving the detection of subtle nuances and variations within medical images [9]. This synergy between temporal understanding and attentive contextual analysis significantly boosts the model’s diagnostic precision, providing a more detailed and comprehensive understanding of patient-specific conditions [10].

Despite the notable advantages of SA and LSTM mechanism, these modules still encounter certain challenges, especially in clinical applications. Firstly, the conventional LSTM

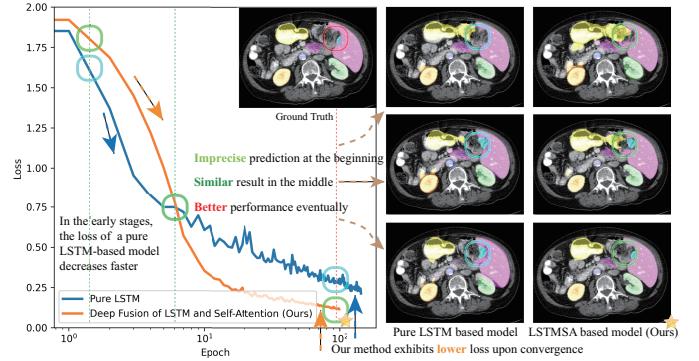


Fig. 1. The comparison of loss curves between LSTM based model and LSTMSA based model (Ours). Accompanying the curves are CT scan images demonstrating the ground truth and the progressive improvement of the deep fusion model’s segmentation accuracy, with the final image marked by a star symbolizing the lowest loss point.

module embedding often results in a rapid initial reduction of loss, which can lead to **premature convergence**. This issue not only impedes the model’s ability to retain long-term memory but also limits further accuracy improvements during later training stages [11]. In medical image segmentation, this premature convergence can severely compromise the model’s capacity to capture intricate anatomical details over extended sequences, thus undermining segmentation performance. The root of this issue lies in the model’s inefficient propagation of past information due to recurrent connections, coupled with a neglect of local details, which ultimately prevents the LSTM from fully realizing its potential. A balanced approach that fosters the synergy between SA and LSTM is essential for sustained learning progress. Fig. 1 illustrates the loss curves of standard LSTM network compared to our hybrid model, where our approach demonstrates a steadier and more uniform reduction in loss, indicative of enhanced learning stability. As training progresses, our method not only achieves a lower loss level in a more stable manner but also effectively integrates LSTM and SA strategies for improved temporal retention and long-range dependency capture. This stability is critical in medical image segmentation, where the ability to accurately segment across varying temporal and spatial contexts directly impacts clinical outcomes. Visual evidence from computed tomography (CT) scan images supports the quantitative improvements, showing enhanced segmentation accuracy that closely mirrors the ground truth over time.

Another significant challenge in medical image segmentation arises from the use of pure SA networks, particularly their susceptibility to a phenomenon where the network converges exponentially (with depth) to a rank-1 matrix, causing all tokens to become indistinguishable [12], a problem we

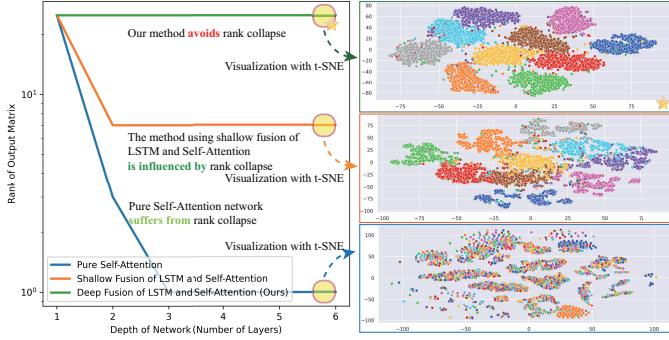


Fig. 2. Comparative analysis of the rank of the output matrix in neural networks of varying depths, contrasting the performance of pure SA, shallow fusion of LSTM and SA, and deep fusion of LSTM and SA (Ours). The three t-SNE visualizations on the right highlight the clustering of features at different layers.

refer to as **rank collapse**. In the context of medical image segmentation, this rank collapse can lead to a loss of critical spatial information, reducing the model's ability to distinguish between different tissue types or pathological features, thereby diminishing segmentation quality. This issue is driven by the SA mechanism functioning as an ensemble of shallow networks, which inherently limits matrix rank. To counteract this, advanced architectural designs are necessary to preserve high-dimensional features. Fig. 2 evaluates rank stability across network depths, demonstrating the effectiveness of our deep fusion method. Our approach consistently maintains a high rank, showcasing resilience against rank collapse, unlike pure SA networks, which exhibit a significant reduction in rank with increasing depth. t-SNE visualizations further validate our model's superior feature representation, leading to smoother loss transitions and enhanced segmentation accuracy by maintaining precise feature delineation.

To address the aforementioned challenges in medical image segmentation, we present a novel deep learning model known as SA-embedded LSTM (LSTMSA). Within the architecture of LSTMSA, we employ a deep embedding of the SA mechanism into the LSTM module, with the overarching goal of achieving a smoother training process and more efficient feature representation. This innovative approach involves integrating the SA mechanism at various time steps within the LSTM module, with the dual objective of maintaining a stronger emphasis on contextual information during the early stages of training, while harnessing the full spatiotemporal modeling capabilities inherent to LSTM throughout the entirety of the training process.

The contributions of this paper can be summarized as follows:

- Deep Fusion of LSTM and SA: Our first innovation lies in the deep integration of the LSTM and SA mechanism within a single module. This integration enhances the model's capability to capture both global and local contextual information effectively, allowing for superior medical image segmentation results.
- Resolving Premature Convergence and Rank Collapse Challenges: The second noteworthy innovation of this

work is the resolution of two prominent challenges in contextual models. Firstly, LSTMSA effectively addresses the issue of premature loss convergence, ensuring a more balanced and stable training process. Secondly, it tackles the problem of rank collapse inherent in pure SA networks, preserving high-dimensional feature representations, and averting information loss, particularly in complex medical image structures.

- State-of-the-art (SOTA) Performance on Four Public Datasets and Improvements on Twenty Baselines: Our third innovation manifests in the validation of LSTMSA's performance on diverse medical image datasets. Through extensive experimentation, LSTMSA consistently outperforms existing methodologies, achieving SOTA results in terms of segmentation accuracy and robustness. Moreover, LSTMSA demonstrates significant improvements over twenty baseline models across four widely used medical image datasets. These superior outcomes underscore the practical significance and clinical utility of our proposed model.

II. BACKGROUND

A. Literature Review

Separated LSTM and SA for Embedded Contextual Information in Medical Image Segmentation: Recent advancements in medical image segmentation have focused on embedding contextual information into Convolutional Neural Networks (CNNs), particularly within encoder-decoder architectures such as UNet [13]. These models capture both local semantic and global contextual features, improving segmentation of fine-grained structures [14]. Moreover, combining CNNs with RNNs, such as LSTM, has been effective for capturing temporal dependencies in 3D medical images, aiding in the segmentation of complex anatomical structures [4]. The inclusion of SA mechanism has further enriched these networks from the other side, enhancing their ability to capture contextual dependencies and improving performance across diverse medical imaging datasets [3]. Collectively, these advancements have significantly improved the accuracy and reliability of medical image segmentation.

Hybrid LSTM and SA Mechanism in Medical Image Segmentation: Several studies have explored hybrid approaches that combine LSTM and SA mechanisms for medical image segmentation. Typically, these methods use shallow fusion strategies, stacking the two modules where one's output serves as the other's input [15]. This approach, however, fails to fully leverage the complementary strengths of LSTM and SA. Alternatively, deeper integrations have been proposed, such as SwinLSTM [16], and RWKV [17], achieving more intricate fusions. Nevertheless, these models frequently fail to adequately tackle key challenges mentioned above, namely, premature convergence and rank collapse, and they do not deeply explore the individual contributions of LSTM and SA. Consequently, they may miss opportunities to optimize the interaction between these components, which in turn limits their effectiveness in medical image segmentation.

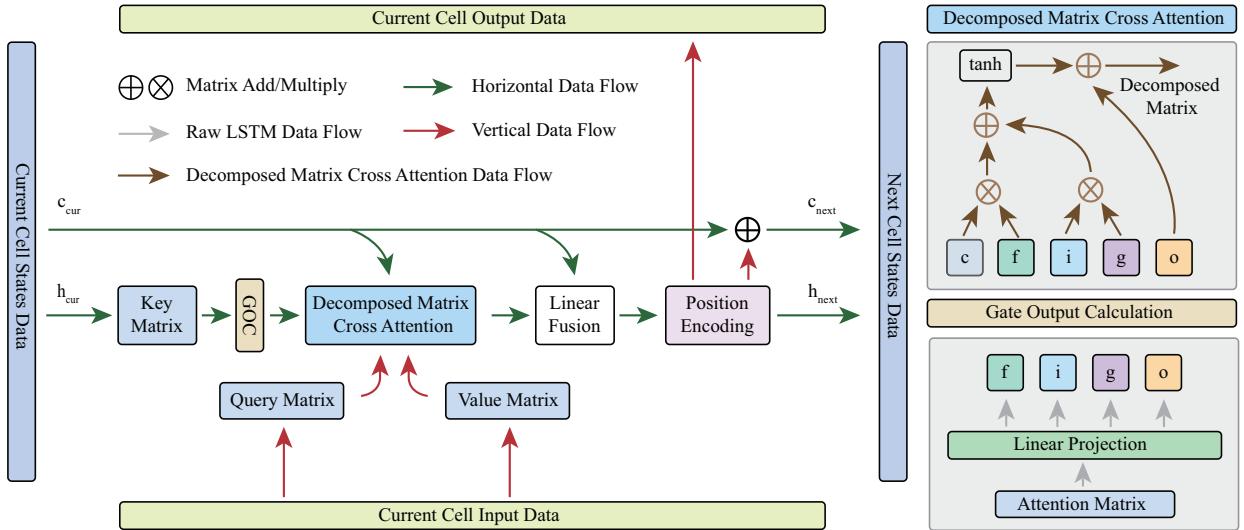


Fig. 3. Overview of LSTMSA module. The vertical data flow represents the processing over time in the input sequence, and the horizontal data flow represents the state updates within each cell.

B. Loss Convergence Analysis

Lemma: *Rapid loss reduction in early training indicates an aggressive weight adjustment strategy.*

Intuitive Proof: During early training, a sharp decline in loss suggests significant weight adjustments in response to gradients, quickly moving the network toward a favorable region in the loss landscape. However, this aggressive adjustment can diminish the network's ability to preserve long-term dependencies [11]. The rapid weight decay prioritizes recent data, undermining the integration and retention of older information. Thus, while fast initial learning improves early performance, it compromises long-term memory, highlighting a critical balance between learning speed and historical data retention in neural network training.

C. Rank Collapse Analysis

Pure SA mechanism can lead to token uniformity, reducing the network's ability to differentiate the importance of tokens in sequence processing, thus lowering the rank of the output matrix [12]. This uniformity problem arises from the inherent decomposition and induction bias of SA mechanism. In medical image segmentation, this issue can cause the model to inadequately attend to specific image regions. By assigning equal importance to all tokens, the network may overlook critical information, reducing its effectiveness in modeling the internal structure of sequences and affecting overall performance.

III. METHOD

A. Module Design Overview

As illustrated in Fig. 3, the LSTMSA module is designed to address two critical challenges: premature convergence in LSTM based models and rank collapse in pure SA mechanism. By deeply embedding SA into the LSTM architecture, LSTMSA effectively balances the strengths and mitigates the limitations of both components, enabling effective capture of

sequential temporal dependencies typical of LSTM and the intricate internal dependencies facilitated by SA for medical image segmentation. The module's design is guided by the following principles:

Firstly, to counteract rapid convergence in LSTM models, which often leads to the loss of effective long-term memory modeling early in training, LSTMSA integrates the SA matrices calculation into LSTM's basic units. This integration enhances the model's capacity to capture long-range dependencies by dynamically adjusting based on global sequence information, thereby mitigating premature convergence. SA allows the model to maintain a better understanding of long-term dependencies by leveraging global contextual information.

Secondly, LSTMSA addresses the rank collapse issue associated with pure SA networks, where SA tends to produce low-rank weight matrices, potentially compromising performance and generalization. To combat this, LSTMSA introduces a novel combination of query, key, and value matrices, decomposed and integrated into the LSTM unit states. This design enhances the model's ability to capture various aspects of the input sequence, thereby maintaining higher rank in the output matrix and improving the network's representational capacity.

B. Module Design Details

Alg. 1 outlines the calculation process of the LSTMSA module. The decision to use the hidden states as the key matrix is based on the necessity of incorporating contextual information that spans across different time steps (patches or slices in medical image). In medical imaging, such temporal and historical context is crucial for accurately segmenting structures that evolve or appear differently over time. This design is supported by our experimental analysis of time-step dependency IV-D, which demonstrates that the key matrix, derived from LSTM's hidden states, provides the most comprehensive representation of this temporal context. On the other hand, the input data, while serving as query and value matrices, allow the model to integrate immediate information

Algorithm 1 LSTM SA Calculation Process

Require: Current cell input x_{cur} , current hidden state h_{cur} , weight matrices W_i, W_f, W_o, W_g , bias vectors b_i, b_f, b_o, b_g , sequence length l , model dimensionality d_{model}

Ensure: Next cell state c_{next} , next hidden state h_{next}

- 1: **Step 1: Gate Output Calculation**
- 2: $xh \leftarrow [x_{cur}; h_{cur}]$
- 3: $i \leftarrow \sigma(W_i \odot xh + b_i)$
- 4: $f \leftarrow \sigma(W_f \odot xh + b_f)$
- 5: $o \leftarrow \sigma(W_o \odot xh + b_o)$
- 6: $g \leftarrow \tanh(W_g \odot xh + b_g)$
- 7: **Step 2: Decomposed Matrix Attention**
- 8: $D \leftarrow o + \tanh(c \cdot f + i \cdot g)$
- 9: **Step 3: Cross Attention**
- 10: $CA \leftarrow \text{Softmax} \left(\frac{D \cdot h_{cur}^T}{\sqrt{d_{h_{cur}}}} \right) \cdot D$
- 11: **Step 4: Feature Fusion**
- 12: $x \leftarrow [c_{cur}; CA]$
- 13: $F \leftarrow \text{ReLU}(W(x))$
- 14: **Step 5: Dynamic Positional Encoding**
- 15: $P_{dyn} \leftarrow \text{Variable}(\text{shape} = [l, d_{model}], \text{trainable} = True)$
- 16: $PF_{dyn} \leftarrow F + P_{dyn}$
- 17: **Step 6: Computing the Next Cell State**
- 18: $c_{next} \leftarrow c_{cur} + PF_{dyn}$
- 19: $h_{next} \leftarrow PF_{dyn}$
- 20: **return** c_{next}, h_{next}

with historical context during decision-making. This balance between using historical context (as keys) and immediate input data (as queries and values) is further validated by our ablation study IV-E1, which show that this configuration optimally enhances the model's performance in medical image segmentation.

By decomposing and reintegrating query, key, and value matrices within the LSTM framework, the model explores input data nuances more effectively, significantly boosting its ability to represent complex dependencies. The cross attention mechanism excels at uncovering dependencies across input segments, leading to a more granular understanding of sequence relationships. The deep fusion of LSTM states with cross attention outcomes preserves data integrity, effectively addressing rank collapse issues common in pure SA networks. Additionally, the incorporation of dynamic positional encoding differentiates sequence element positions, enhancing the model's understanding of order and structure within sequences, making LSTM SA a robust solution for complex, sequential data tasks.

IV. EXPERIMENTS AND RESULTS

A. Implementation Details

For the assessment of our proposed modules on 2D inputs, we conduct experiments on the following datasets:

- Synapse dataset: This dataset encompasses various organ segmentation tasks, including eight distinct abdominal organs: the aorta, gallbladder (GB), spleen (SP), left

kidney (KL), right kidney (KR), liver, pancreas (PC), and stomach (SM)¹.

- ISIC2018 dataset: Focusing on skin lesion segmentation, this dataset has been employed in previous study by Codella *et al.* [18].

In scenarios involving 3D inputs, we perform experiments on the following datasets:

- ACDC dataset: Dataset for heart segmentation, including the right ventricle (RV), left ventricle (LV), and myocardium (Myo)².
- CVC-ClinicDB dataset: Addressing polyp segmentation in colonoscopy videos, this dataset has been used for benchmarking segmentation methods in the work by Bernal *et al.* [19].

We assess the performance of our proposed modules on these datasets using the DICE score and Hausdorff distance 95% (HD95) as evaluation metrics. The DICE score quantifies the overlap between predicted and ground truth masks, defined as $\text{DICE} = \frac{2 \times |X \cap Y|}{|X| + |Y|}$, where X and Y are the predicted and ground truth masks. The HD95 measures the 95th percentile of the maximum distance between the two masks, offering a comprehensive evaluation of segmentation quality and edge accuracy.

Across all the aforementioned datasets, our approach integrates the U-shaped network with our proposed modules as a baseline model. Importantly, we introduce the proposed modules into the baseline models' architecture without any additional alterations to the original network structure and employ the same training strategies as the baseline models, indicating the seamless integration of our proposed modules into existing network structures for enhanced model performance, without necessitating extra training pipeline configuration or hyperparameter adjustments.

B. Results on 2D Input

The data in Table I shows that networks incorporating the LSTM SA module consistently achieve SOTA results across both 2D datasets. For instance, in the ISIC2018 dataset, LSTM SA achieves a precision score of 91.42, significantly outperforming other baseline methods, indicating superior segmentation performance on 2D inputs. This performance improvement is linked to the LSTM SA's impact on loss attenuation (Fig. 1) and rank preservation. The module stabilizes loss reduction during early training, reducing the risk of rapid loss decline, which facilitates better convergence. Moreover, the cross attention mechanism in LSTM SA accurately captures key image features, reducing rank collapse and enhancing overall segmentation precision.

We also visualize network performance with the LSTM SA module, as depicted in Fig. 4 and Fig. 5. These figures compare LSTM SA with other baseline methods in organ segmentation tasks, showing that LSTM SA produces notably precise results, particularly in small organ segmentation. The segmentation boundaries generated by LSTM SA are clearer

¹<https://www.synapse.org/#/Synapse:syn3193805/wiki/217789>

²<https://www.creatis.insa-lyon.fr/Challenge/acdc/>

TABLE I

QUANTITATIVE RESULTS ON SYNAPSE MULTI-ORGAN DATASET AND ISIC2018 SEGMENTATION DATASET. DICE SCORES (%), HD95 ARE REPORTED. THE BEST RESULTS ARE BOLDED. THE SECOND BEST RESULTS ARE UNDERLINED. \uparrow DENOTES HIGHER VALUE INDICATING BETTER PERFORMANCE, \downarrow DENOTES LOWER VALUE INDICATING BETTER PERFORMANCE. BASELINES ARE STARRED.

| Method | Synapse | | | | | | | | | | ISIC2018 | | |
|-------------------|-----------------|-------------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|--------------|----------------|-----------------|-------------------|
| | DICE \uparrow | HD95 \downarrow | Aorta | GB | KL | KR | Liver | PC | SP | SM | Method | DICE \uparrow | HD95 \downarrow |
| UNet [13] | 70.11 | 44.69 | 84.00 | 56.70 | 72.41 | 62.64 | 86.98 | 48.73 | 81.48 | 67.96 | UNet [13] | 87.41 | 4.03 |
| TransUNet [1] | 77.48 | 31.69 | 87.23 | 63.13 | 81.87 | 77.02 | 94.08 | 55.86 | 85.08 | 75.62 | DWUNet [20] | 87.47 | 4.55 |
| MT-UNet [2] | 78.59 | 26.59 | 87.92 | 64.99 | 81.47 | 77.29 | 93.06 | 59.46 | 87.75 | 76.81 | ResUNet [21] | 87.91 | 3.49 |
| SwinUNet [22] | 79.13 | 21.55 | 85.47 | 66.53 | 83.28 | 79.61 | 94.29 | 56.58 | 90.66 | 76.60 | UNet++ [14] | 88.32 | 3.83 |
| MISSFormer [3] | 81.96 | 18.20 | 86.99 | 68.65 | 85.21 | 82.00 | 94.41 | 65.67 | 91.92 | 80.81 | R2UNet [23] | 90.13 | 3.62 |
| TransCASCADE [24] | 82.68 | 17.34 | 86.63 | 68.48 | 87.66 | 84.56 | 94.43 | 65.33 | 90.79 | 83.52 | DCSAU-Net [7]* | 90.41 | 2.21 |
| MERIT [8]* | <u>84.22</u> | <u>16.51</u> | <u>88.38</u> | <u>73.48</u> | <u>87.21</u> | <u>84.31</u> | <u>95.06</u> | 69.97 | <u>91.21</u> | <u>84.15</u> | MSCA-Net [25] | <u>90.52</u> | 2.79 |
| LSTMSA (Ours) | 84.85 | <u>15.83</u> | 89.51 | <u>73.32</u> | 85.88 | 84.88 | 95.44 | <u>69.20</u> | 90.99 | 84.49 | LSTMSA (Ours) | 91.42 | <u>2.53</u> |

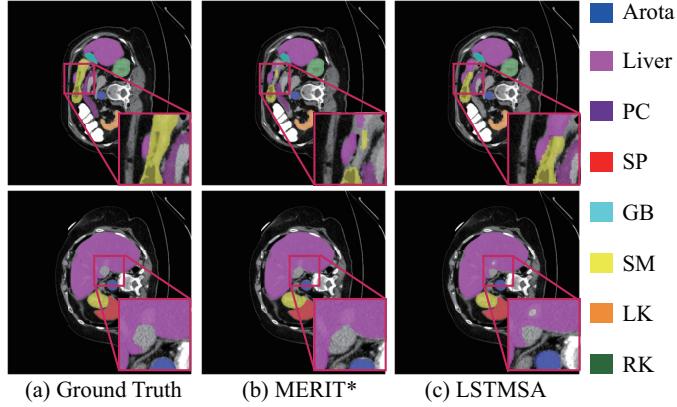


Fig. 4. 2D visualization results on the Synapse dataset. (a) Ground truths, (b) MERIT* (baseline are starred), (c) LSTMSA (Ours). The red rectangular box indicates the zoomed-in region.

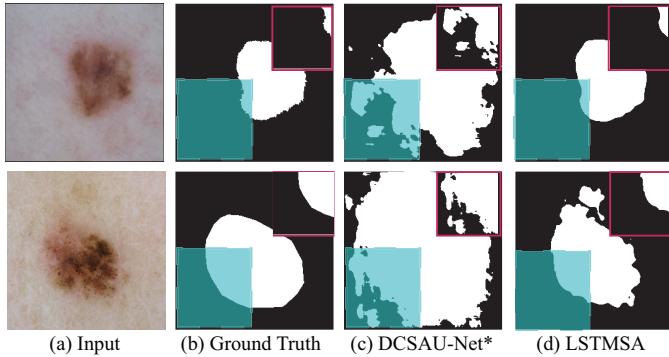


Fig. 5. 2D visualization results on the ISIC2018 dataset. (a) Input, (b) Ground Truth, (c) DCSAU-Net* (baseline are starred), (d) LSTMSA (Ours). The red rectangular box indicates the zoomed-in region.

and more accurate, effectively eliminating blurriness and imprecision. Detailed analysis reveals that the LSTMSA module exhibits heightened attention accuracy around small organ edges, capturing even subtle features within the structures.

C. Results on 3D input

Our analysis of 3D medical image segmentation, as shown in Table II, demonstrates that the integration of the LSTMSA module significantly improves performance over SOTA methods, particularly those relying solely on attention mechanisms. For example, on the ACDC dataset, the LSTMSA-enhanced

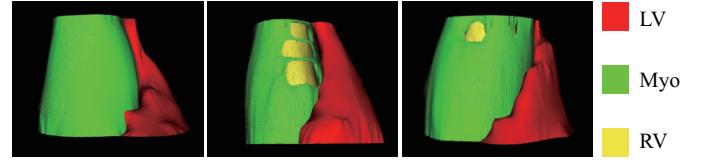


Fig. 6. 3D visualization results on the ACDC dataset. ED is End Diastolic, and ES is End Systolic.

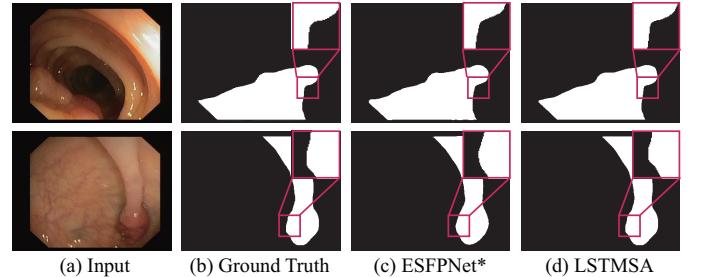


Fig. 7. 2D visualization results on the CVC-ClinicDB dataset. (a) Input, (b) Ground Truth, (c) ESFPNet* (baseline are starred), (d) LSTMSA. The red rectangular box indicates the zoomed-in region.

model achieves a DICE score of 92.86, surpassing the MT-UNet baseline by nearly 3%, and delivers high precision in segmenting various cardiac structures, with scores of 91.54, 90.77, and 96.27 for RV, Myo, and LV, respectively. These results highlight LSTMSA's effectiveness in capturing complex features and offering nuanced improvements.

On the CVC-ClinicDB dataset, the LSTMSA model achieves a DICE score of 96.11, outperforming advanced models like ESFPNet and TGANet. This performance underscores LSTMSA's capability in handling intricate segmentation tasks, pushing beyond current SOTA levels in 3D segmentation. The model's edge lies in its integration of SA with sequential modeling through LSTM, which enhances the network's ability to capture both spatial and temporal dependencies. This combination refines segmentation accuracy and generalization, particularly visible in the sharper and more accurate boundaries in Fig. 6 and Fig. 7.

D. Validation of Module Structure Design

We conduct a time-step dependency analysis to validate that using the LSTM states combination as the key matrix in cross attention, rather than as the query and value matrices, better manages both long-term and short-term memory. This

TABLE II

QUANTITATIVE RESULTS ON ACDC AND CVC-CLINICDB DATASET. DICE SCORES (%) AND HD95 ARE REPORTED. THE BEST RESULTS ARE BOLDED. THE SECOND BEST RESULTS ARE UNDERLINED. \uparrow DENOTES HIGHER VALUE INDICATING BETTER PERFORMANCE, \downarrow DENOTES LOWER VALUE INDICATING BETTER PERFORMANCE. BASELINES ARE STARRED.

| Method | ACDC | | | | | CVC-ClinicDB | | |
|-------------------|-----------------|-------------------|--------------|--------------|--------------|-------------------|-----------------|-------------------|
| | DICE \uparrow | HD95 \downarrow | RV | Myo | LV | Method | DICE \uparrow | HD95 \downarrow |
| TransUNet [1] | 89.71 | 2.54 | 88.86 | 84.53 | 95.73 | ColonSegNet [26] | 88.62 | 4.56 |
| SwinUNet [22] | 90.00 | 4.52 | 88.55 | 85.62 | 95.83 | FCBFormer [9] | 92.53 | 3.21 |
| MT-UNet [2]* | 90.43 | 2.23 | 86.64 | 89.04 | 95.62 | SSFormer-S [5] | 92.68 | 1.45 |
| MISSFormer [3] | 90.86 | 2.13 | 89.55 | 88.04 | 94.99 | HarDNet-DFUS [27] | 93.32 | 1.29 |
| PVT-CASCADE [24] | 91.46 | 1.09 | 88.9 | 89.97 | 95.50 | FANet [28] | 93.55 | 1.15 |
| TransCASCADE [24] | 91.63 | 1.09 | 89.14 | 90.25 | 95.50 | TGANet [29] | 94.57 | 1.47 |
| MERIT [8] | 92.32 | <u>1.08</u> | 90.87 | 90.00 | <u>96.08</u> | SSFormer-L [5] | 94.72 | <u>0.73</u> |
| FCT [30] | 92.84 | 5.29 | 92.02 | 90.61 | 95.89 | ESFPNet [10]* | 94.90 | 1.21 |
| LSTMSA (Ours) | 92.86 | 1.07 | 91.54 | 90.77 | 96.27 | LSTMSA (Ours) | 96.11 | 0.53 |

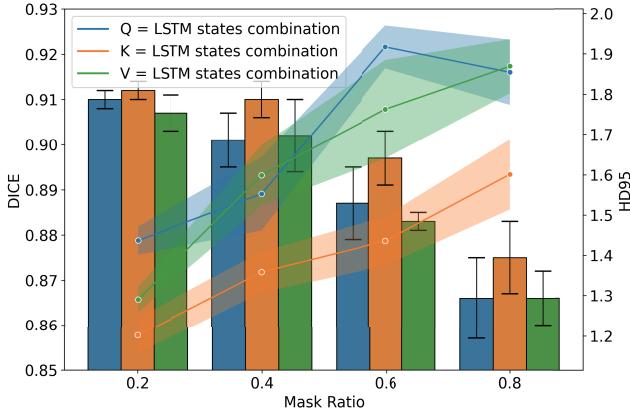


Fig. 8. Time-step dependency analysis on ACDC dataset. The x-axis represents the percentage of slices masked, and the y-axis represents the performance in terms of DICE and HD95.

analysis involves selectively masking time steps in the input sequence and evaluating the resultant performance changes. Accordingly, we configured three distinct variants, using the LSTM states combination as the query, key, and value matrices respectively.

Fig. 8 shows the DICE and HD95 performance on the ACDC dataset under different masking proportions (0.2, 0.4, 0.6, 0.8). Results, averaged over three experiments, reveal that using the LSTM states combination as the key matrix results in the least performance degradation. This supports the design choice and indicates that long-term and short-term memory is more effectively managed when encoded in the key matrix.

E. Abalation Study and Efficiency Analysis

1) *Ablation Study:* The ablation study on the ACDC dataset, as seen in Table III, examines the individual components of the LSTMSA module and their impact on segmentation. Removing either the SA or LSTM components leads to a significant drop in performance, highlighting their combined importance. The study also reveals that excluding the cross attention mechanism (including QKV selection, only using input data and only using hidden states) or dynamic positional encoding reduces segmentation accuracy, emphasizing their critical role in the LSTMSA module's performance.

TABLE III

ABLATION STUDY ON ACDC DATASET. DICE SCORES (%) AND HD95 ARE REPORTED. \checkmark DONATES USED, AND \times DONATES NOT USED. THE BEST RESULTS ARE BOLDED. THE SECOND BEST RESULTS ARE UNDERLINED. SA(X,X,X) DENOTES USING ALL INPUT DATA FOR QKV SELECTION, SA(H,H,H) DENOTES ALL USING HIDDEN STATES FOR QKV SELECTION.

| Method | SA | LSTM | DICE \uparrow | HD95 \downarrow |
|--------------------------------|--------------|--------------|-----------------|-------------------|
| MT-UNet [2] (Baseline) | \times | \times | 90.43 | 2.23 |
| Only SA | \checkmark | \times | 91.17 | 1.77 |
| Only LSTM | \times | \checkmark | 91.69 | 1.76 |
| No Cross Attention (SA(x,x,x)) | \checkmark | \checkmark | 92.43 | 1.22 |
| No Cross Attention (SA(h,h,h)) | \checkmark | \checkmark | 92.34 | 1.17 |
| No Dynamic Positional Encoding | \checkmark | \checkmark | 92.38 | 1.18 |
| LSTMSA (Ours) | \checkmark | \checkmark | 92.86 | 1.07 |

TABLE IV

EFFICIENCY COMPARISONS BETWEEN OUR PROPOSED MODULE AND SEPARATED SA AND LSTM WITHIN SOTA METHOD ON EACH DATASET.

THE RESULTS ARE OBTAINED BY AVERAGING THE OUTCOMES OF 10 EXPERIMENT RUNS AND THE RATE OF INCREASE IS ALSO CALCULATED. PARAMS DENOTES THE NUMBER OF PARAMETERS, FLOPS DENOTES THE NUMBER OF FLOATING-POINT OPERATIONS, AND TIME DENOTES THE INFERENCE TIME TAKEN FOR GPU INFERENCE.

| Matrices | Dataset | Baseline | Separated δ | LSTMSA δ |
|------------|--------------|----------|--------------------|-----------------|
| Params (M) | Synapse | 146.51 | +0.75% | +0.56% |
| | ISIC2018 | 2.59 | +4.63% | +4.15% |
| | ACDC | 49.31 | +0.41% | +1.65% |
| | CVC-ClinicDB | 3.53 | +3.11% | +8.97% |
| Flops (G) | Synapse | 28.32 | +3.95% | +0.28% |
| | ISIC2018 | 15.88 | +6.92% | +1.50% |
| | ACDC | 14.35 | +5.57% | +0.56% |
| | CVC-ClinicDB | 0.48 | +41.67% | +34.29% |
| Time (ms) | Synapse | 45.41 | +2.42% | +1.44% |
| | ISIC2018 | 10.95 | +9.31% | +5.51% |
| | ACDC | 33.91 | +3.83% | +0.56% |
| | CVC-ClinicDB | 5.19 | +13.29% | +10.05% |

2) *Efficiency Analysis:* Table IV presents an efficiency analysis on the Synapse, ISIC2018, ACDC, and CVC-ClinicDB datasets, comparing the computational costs of our LSTMSA module and separated SA and LSTM components within baseline models, using a tensor with shape of (1, 3, 256, 256). While the LSTMSA module shows a slight increase in parameters, FLOPs and GPU time, it does not introduce significant overhead compared to using SA and LSTM separately, making it an efficient solution for medical image segmentation.

3) *Input Perturbation Analysis:* Table V shows the input perturbation analysis results on the ACDC datasets, assessing the LSTMSA module's robustness to noise. The module

TABLE V

INPUT PERTURBATION ANALYSIS ON ACDC DATASET. DICE SCORES (%) AND HD95 ARE REPORTED. LOW AND MODERATE NOISE IS ADDED USING GAUSSIAN NOISE WITH $\sigma = 0.1$ AND $\sigma = 0.2$, RESPECTIVELY. STRONG NOISE IS ADDED USING RANDOM CONTRAST, SATURATION, HUE AND BRIGHTNESS TRANSFORMATIONS.

| Method | Noise level | DICE δ | HD95 δ |
|---------------|-------------|---------------|---------------|
| MT-UNet [2] | low | -4.22 | +0.81 |
| LSTMSA (Ours) | low | -2.63 | +0.29 |
| MT-UNet [2] | moderate | -9.41 | +2.22 |
| LSTMSA (Ours) | moderate | -6.98 | +0.90 |
| MT-UNet [2] | strong | -77.89 | +27.22 |
| LSTMSA (Ours) | strong | -72.45 | +20.11 |

TABLE VI

BASELINE COMPARISON AND STATISTICAL SIGNIFICANCE TESTING RESULTS ACROSS DATASETS AND BASELINES. DICE SCORES (%) ARE REPORTED. \uparrow INDICATES AN INCREASE IN DICE OF LESS THAN 1% AND $\uparrow\uparrow$ INDICATES AN INCREASE IN DICE OF GREATER THAN 1%. A P-VALUE LESS THAN 0.05 INDICATES STATISTICAL SIGNIFICANCE.

| Dataset | Method | DICE δ | Significance |
|--------------|-------------------|--------------------------|---|
| Synapse | TransUNet [1] | 78.89 $\uparrow\uparrow$ | t-statistic: 3.3665 p-value: 0.0281 Significant: Yes |
| | MT-UNet [2] | 78.41 \uparrow | |
| | MISSFormer [3] | 80.49 \uparrow | |
| | DAEFormer [6] | 82.52 \uparrow | |
| | MERIT [8] | 84.85 \uparrow | |
| ISIC2018 | UNet [13] | 88.37 \uparrow | t-statistic: 19.5116 p-value: 0.0004 Significant: Yes |
| | DWUNet [20] | 88.56 \uparrow | |
| | ResUNet [21] | 89.01 \uparrow | |
| | UNet++ [14] | 89.14 \uparrow | |
| | DCSAU-Net [7] | 91.42 $\uparrow\uparrow$ | |
| ACDC | MISSFormer [3] | 90.75 \uparrow | t-statistic: 3.7350 p-value: 0.0202 Significant: Yes |
| | TransUNet [1] | 91.64 $\uparrow\uparrow$ | |
| | DAEFormer [6] | 91.82 \uparrow | |
| | MT-UNet [2] | 92.86 \uparrow | |
| | MERIT [8] | 92.45 \uparrow | |
| CVC-ClinicDB | FCBFormer [9] | 95.63 \uparrow | t-statistic: 6.0501 p-value: 0.0037 Significant: Yes |
| | SSFormer-S [5] | 94.10 \uparrow | |
| | HarDNet-DFUS [27] | 95.65 \uparrow | |
| | FANet [28] | 95.88 \uparrow | |
| | ESFPNet [10] | 96.11 \uparrow | |

demonstrates superior robustness compared to baseline models, with minimal performance decline under varying noise levels, indicating its ability to maintain segmentation accuracy in the presence of noise.

4) *Effect on Different Baselines:* To assess the LSTMSA module's universality, Table VI compares performance across different baseline networks with and without the LSTMSA module. It is worth noting that for some models without publicly available code, we reimplement the corresponding models according to the respective papers. The results demonstrate that our module consistently improves segmentation accuracy across datasets, confirming its versatility. Moreover, we adopt significance testing to evaluate the improvements in segmentation performance achieved by the LSTMSA module across different datasets and baseline models. The paired t-test is chosen due to the nature of the experimental setup, where each dataset is evaluated under two different conditions before and after the application of the proposed module. The results show that the improvements in segmentation performance are statistically significant, indicating the robustness and reliability of the LSTMSA module.

V. CONCLUSION

In this paper, we present LSTMSA, a novel deep learning model that integrates LSTM with SA to improve medical image segmentation. LSTMSA effectively balances global and local context modeling, addressing issues about premature loss convergence and rank collapse found in pure SA networks. Extensive evaluations show that LSTMSA consistently outperforms existing methods across various datasets and baselines.

However, the current sequential fusion of LSTM and SA limits training speed. Future work will explore parallel training structures to optimize both efficiency and performance, further enhancing LSTMSA's practicality in medical image segmentation.

REFERENCES

- J. Chen, Y. Lu, Q. Yu, X. Luo, E. Adeli, Y. Wang, L. Lu, A. L. Yuille, and Y. Zhou, "Transunet: Transformers make strong encoders for medical image segmentation," *arXiv preprint arXiv:2102.04306*, 2021.
- H. Wang, S. Xie, L. Lin, Y. Iwamoto, X.-H. Han, Y.-W. Chen, and R. Tong, "Mixed transformer u-net for medical image segmentation," in *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2022, pp. 2390–2394.
- X. Huang, Z. Deng, D. Li, X. Yuan, and Y. Fu, "Missformer: An effective transformer for 2d medical image segmentation," *IEEE Transactions on Medical Imaging*, vol. 42, no. 5, pp. 1484–1494, 2023.
- J.-H. Chang, K.-H. Lin, T.-H. Wang, Y.-K. Zhou, and P.-C. Chung, "Image segmentation in 3d brachytherapy using convolutional lstm," *Journal of Medical and Biological Engineering*, vol. 41, pp. 636–651, 2021.
- J. Wang, Q. Huang, F. Tang, J. Meng, J. Su, and S. Song, "Stepwise feature fusion: Local guides global," in *Medical Image Computing and Computer Assisted Intervention-MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part III*. Springer, 2022, pp. 110–120.
- R. Azad, R. Arimond, E. K. Aghdam, A. Kazerouni, and D. Merhof, "Dae-former: Dual attention-guided efficient transformer for medical image segmentation," *arXiv preprint arXiv:2212.13504*, 2022.
- Q. Xu, Z. Ma, H. Na, and W. Duan, "Dcsau-net: A deeper and more compact split-attention u-net for medical image segmentation," *Computers in Biology and Medicine*, vol. 154, p. 106626, 2023.
- M. M. Rahman and R. Marculescu, "Multi-scale hierarchical vision transformer with cascaded attention decoding for medical image segmentation," *arXiv preprint arXiv:2303.16892*, 2023.
- E. Sanderson and B. J. Matuszewski, "Fcn-transformer feature fusion for polyp segmentation," in *Medical Image Understanding and Analysis: 26th Annual Conference, MIUA 2022, Cambridge, UK, July 27–29, 2022, Proceedings*. Springer, 2022, pp. 892–907.
- Q. Chang, D. Ahmad, J. Toth, R. Bascom, and W. E. Higgins, "EsfpNet: efficient deep learning architecture for real-time lesion segmentation in autofluorescence bronchoscopy video," *arXiv preprint arXiv:2207.07759*, 2022.
- J. Zhao, F. Huang, J. Lv, Y. Duan, Z. Qin, G. Li, and G. Tian, "Do rnn and lstm have long memory?" in *International Conference on Machine Learning*. PMLR, 2020, pp. 11365–11375.
- Y. Dong, J.-B. Cordonnier, and A. Loukas, "Attention is not all you need: Pure attention loses rank doubly exponentially with depth," in *International Conference on Machine Learning*. PMLR, 2021, pp. 2793–2803.
- O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Medical Image Computing and Computer-Assisted Intervention-MICCAI 2015: 18th International Conference, Munich, Germany, October 5–9, 2015, Proceedings, Part III*. Springer, 2015, pp. 234–241.
- Z. Zhou, M. M. Rahman Siddiquee, N. Tajbakhsh, and J. Liang, "Unet++: A nested u-net architecture for medical image segmentation," in *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: 4th International Workshop, DLMIA 2018, and 8th International Workshop, ML-CDS 2018, Held in Conjunction with MICCAI 2018, Granada, Spain, September 20, 2018, Proceedings 4*. Springer, 2018, pp. 3–11.

- [15] H. Zhang, Z. Wang, and H. Vallery, "Learning-based nlos detection and uncertainty prediction of gnss observations with transformer-enhanced lstm network," in *2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC)*. IEEE, 2023, pp. 910–917.
- [16] S. Tang, C. Li, P. Zhang, and R. Tang, "Swinlstm: Improving spatiotemporal prediction accuracy using swin transformer and lstm," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2023, pp. 13 470–13 479.
- [17] B. Peng, E. Alcaide, Q. Anthony, A. Albalak, S. Arcadinho, H. Cao, X. Cheng, M. Chung, M. Grella, K. K. GV *et al.*, "Rwkv: Reinventing rnns for the transformer era," *arXiv preprint arXiv:2305.13048*, 2023.
- [18] N. C. Codella, D. Gutman, M. E. Celebi, B. Helba, M. A. Marchetti, S. W. Dusza, A. Kalloo, K. Liopyris, N. Mishra, H. Kittler *et al.*, "Skin lesion analysis toward melanoma detection: A challenge at the 2017 international symposium on biomedical imaging (isbi), hosted by the international skin imaging collaboration (isic)," in *2018 IEEE 15th international symposium on biomedical imaging (ISBI 2018)*. IEEE, 2018, pp. 168–172.
- [19] J. Bernal, F. J. Sánchez, G. Fernández-Esparrach, D. Gil, C. Rodríguez, and F. Vilariño, "Wm-dova maps for accurate polyp highlighting in colonoscopy: Validation vs. saliency maps from physicians," *Computerized medical imaging and graphics*, vol. 43, pp. 99–111, 2015.
- [20] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
- [21] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [22] H. Cao, Y. Wang, J. Chen, D. Jiang, X. Zhang, Q. Tian, and M. Wang, "Swin-unet: Unet-like pure transformer for medical image segmentation," in *Computer Vision–ECCV 2022 Workshops: Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part III*. Springer, 2023, pp. 205–218.
- [23] M. Z. Alom, M. Hasan, C. Yakopcic, T. M. Taha, and V. K. Asari, "Recurrent residual convolutional neural network based on u-net (r2u-net) for medical image segmentation," *arXiv preprint arXiv:1802.06955*, 2018.
- [24] M. M. Rahman and R. Marculescu, "Medical image segmentation via cascaded attention decoding," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 6222–6231.
- [25] Y. Sun, D. Dai, Q. Zhang, Y. Wang, S. Xu, and C. Lian, "Msca-net: Multi-scale contextual attention network for skin lesion segmentation," *Pattern Recognition*, vol. 139, p. 109524, 2023.
- [26] D. Jha, S. Ali, N. K. Tomar, H. D. Johansen, D. Johansen, J. Rittscher, M. A. Riegler, and P. Halvorsen, "Real-time polyp detection, localization and segmentation in colonoscopy using deep learning," *Ieee Access*, vol. 9, pp. 40 496–40 510, 2021.
- [27] T.-Y. Liao, C.-H. Yang, Y.-W. Lo, K.-Y. Lai, P.-H. Shen, and Y.-L. Lin, "Hardnet-dfus: An enhanced harmonically-connected network for diabetic foot ulcer image segmentation and colonoscopy polyp segmentation," *arXiv preprint arXiv:2209.07313*, 2022.
- [28] N. K. Tomar, D. Jha, M. A. Riegler, H. D. Johansen, D. Johansen, J. Rittscher, P. Halvorsen, and S. Ali, "Fanet: A feedback attention network for improved biomedical image segmentation," *IEEE Transactions on Neural Networks and Learning Systems*, 2022.
- [29] N. K. Tomar, D. Jha, U. Bagci, and S. Ali, "Tganet: text-guided attention for improved polyp segmentation," in *Medical Image Computing and Computer Assisted Intervention–MICCAI 2022: 25th International Conference, Singapore, September 18–22, 2022, Proceedings, Part III*. Springer, 2022, pp. 151–160.
- [30] A. Tragakis, C. Kaul, R. Murray-Smith, and D. Husmeier, "The fully convolutional transformer for medical image segmentation," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, 2023, pp. 3660–3669.