

GDKVM: Echocardiography Video Segmentation via Spatiotemporal Key-Value Memory with Gated Delta Rule

Rui Wang¹ Yimu Sun¹ Jingxing Guo¹ Huisi Wu¹ * Jing Qin²

¹College of Computer Science and Software Engineering, Shenzhen University

²Centre for Smart Health, School of Nursing, The Hong Kong Polytechnic University

2400101058@mails.szu.edu.cn, hswu@szu.edu.cn

Abstract

Accurate segmentation of cardiac chambers in echocardiography sequences is crucial for the quantitative analysis of cardiac function, aiding in clinical diagnosis and treatment. The imaging noise, artifacts, and the deformation and motion of the heart pose challenges to segmentation algorithms. While existing methods based on convolutional neural networks, Transformers and space-time memory networks, have improved segmentation accuracy, they often struggle with the trade-off between capturing long-range spatiotemporal dependencies and maintaining computational efficiency with fine-grained feature representation. In this paper, we introduce GDKVM, a novel architecture for echocardiography video segmentation. The model employs Linear Key-Value Association (LKVA) to effectively model inter-frame correlations, and introduces Gated Delta Rule (GDR) to efficiently store intermediate memory states. Key-Pixel Feature Fusion (KPFF) module is designed to integrate local and global features at multiple scales, enhancing robustness against boundary blurring and noise interference. We validated GDKVM on two mainstream echocardiography video datasets (CAMUS and EchoNet-Dynamic) and compared it with various state-of-the-art methods. Experimental results show that GDKVM outperforms existing approaches in terms of segmentation accuracy and robustness, while ensuring real-time performance. Codes are available at <https://github.com/wangrui2025/GDKVM>.

1. Introduction

Echocardiography is a vital imaging modality for cardiac assessment, valued for its non-invasive nature and real-time imaging capabilities. Accurate segmentation of cardiac chambers from echocardiography videos is a foundational step for quantitative analysis [11]. Delineating the left ven-

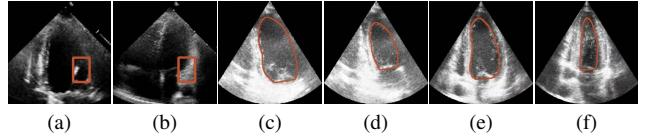


Figure 1. Illustrative challenges for echocardiography video segmentation: (a) speckle noise, (b) indistinct or blurred contours, and (c-f) the substantial changes in the target’s shape and scale throughout the cardiac cycle.

tricle is crucial for calculating the left ventricular ejection fraction, a cornerstone metric for diagnosing cardiovascular diseases and guiding treatment strategies [1, 43].

However, achieving precise segmentation in echocardiography faces severe challenges stemming from poor image quality and complex cardiac dynamics. Figures 1a and 1b show that ultrasound images are characterized by high speckle noise and low contrast, which obscure tissue structures and lead to weak or incomplete anatomical boundaries. These artifacts hinder the model’s ability to learn robust features and can result in the inaccurate segmentation. In the temporal dimension, the heart undergoes significant non-rigid deformation throughout the cardiac cycle. Figures 1c to 1f show that, the shape and scale of the LV change dramatically between systole and diastole. This substantial dynamic variation demands that segmentation models possess strong temporal modeling capabilities to accurately track the changing appearance of the target across the video sequence [20].

Convolutional Neural Networks [29], particularly U-Net architectures [14, 24, 33], marked a significant milestone in echocardiography segmentation. Their proficiency in learning hierarchical spatial features improved delineation robustness against noise [15]. Yet, their inherently local receptive fields limit performance where boundaries are ambiguous and preclude the direct modeling of inter-frame temporal dependencies [31]. Vision Transformers address the local-context limitation by providing a global field of

*Corresponding Author

view [21, 38], but their quadratic computational complexity and large data requirements present practical barriers. Spatiotemporal models incorporating recurrent units like ConvLSTM [4] aim for temporal coherence, yet they risk propagating errors through sequences and add significant computational overhead. A gap thus remains for a model that can integrate global temporal context efficiently while preserving local spatial precision.

This paper introduces **GDKVM**, a new architecture for echocardiography video segmentation designed to synthesize global temporal modeling with local feature precision in a computationally efficient manner. The core of GDKVM is a Linear Key-Value Association (LKVA) mechanism that captures inter-frame relationships. To efficiently update and preserve essential temporal information, we propose Gated Delta Rule (GDR) module that discards irrelevant historical memory while retaining features pertinent to the current frame. Key-Pixel Feature Fusion (KPFF) module enhances spatiotemporal representation by integrating local key features, global key features, and pixel-level features. We evaluate GDKVM on two widely adopted echocardiography video datasets, CAMUS [18] and EchoNet-Dynamic [30], and compare it with various task-specific state-of-the-art methods as well as recent STM models. Experimental results indicate that **GDKVM** achieves consistently higher segmentation accuracy than existing approaches, demonstrating its effectiveness and robustness on both datasets. Our main contributions are summarized as follows:

- We propose a novel model **GDKVM** for echocardiography video segmentation, which fully leverages the powerful representational capacity of linear key-value association with new adaptations to ultrasound videos.
- We introduce Gated Delta Rule (GDR) and Key-Pixel Feature Fusion (KPFF) module; GDR employs interactions between the current frame and memory bias, helping the model manage its memory. KPFF enables the model to acquire global spatiotemporal features.
- We conduct extensive experiments on the CAMUS and EchoNet-Dynamic datasets to validate the superior performance of our approach over current state-of-the-art methods.

2. Related Work

2.1. Echocardiography Video Segmentation

Echocardiography video segmentation aims for the temporally coherent segmentation of key cardiac structures, such as the left ventricle, which is crucial for the accurate assessment of cardiac functions like ejection fraction. Fully-supervised methods exhibit excellent performance, but they rely on dense, frame-by-frame annotation that is both time-consuming and costly. Sparsely-supervised methods, which require only a few annotated frames, have thus become

an active area of research, aiming to effectively balance segmentation accuracy with annotation overhead. Multi-frame aggregation methods [23, 35] and Space-Time Memory Networks (STM) [28] are widely used in video segmentation tasks to leverage temporal information. The former processes multiple frames simultaneously to improve segmentation performance in long sequences, but demands substantial GPU resources in large-scale scenarios. By storing the features and segmentation results of past frames in a memory bank and retrieving only the necessary semantic information when a new frame arrives, STM networks significantly reduce memory overhead while maintaining temporal consistency. Methods such as XMem [5] and XMem++ [3] introduce multi-level memory mechanisms, drawing on human cognitive processes to sustain effective segmentation in extremely long videos while controlling memory usage. Existing Space-Time Memory Networks [3, 5–7, 28] typically rely on mask propagation from fully annotated reference frames, yet they lack functionality to identify phases of the cardiac cycle. In medical applications, where it may be necessary to automatically locate end-systole and end-diastole frames, this limitation can reduce segmentation accuracy at critical moments. MemSAM [10] incorporates richer temporal cues into the memory bank, reducing misinterpretation of ambiguous structures and mitigating error propagation, but its computational overhead remains high, hindering direct deployment in clinical settings. Consequently, achieving both low computational cost and precise cardiac cycle recognition is a critical challenge that spatiotemporal memory methods urgently need to address in medical video segmentation.

2.2. Linear Attention Methods

Traditional self-attention [37] exhibits quadratic growth in computational complexity as sequence length increases, leading to enormous resource overhead. To address this issue, Linear Attention approximates the original dot-product similarity calculation by introducing techniques such as a kernel function [8, 16] or by employing low-rank decomposition [39]. By changing the order of computation or using matrix decomposition, these methods avoid the explicit construction of the attention matrix, which has a quadratic complexity. This approach significantly reduces the computational and memory burden of Transformers for long-sequence tasks. Vision LSTM [2], a linear RNN methods, uses a stacked LSTM structure to extract multi-scale visual features and integrates temporal information through gating mechanisms, achieving strong performance in tasks such as video understanding. Mamba [12] and its variant Video Mamba [19], based on state-space models (SSMs), efficiently model videos and capture rich spatiotemporal dependencies in long-sequence scenarios, demonstrating excellent scalability and expressive power. However, these

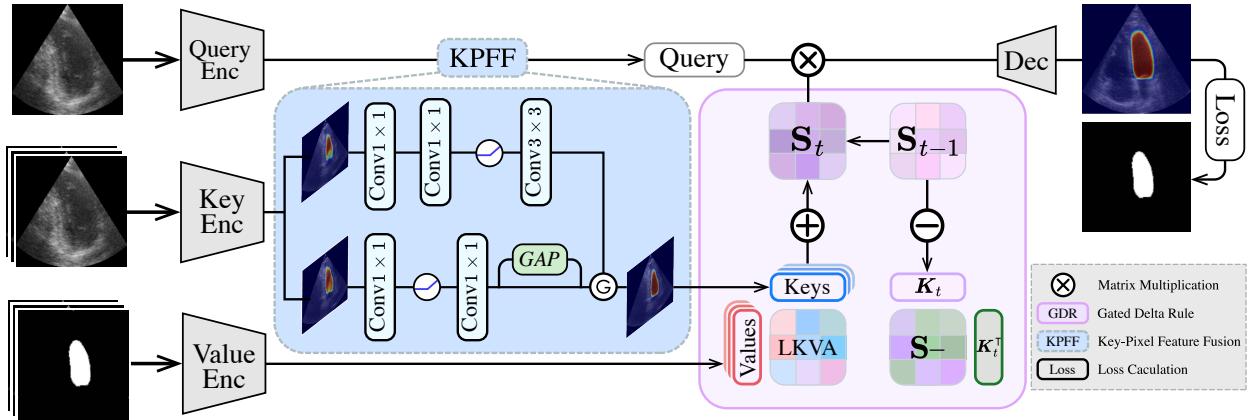


Figure 2. An illustration of GDKVM architecture. Linear Key-Value Association defines frame-to-frame causal relations as the state transition matrix. Gated Delta Rule helps in dynamically managing memory. Key-Pixel Feature Fusion fuses the local key feature, the global key feature with the pixel feature.

methods have not yet been thoroughly validated in the field of medical ultrasound imaging. This domain is particularly challenging due to its complex noise and low contrast. Optimizing model structures and training strategies for medical scenarios remains key to further unleashing the potential of linear attention.

3. Method

3.1. Overview

We propose **GDKVM**, a linear key-value memory network designed to address the challenges of high computational cost and poor noise robustness in echocardiography video segmentation. The architecture is illustrated in Fig. 2.

The key encoder uses the ResNet-50 [13] backbone to extract raw features. The KPFF module integrates the fused features obtained by combining key features, global key features, and pixel features to get the **Keys**. The value encoder encodes the original frame and the previous prediction mask into the **Values**. We combine **Keys** and **Values** into a recurrent hidden **State**. During the linear update of the state, we use 2 data-dependent matrices, α_t and β_t , to control the state decay, both of which are projected from the previous state S_{t-1} . β_t balances how much new information is incorporated into the state, while α_t adaptively forgets old memories. Using the **Query** frame, the readout is derived from the state S_t , and the predicted segmentation result is obtained through the Decoder. To simulate a clinical scenario, the model does not have access to the ground truth when making a prediction. During each training step, the model first predicts the morphology of the first and last frames; only then is the loss calculated by comparing the prediction against the ground truth [10].

3.2. Linear Key-Value Association

Softmax matching. The input frames of an echocardiography video can be described as $I_{1:t} \in \mathbb{R}^{tHW \times C_d}$. Current Transformers and Space-Time Memory Networks [3, 5–7] use softmax attention matching to maintain the memory bank for segmenting the object:

$$O_t = \sum_{i=1}^t \frac{\exp(\mathbf{K}_i^\top \mathbf{Q}_t)}{\sum_{j=1}^t \exp(\mathbf{K}_j^\top \mathbf{Q}_t)} V_i, \quad (1)$$

where $\mathbf{Q}_t, \mathbf{K}_t, \mathbf{V}_t \in \mathbb{R}^{tHW \times C_d}$. The $\mathcal{O}(t^2 C_d)$ computational cost required by attention significantly increases inference time and creates significant pressure for real-time application requirements.

Linear matching. Linear attention [16, 36] uses the kernel method to replace the exponential kernel $\exp(\mathbf{K}_i^\top \mathbf{Q}_t)$ with a non-negative dot-product $\phi(\mathbf{K}_i)^\top \phi(\mathbf{Q}_t)$. Figure 2 shows that LKVA allows us to express the softmax matching as a linear RNN with a matrix-form hidden state S_t :

$$\begin{aligned} O_t &= \sum_{i=1}^t \frac{\phi(\mathbf{K}_i)^\top \phi(\mathbf{Q}_t)}{\sum_{j=1}^t \phi(\mathbf{K}_j)^\top \phi(\mathbf{Q}_t)} V_i \\ &= \frac{(\sum_{i=1}^t \mathbf{V}_i \phi(\mathbf{K}_i)^\top) \phi(\mathbf{Q}_t)}{(\sum_{j=1}^t \phi(\mathbf{K}_j)^\top) \phi(\mathbf{Q}_t)} \\ &= \frac{S_t \phi(\mathbf{Q}_t)}{Z_t^\top \phi(\mathbf{Q}_t)}, \end{aligned} \quad (2)$$

where $S_t = \sum_{i=1}^t \mathbf{V}_i \phi(\mathbf{K}_i)^\top \in \mathbb{R}^{C_v \times C_k}$ can be regarded as a two-dimensional recurrent hidden state of fixed size, and $Z_t = \sum_{j=1}^t \phi(\mathbf{K}_j) \in \mathbb{R}^{C_k}$. LKVA maps the key space \mathbb{R}^{C_k} and value space \mathbb{R}^{C_v} of each timestep into a

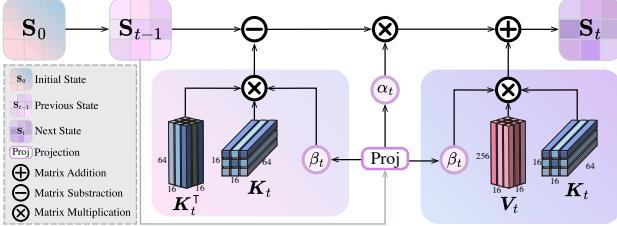


Figure 3. An illustration of Gated Delta Rule (GDR). It leverages a projected β_t for efficient key–value associations and frame updates, while a projected α_t quickly adapts to drastic heart shape changes and preserves crucial long-term information.

more expressive state space $\mathbb{R}^{C_v \times C_k}$. Some recent linear RNN models [27, 32, 34, 41] have further provided a simplified linear Transformer formulation:

$$\begin{aligned} \mathbf{S}_t &= \mathbf{S}_{t-1} + \mathbf{V}_t \mathbf{K}_t^\top \in \mathbb{R}^{C_v \times C_k}, \\ \mathbf{O}_t &= \mathbf{S}_t \mathbf{K}_t \in \mathbb{R}^{HW \times C_v}. \end{aligned} \quad (3)$$

Because this recurrent form of linear attention only needs to accumulate and update \mathbf{S}_t along with the associated key–value representations, its computational complexity drops to $\mathcal{O}(tC_vC_k)$.

The LKVA method possesses linear computational complexity, a feature of critical importance for clinical applications. This ensures the speed and scalability required for processing long-duration echocardiography videos, making real-time analysis possible without sacrificing the model’s ability to capture long-range dependencies. The sequential processing approach of LKVA, combined with its highly expressive state representation, is exceptionally well-suited for capturing the continuous and often periodic motion of cardiac structures by progressively accumulating contextual information frame by frame.

3.3. Gated Delta Rule

Delta rule. Pure linear attention superimposes all historical information with equal weight. When a sufficient number of tokens are superimposed, the proportion of information from each individual token becomes extremely small. Relying solely on a fixed-size state matrix \mathbf{S}_t makes it impossible to accurately reconstruct even an arbitrary output \mathbf{O}_t , causing the state representations to become indistinct and blurred [34, 36, 41]. To mitigate this problem, the model should discard less important associations based on how newly arriving keys interact with existing contents.

In video tasks, from the perspective of key–value retrieval, we first remove the old association for the current frame’s key \mathbf{K}_t from the previous time-step’s state \mathbf{S}_{t-1} , which is represented as $\mathbf{V}_t^{\text{old}} = \mathbf{S}_{t-1} \mathbf{K}_t$. A new value $\mathbf{V}_t^{\text{new}}$ is obtained by interpolating between the old value and the current frame’s target \mathbf{V}_t , replacing the old value in the

state:

$$\begin{aligned} \mathbf{S}_t &= \mathbf{S}_{t-1} \\ &\quad - \underbrace{(\mathbf{S}_{t-1} \mathbf{K}_t)}_{\mathbf{V}_t^{\text{old}}} \mathbf{K}_t^\top + \underbrace{(\beta_t \mathbf{V}_t + (\mathbf{I} - \beta_t) \mathbf{S}_{t-1} \mathbf{K}_t)}_{\mathbf{V}_t^{\text{new}}} \mathbf{K}_t^\top \\ &= \mathbf{S}_{t-1} (\mathbf{I} - \beta_t \mathbf{K}_t \mathbf{K}_t^\top) + \beta_t \mathbf{V}_t \mathbf{K}_t^\top, \end{aligned} \quad (4)$$

where $\beta_t \in \mathbb{R}^{C_v \times C_k}$ is a data-dependent matrix, projected from the previous state \mathbf{S}_{t-1} , and it represents the soft writing strength. When rapid cardiac structure motion or clear boundaries appear in the video, the model is assigned a larger β_t value to strengthen the learning of key information in the current frame. Conversely, when the image quality is poor or when images are filled with speckle noise, the model can use it to erase interfering information from the previous frame caused by probe movement or artifacts.

Data-dependent decay. In echocardiography video segmentation, the cardiac structures undergo rapid and dramatic dynamic changes. As a result, the model’s memory state can easily become saturated with a large volume of complex features, which in turn leads to difficulties in key–value retrieval [9, 42]. We project a data-dependent decay matrix $\alpha_t \in \mathbb{R}^{C_v \times C_k}$ from the previous state to control the attenuation of past memory:

$$\mathbf{S}_t = \mathbf{S}_{t-1} (\alpha_t (\mathbf{I} - \beta_t \mathbf{K}_t \mathbf{K}_t^\top)) + \beta_t \mathbf{V}_t \mathbf{K}_t^\top. \quad (5)$$

Figure 3 and Eq. (5) show that the introduction of α_t weakens the strong regularization constraint between \mathbf{S}_t and \mathbf{S}_{t-1} . This allows the model to selectively forget certain past information and thus dynamically allocate memory resources to the most critical cardiac cycle features. The α_t enables the model to free up sufficient memory space for new key frames or structures, including frames corresponding to abnormal cardiac rhythms, while preventing irrelevant or redundant historical features from occupying memory.

For echocardiography video segmentation tasks with strong temporal correlations, GDR offers a better balance between memory retention and new information integration compared to simple additive updates. It preserves long-term important information while more rapidly adapting to drastic changes in heart shape, thereby improving the accuracy of spatiotemporal feature extraction.

3.4. Key-Pixel Feature Fusion

In echocardiography video segmentation, issues such as speckle noise, artifacts caused by signal attenuation, and low contrast between tissue and blood pose significant challenges to the feature extraction capabilities of models. Linear models demonstrate rapid and stable memory capabilities in forming key–value associations. However, for inputs with rich spatial structures, such as images, relying

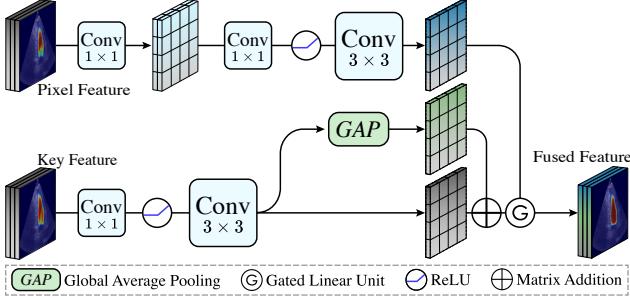


Figure 4. An illustration of Key-Pixel Feature Fusion (KPFF). KPFF fuses the original coarse key features, global key features, and extracted pixel features into key-pixel features, providing enhanced spatiotemporal robustness.

solely on this type of associative memory may be insufficient for complex spatial reasoning. Sub-quadratic attention models like xLSTM [2] widely utilize small-kernel 1D depth-wise separable convolutions to extract local semantic information. However, these spatial features often suffer from over-smoothing, leading to the loss of some global context. Figure 4 shows that KPFF merges the convolution outputs learned from local key feature \mathbf{F}_K , global key feature $\mathbf{F}_{\text{Global}}$ and pixel-based feature \mathbf{F}_{Pix} , each capturing information at different spatial scales, into a single network branch:

$$\begin{aligned} \mathbf{F}_{\text{Global}} &= \text{Expand}(\text{GAP}(\mathbf{F}_K)), \\ \mathbf{G} &= \sigma(\text{Conv}_{\text{gate}}(\mathbf{F}_K + \mathbf{F}_{\text{Global}})), \\ \mathbf{F}_{\text{fused}} &= \mathbf{G} \cdot (\mathbf{F}_K + \mathbf{F}_{\text{Global}}) + (1 - \mathbf{G}) \cdot \mathbf{F}_{\text{Pix}}, \end{aligned} \quad (6)$$

where $\sigma(\cdot)$ denotes the sigmoid function. When local features \mathbf{F}_K become unreliable due to speckle noise or signal dropout, $\mathbf{F}_{\text{Global}}$ can provide a stable global context as a supplement. By fusing the smoothed global and local features with pixel-level features \mathbf{F}_{Pix} , which retain high-frequency information, KPFF effectively preserves critical diagnostic details and prevents the loss of fine-grained information, such as valve opening and closing or regional wall motion abnormalities.

4. Experiments

4.1. Setup

Datasets. We evaluate our proposed model on two public echocardiography video datasets, CAMUS [18] and EchoNet-Dynamic [30]. **CAMUS** contains 500 patient cases, each with apical four-chamber and two-chamber view videos. For each video, it provides frame-by-frame annotations spanning from end-diastole (*ED*) to end-systole (*ES*). These annotations include labeled contours of the left ventricular endocardium, left ventricular epicardium, and the left atrium. **EchoNet-Dynamic** contains 10,030 api-

cal four-chamber view videos. It provides annotations for each video only at two key frames: *ED* and *ES*. We used the original data splits for CAMUS and EchoNet-Dynamic, resizing them to resolutions of 256×256 and 128×128 , respectively. From each cardiac video, we uniformly sampled 10 frames.

Evaluation metrics. To ensure a fair and comprehensive evaluation, we assessed the model’s performance using both standard segmentation metrics and key clinical indices. The geometric accuracy of the segmentation was evaluated with four common metrics: the mean Dice coefficient (mDice), mean Intersection over Union (mIoU), Hausdorff Distance (HD), and Average Surface Distance (ASD). To evaluate the model’s clinical utility for left ventricular analysis, we estimated the Left Ventricular Ejection Fraction (LV_{EF}) from the segmentation masks using Simpson’s rule [17]. We then assessed the quality of these estimations by calculating the Pearson correlation coefficient (corr), mean bias (bias), and standard deviation (std) between the predicted and ground truth LV_{EF} values.

Implementation details. All training runs on a single RTX 3090 GPU. For CAMUS and EchoNet-Dynamic, we conducted training over 1500 iterations. The AdamW optimizer [25] was used with a learning rate of 1e-4, a batch size of 10 and a variety of data augmentation techniques were applied, including gamma augmentation, random scaling, random rotation, and random contrast adjustment, each with a probability of 0.5. We clip the global gradient norm to $\lambda = 3$ and use stable data augmentation. We use a combined loss function of cross-entropy and soft dice loss with equal weighting following Cheng et al. [7].

4.2. Comparison with State-of-the-art Methods

We compare GDKVM with four task-specific methods (PKEchoNet [40], DSA [22], MemSAM [10], Sim-LVSeg [26]) and four related approaches (Xmem++ [3], Cutie [7], VideoMamba [19], Vision LSTM [2]).

These prior works incorporate memory mechanisms, time-series modeling, or efficient attention modules and have advanced medical video segmentation. However, they still face challenges in capturing complex spatial variations and retrieving temporal information accurately. GDKVM addresses these limitations by introducing key-value associations within a linear matching framework, which reduces the typical high computational overhead of conventional attention. In Tab. 1, GDKVM delivers higher overall metrics on two datasets and outperforms four widely used task-specific models. In Fig. 5, GDKVM achieves an mDice of 95.11 at 37 FPS with 35.2 M parameters. GDKVM leverages GDR for dynamic memory management to discard irrelevant information promptly and update essential mem-

Method	Venue & Year	CAMUS				EchoNet-Dynamic			
		mDice	mIoU	HD	ASD	mDice	mIoU	HD	ASD
XMem++ [3]	ICCV'23	89.38	85.81	4.03	4.87	87.51	83.57	3.14	2.69
Cutie [7]	CVPR'24	91.09	87.97	3.89	3.74	88.96	85.63	2.89	2.24
VideoMamba [19]	ECCV'24	91.96	89.04	3.48	3.31	90.22	87.03	2.79	2.05
Vision LSTM [2]	ICLR'25	92.14	89.11	3.79	3.39	90.24	89.14	2.65	1.69
PKEchoNet [40]	AAAI'23	93.49	90.95	3.42	2.93	92.60	89.89	2.53	1.48
DSA [22]	TMI'24	94.25	91.80	3.27	2.37	92.91	90.26	2.46	1.44
MemSAM [10]	CVPR'24	93.63	90.97	3.47	2.60	92.71	89.90	2.56	1.51
SimLVSeg [26]	UMB'24	92.54	89.71	3.65	3.12	91.91	89.08	2.65	1.65
GDKVM	-	95.11	92.97	3.05	1.98	93.46	90.86	2.38	1.36

Table 1. Segmentation performance of GDKVM with state-of-the-art methods on CAMUS and EchoNet-Dynamic.

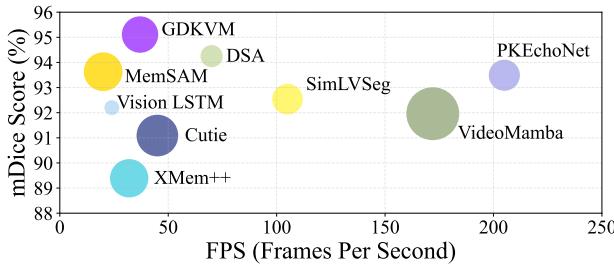


Figure 5. Comparison of different methods in terms of speed (FPS) and Dice score on CAMUS, with bubble size representing the number of parameters.

ory content efficiently, matching or surpassing specialized models in segmentation performance. This design avoids bottlenecks arising from incomplete memory mechanisms or overly complex training procedures, allowing GDKVM to maintain both accuracy and speed.

Clinical metric. Based on the segmentation results, we estimated the clinical metric LV_{EF} . Table 2 shows that GDKVM obtains the highest Pearson correlation coefficient of 0.904 and the lowest bias and standard deviation ($-0.19 \pm 11.3\%$). This indicates a strong agreement between our predictions and the ground truth. Figure 6 provides a visual representation, where the linear regression plot shows a strong correlation, and the Bland–Altman plot demonstrates high consistency.

Visual comparison with SOTA A visual comparison of GDKVM with other state-of-the-art methods is presented in Fig. 7, illustrating the ability of GDKVM to preserve structural integrity and delineate precise boundaries in difficult frames. Although prior approaches have significantly advanced heart cavity segmentation, they can still produce

Method	CAMUS	
	corr	bias \pm std (%)
XMem++ [3]	0.746	1.70 ± 21.9
Cutie [7]	0.787	1.67 ± 21.7
VideoMamba [19]	0.780	-4.49 ± 19.4
Vision LSTM [2]	0.806	-0.31 ± 18.8
PKEchoNet [40]	0.862	-1.53 ± 16.4
DSA [22]	0.891	0.86 ± 13.4
MemSAM [10]	0.878	-0.89 ± 12.3
SimLVSeg [26]	0.895	1.83 ± 13.8
GDKVM	0.904	-0.19 ± 11.3

Table 2. Clinical metrics comparison against different state-of-the-art methods on CAMUS.

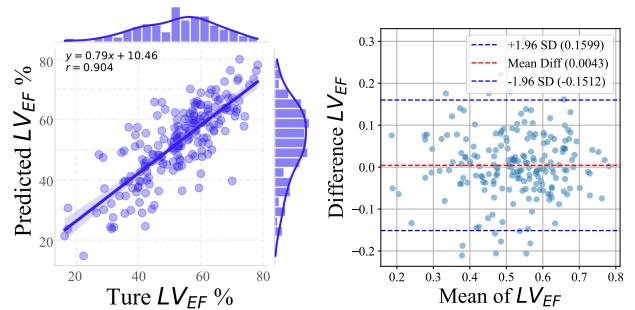


Figure 6. Linear regression and Bland–Altman plots for clinical metric LV_{EF} on CAMUS.

over-segmentation (including extraneous areas beyond the heart cavity) or under-segmentation (missing parts of the ventricle). In these challenging cases, existing methods sometimes deviate from the ground truth shape by forming gaps or protrusions. GDKVM, by contrast, adheres more closely to the anatomical contours, generating cleaner masks with noticeably fewer artifacts. This property yields

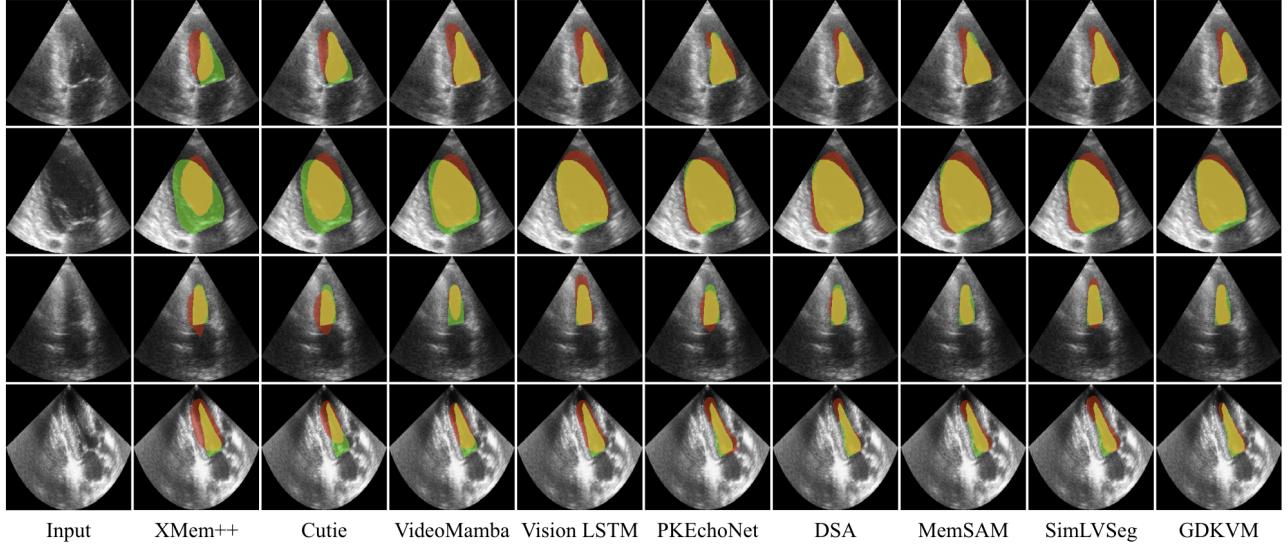


Figure 7. Visual comparison with state-of-the-art methods on the CAMUS test set. Green, red, and yellow regions represent the ground truth, prediction, and overlapping regions, respectively.

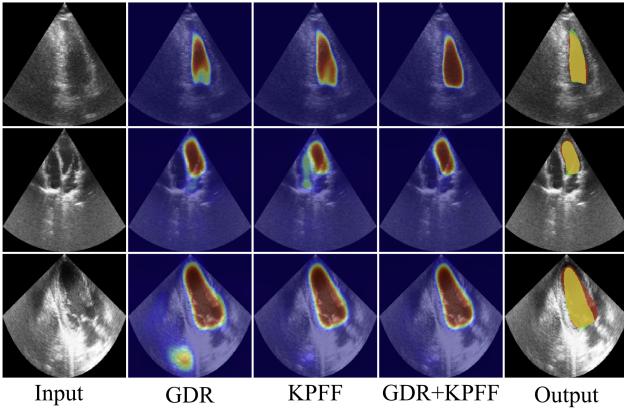


Figure 8. Visualization of segmentation results using different ablation strategies of GDKVM on CAMUS.

LKVA	GDR	KPFF	mDice	mIoU	HD	ASD
✓			93.10	90.46	3.65	2.85
✓	✓		94.49	92.11	3.21	2.19
✓		✓	93.30	90.78	3.55	2.74
✓	✓	✓	95.11	92.97	3.05	1.98

Table 3. Ablation study on different components of GDKVM on CAMUS.

higher segmentation fidelity on frames where earlier approaches often falter.

4.3. Ablation Study

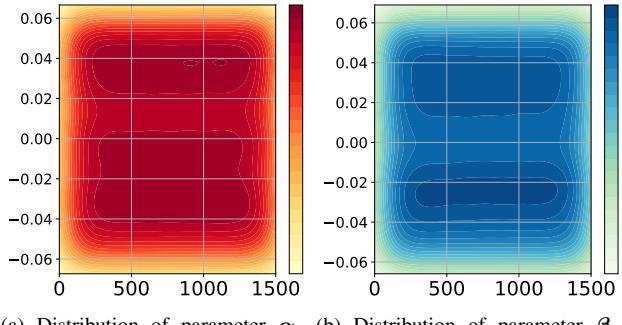
To evaluate each component of GDKVM, we performed ablation studies on CAMUS, as shown in Fig. 8 and Tab. 3.

Strategy	Recurrence Equation	mDice	Inf. Time
Baseline	$\mathbf{S}_t = \mathbf{S}_{t-1} + \mathbf{V}_t \mathbf{K}_t^\top$	93.30	151.61 ms
Sanity Check	$\mathbf{S}_t = \mathbf{S}_{t-1} - (\mathbf{S}_{t-1} \mathbf{K}_t) \mathbf{K}_t^\top + \mathbf{V}_t \mathbf{K}_t^\top$	74.68	155.09 ms
w/o α_t	$\mathbf{S}_t = \mathbf{S}_{t-1} (\mathbf{I} - \beta_t \mathbf{K}_t \mathbf{K}_t^\top) + \beta_t \mathbf{V}_t \mathbf{K}_t^\top$	94.57	158.77 ms
w/o β_t	$\mathbf{S}_t = \alpha_t \mathbf{S}_{t-1} + \mathbf{V}_t \mathbf{K}_t^\top$	94.26	156.90 ms
GDR	$\mathbf{S}_t = \mathbf{S}_{t-1} (\alpha_t (\mathbf{I} - \beta_t \mathbf{K}_t \mathbf{K}_t^\top)) + \beta_t \mathbf{V}_t \mathbf{K}_t^\top$	95.11	160.62 ms

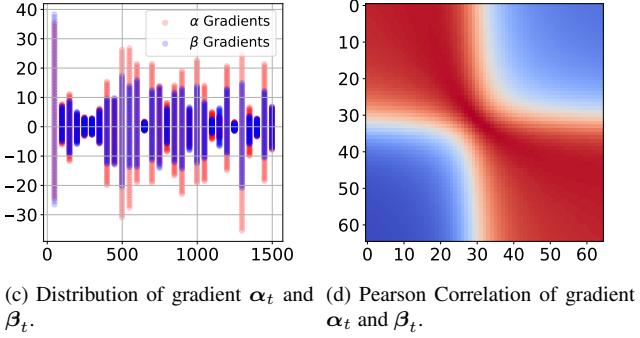
Table 4. Comparison of different strategies to update the state.

Effectiveness of Key-Pixel Feature Fusion KPFF mechanism is critical for handling the significant scale variations and spatial transformations inherent in echocardiography videos. Disabling this mechanism and relying on a single key feature extraction approach causes a marked degradation in segmentation performance (Tab. 3). The model becomes particularly vulnerable to image noise and distortion, leading to less accurate boundary delineation, as visually evidenced in Fig. 8. These results establish that the proposed feature fusion is essential for achieving robust and precise segmentation in such dynamic clinical data.

Effectiveness of Gated Delta Rule In Tab. 4, we compare several linear update strategies for the memory state. Directly replacing old memory with new information yields a Dice score of 74.68, indicating that a simple structural change in the update formula is insufficient. Introducing either α_t or β_t alone slightly improves performance, suggesting that selectively forgetting old content or reinforcing new features helps capture dynamic cardiac boundaries. Combining α_t and β_t further provides bidirectional control: α_t governs the decay of stored shapes to remove noise or er-



(a) Distribution of parameter α_t . (b) Distribution of parameter β_t . Darker regions indicate higher density.



(c) Distribution of gradient α_t and β_t . (d) Pearson Correlation of gradient α_t and β_t .

Figure 9. Weights of parameters α_t and β_t over training steps on CAMUS.

rors, and β_t sets the strength of newly written features to accommodate abrupt structural changes. As a result, GDR reaches a Dice score of 95.11, surpassing the baseline by about 1.7 percentage points.

In Figs. 9a and 9b, α_t and β_t concentrate at specific intervals. During stable frames, both decay and writing intensities remain moderate. In frames with blurred boundaries or heavy noise, α_t increases forgetting to discard flawed information. In cases of sudden contour changes, β_t boosts new feature updates. Large gradient jumps, such as +40 or -30 in Fig. 9c, illustrate the importance of robust gating. In Fig. 9d, the correlation pattern forms block-like and diagonal structures, reflecting the coordinated adjustment of α_t and β_t . This mechanism refines echocardiographic boundaries, filters noise during complex cardiac motion, and retains useful features in stable periods for more precise segmentation.

4.4. Discussions and Limitations

Our work has some limitations. While the model performs well in ultrasound cardiac video segmentation, its full potential in other medical video object segmentation scenarios remains to be explored. This includes testing on fully-supervised and semi-supervised benchmarks, which would allow us to evaluate the model’s performance in accurately

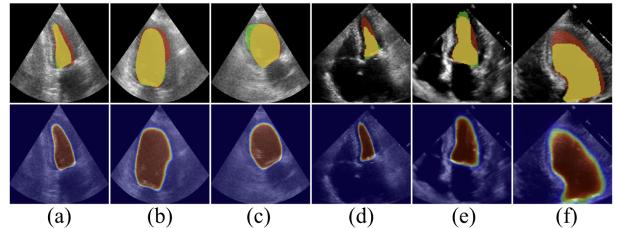


Figure 10. Failure cases on CAMUS (a-c) and EchoNet-Dynamic (d-f) test sets.

modeling long-term temporal sequences. Future work could also focus on developing metrics that better assess the temporal consistency of sequences. Figure 10 shows that the model sometimes over-relies on its own predictions during the final optimization stage. For some difficult samples, this hinders the precise delineation of contours. A potential direction is to develop more flexible boundary definitions or adaptive corrections to address these artifacts. Our current intermediate state matrix is not a standard square matrix, which restricts chunkwise parallelization and hardware acceleration. Future work could explore a square matrix design to enable more efficient parallel operators, reduce training time, and improve inference speed.

5. Conclusion

This work describes GDKVM, a novel linear key-value memory network designed for efficient and robust echocardiography video segmentation. GDKVM directly addresses the clinical need for real-time analysis by leveraging a Linear Key-Value Association (LKVA) to achieve linear computational complexity. To accurately model cardiac dynamics, Gated Delta Rule (GDR) employs two data-dependent decays to intelligently update and forget memory, capturing rapid motion while maintaining temporal consistency. Key-Pixel Feature Fusion (KPFF) module enhances robustness against image artifacts and noise by merging features from different receptive fields. Experiments demonstrate that GDKVM consistently achieves state-of-the-art segmentation accuracy while remaining highly efficient across two echocardiography video datasets (CAMUS and EchoNet-Dynamic).

Acknowledgements

This work was supported partly by the National Natural Science Foundation of China (No. 62273241), the Natural Science Foundation of Guangdong Province, China (No. 2024A1515011946), and the Hong Kong RGC Collaborative Research Fund (project no. C5055-24G).

References

- [1] Alyaa Amer, Xujiong Ye, and Faraz Janan. Resdunet: A deep learning-based left ventricle segmentation method for echocardiography. *IEEE Access*, 9:159755–159763, 2021. 1
- [2] Maximilian Beck, Korbinian Pöppel, Markus Spanring, Andreas Auer, Oleksandra Prudnikova, Michael Kopp, Günter Klambauer, Johannes Brandstetter, and Sepp Hochreiter. xlstm: Extended long short-term memory. In *Thirty-eighth Conference on Neural Information Processing Systems*, 2024. 2, 5, 6
- [3] Maksym Bekuzarov, Ariana Bermudez, Joon-Young Lee, and Hao Li. Xmem++: Production-level video segmentation from few annotated frames, 2023. 2, 3, 5, 6
- [4] Yida Chen, Xiaoyan Zhang, Christopher M Haggerty, and Joshua V Stough. Assessing the generalizability of temporally coherent echocardiography video segmentation. In *Medical Imaging 2021: Image Processing*, pages 463–469. SPIE, 2021. 2
- [5] Ho Kei Cheng and Alexander G. Schwing. Xmem: Long-term video object segmentation with an atkinson-shiffrin memory model, 2022. 2, 3
- [6] Ho Kei Cheng, Yu-Wing Tai, and Chi-Keung Tang. Rethinking space-time networks with improved memory coverage for efficient video object segmentation, 2021.
- [7] Ho Kei Cheng, Seoung Wug Oh, Brian Price, Joon-Young Lee, and Alexander Schwing. Putting the object back into video object segmentation, 2024. 2, 3, 5, 6
- [8] Krzysztof Choromanski, Valerii Likhosherstov, David Dohan, Xingyou Song, Andreea Gane, Tamas Sarlos, Peter Hawkins, Jared Davis, Afroz Mohiuddin, Lukasz Kaiser, et al. Rethinking attention with performers. *arXiv preprint arXiv:2009.14794*, 2020. 2
- [9] Tri Dao and Albert Gu. Transformers are ssms: Generalized models and efficient algorithms through structured state space duality, 2024. 4
- [10] Xiaolong Deng, Huisi Wu, Runhao Zeng, and Jing Qin. Memsam: taming segment anything model for echocardiography video segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9622–9631, 2024. 2, 3, 5, 6
- [11] Marwa Chendeb El Rai, Muna Darweesh, and Mina Al-Saad. Semi-supervised segmentation of echocardiography videos using graph signal processing. *Electronics*, 11(21):3462, 2022. 1
- [12] Albert Gu and Tri Dao. Mamba: Linear-time sequence modeling with selective state spaces. *arXiv preprint arXiv:2312.00752*, 2023. 2
- [13] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition, 2015. 3
- [14] Vladimir Iglovikov and Alexey Shvets. Ternausnet: U-net with vgg11 encoder pre-trained on imagenet for image segmentation. *arXiv preprint arXiv:1801.05746*, 2018. 1
- [15] Debesh Jha, Pia H Smedsrød, Michael A Riegler, Dag Johansen, Thomas De Lange, Pål Halvorsen, and Håvard D Johansen. Resunet++: An advanced architecture for medical image segmentation. In *2019 IEEE International Symposium on Multimedia (ISM)*, pages 225–2255. IEEE, 2019. 1
- [16] Angelos Katharopoulos, Apoorv Vyas, Nikolaos Pappas, and François Fleuret. Transformers are rnns: Fast autoregressive transformers with linear attention. In *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 13-18 July 2020, Virtual Event*, pages 5156–5165. PMLR, 2020. 2, 3
- [17] Roberto M. Lang, Luigi P. Badano, Victor Mor-Avi, Jonathan Afilalo, Anderson Armstrong, Laura Ernande, Frank A. Flachskampf, Elyse Foster, Steven A. Goldstein, Tatiana Kuznetsova, Patrizio Lancellotti, Denisa Muraru, Michael H. Picard, Ernst R. Rietzschel, Lawrence Rudski, Kirk T. Spencer, Wendy Tsang, and Jens-Uwe Voigt. Recommendations for cardiac chamber quantification by echocardiography in adults: An update from the american society of echocardiography and the european association of cardiovascular imaging. *Journal of the American Society of Echocardiography*, 28(1):1–39.e14, 2015. 5
- [18] Sarah Leclerc, Erik Smistad, Joao Pedrosa, Andreas Østvik, Frederic Cervenansky, Florian Espinosa, Torvald Espeland, Erik Andreas Rye Berg, Pierre-Marc Jodoin, Thomas Grennier, et al. Deep learning for segmentation using an open large-scale dataset in 2d echocardiography. *IEEE Transactions on Medical Imaging*, 38(9):2198–2210, 2019. 2, 5
- [19] Kunchang Li, Xinhao Li, Yi Wang, Yinan He, Yali Wang, Limin Wang, and Yu Qiao. Videomamba: State space model for efficient video understanding, 2024. 2, 5, 6
- [20] Xiaoshan Li, Lisi Liao, Kai Wu, Alexander Thomas Meng, Yitao Jiang, Yuan Zhu, Chen Cui, Xiaowei Xu, Bobo Shi, and Hongwen Fei. An automatic and real-time echocardiography quality scoring system based on deep learning to improve reproducible assessment of left ventricular ejection fraction. *Quantitative Imaging in Medicine and Surgery*, 15(1):770, 2024. 1
- [21] Minqi Liao, Yifan Lian, Yongzhao Yao, Lihua Chen, Fei Gao, Long Xu, Xin Huang, Xinxing Feng, and Suxia Guo. Left ventricle segmentation in echocardiography with transformer. *Diagnostics*, 13(14):2365, 2023. 2
- [22] Jingyin Lin, Wende Xie, Li Kang, and Huisi Wu. Dynamic-guided spatiotemporal attention for echocardiography video segmentation. *IEEE Transactions on Medical Imaging*, 2024. 5, 6
- [23] Shu Liu, Lu Qi, Haifang Qin, Jianping Shi, and Jiaya Jia. Path aggregation network for instance segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8759–8768, 2018. 2
- [24] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3431–3440, 2015. 1
- [25] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization, 2019. 5
- [26] Fadillah Maani, Asim Ukaye, Nada Saadi, Numan Saeed, and Mohammad Yaqub. Simlvseg: Simplifying left ventricular segmentation in 2-d+ time echocardiograms with self- and weakly supervised learning. *Ultrasound in Medicine & Biology*, 50(12):1945–1954, 2024. 5, 6
- [27] Eric Martin and Chris Cundy. Parallelizing linear recurrent neural nets over sequence length. In *6th International Con-*

- ference on Learning Representations, ICLR 2018, Vancouver, BC, Canada, April 30 - May 3, 2018, Conference Track Proceedings.* OpenReview.net, 2018. [4](#)
- [28] Seoung Wug Oh, Joon-Young Lee, Ning Xu, and Seon Joo Kim. Video object segmentation using space-time memory networks, 2019. [2](#)
- [29] Keiron O’Shea and Ryan Nash. An introduction to convolutional neural networks, 2015. [1](#)
- [30] David Ouyang, Bryan He, Amirata Ghorbani, Matt P Lungren, Euan A Ashley, David H Liang, and James Y Zou. Echonet-dynamic: a large new cardiac motion video data resource for medical machine learning. In *NeurIPS ML4H Workshop*, pages 1–11, 2019. [2, 5](#)
- [31] Nathan Painchaud, Nicolas Duchateau, Olivier Bernard, and Pierre-Marc Jodoin. Echocardiography segmentation with enforced temporal consistency. *IEEE Transactions on Medical Imaging*, 41(10):2867–2878, 2022. [1](#)
- [32] Zhen Qin, Xiaodong Han, Weixuan Sun, Dongxu Li, Lingpeng Kong, Nick Barnes, and Yiran Zhong. The devil in linear transformer. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7025–7041, Abu Dhabi, United Arab Emirates, 2022. Association for Computational Linguistics. [4](#)
- [33] Olaf Ronneberger, Philipp Fischer, and Thomas 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 18*, pages 234–241. Springer, 2015. [1](#)
- [34] Imanol Schlag, Kazuki Irie, and Jürgen Schmidhuber. Linear transformers are secretly fast weight programmers. In *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, pages 9355–9366. PMLR, 2021. [4](#)
- [35] Hajar Sadeghi Sokeh, Vasileios Argyriou, Dorothy Mon-ekosso, and Paolo Remagnino. Superframes, a temporal video segmentation. In *2018 24th International Conference on Pattern Recognition (ICPR)*, pages 566–571. IEEE, 2018. [2](#)
- [36] Yutao Sun, Li Dong, Shaohan Huang, Shuming Ma, Yuqing Xia, Jilong Xue, Jianyong Wang, and Furu Wei. Retentive network: A successor to transformer for large language models. *ArXiv preprint*, abs/2307.08621, 2023. [3, 4](#)
- [37] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017. [2](#)
- [38] Huina Wang, Lan Wei, Bo Liu, Jianqiang Li, Jinshu Li, Juan Fang, and Catherine Mooney. Transformer-based explainable model for breast cancer lesion segmentation. *Applied Sciences*, 15(3):1295, 2025. [2](#)
- [39] Sinong Wang, Belinda Z Li, Madien Khabsa, Han Fang, and Hao Ma. Linformer: Self-attention with linear complexity. *arXiv preprint arXiv:2006.04768*, 2020. [2](#)
- [40] Huisi Wu, Jingyin Lin, Wende Xie, and Jing Qin. Super-efficient echocardiography video segmentation via proxy- and kernel-based semi-supervised learning. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2803–2811, 2023. [5, 6](#)
- [41] Songlin Yang, Bailin Wang, Yu Zhang, Yikang Shen, and Yoon Kim. Parallelizing linear transformers with the delta rule over sequence length. *NeurIPS*, 2024. [4](#)
- [42] Songlin Yang, Jan Kautz, and Ali Hatamizadeh. Gated delta networks: Improving mamba2 with delta rule. In *The Thirteenth International Conference on Learning Representations*, 2025. [4](#)
- [43] Yingyu Yang, Marie Rocher, Pamela Moceri, Maxime Sermesant, et al. Echocardiography analysis with deep learning using priors: Multi-centric evaluation of generalisation. *Machine Learning for Biomedical Imaging*, 2(November 2024 issue):2293–2325, 2024. [1](#)