

EfficientMorph: Parameter-Efficient Transformer-Based Architecture for 3D Image Registration

Abu Zahid Bin Aziz**, Mokshagna Sai Teja Karanam**, Tushar Kataria, and
Shireen Y. Elhabian³

¹ Kahlert School of Computing, University of Utah

² Scientific Computing and Imaging Institute, University of Utah

³ Corresponding author, **Equal Contribution

{mkaranam,zahid,tushar.kataria,shireen}@sci.utah.edu

Abstract. Transformers have emerged as the state-of-the-art architecture in medical image registration, outperforming convolutional neural networks (CNNs) by addressing their limited receptive fields and overcoming gradient instability in deeper models. Despite their success, transformer-based models require substantial resources for training, including data, memory, and computational power, which may restrict their applicability for end users with limited resources. In particular, existing transformer-based 3D image registration architectures face three critical gaps that challenge their efficiency and effectiveness. Firstly, while mitigating the quadratic complexity of full attention by focusing on local regions, window-based attention mechanisms often fail to adequately integrate local and global information. Secondly, feature similarities across attention heads that were recently found in multi-head attention architectures indicate a significant computational redundancy, suggesting that the capacity of the network could be better utilized to enhance performance. Lastly, the granularity of tokenization, a key factor in registration accuracy, presents a trade-off; smaller tokens improve detail capture at the cost of higher computational complexity, increased memory demands, and a risk of overfitting. Here, we propose EfficientMorph, a transformer-based architecture for unsupervised 3D image registration. It optimizes the balance between local and global attention through a plane-based attention mechanism, reduces computational redundancy via cascaded group attention, and captures fine details without compromising computational efficiency, thanks to a Hi-Res tokenization strategy complemented by merging operations. We compare the effectiveness of EfficientMorph on two public datasets, OASIS and IXI, against other state-of-the-art models. Notably, EfficientMorph sets a new benchmark for performance on the OASIS dataset with $\sim 16\text{-}27\times$ fewer parameters.

Keywords: Windowed Attention · 3D Image Registration · Unsupervised Learning · Parameter-Efficient Transformer Architectures

1 Introduction

Image registration is a critical task for various medical imaging applications in fields such as image-guided surgery [1], radiation therapy planning [20], image fusion for multimodality imaging [11], and quality enhancement [3]. Registration entails determining the spatial alignment between two volumes, typically referred to as the *fixed* and *moving* images, by identifying correspondences among similar structures or features and their relative positions. Conventional approaches such as ANTs [2], Elastix [15], and NiftiReg [19] employ optimization-based frameworks. This iterative search for the optimal transformation makes these methods inherently slow, especially when dealing with large datasets or high-resolution images [13]. In light of these challenges, there has been a growing interest in transitioning toward learning-based methods. Specifically, deep learning methods are significantly faster during inference and currently provide state-of-the-art performance for 3D image registration [6,12].

Learning-based approaches for image registration can generally be divided into two main categories: supervised and unsupervised methods. *Supervised methods* (e.g., [21,24]) require estimates of deformation fields derived from traditional optimization-based approaches, the acquisition of which can be prohibitively costly for extensive datasets. Moreover, the efficacy of supervised approaches is contingent upon the availability of high-quality deformation fields for supervised training, with their performance capped by the accuracy of the method used to obtain these fields. In contrast, *unsupervised methods* do not require deformation fields and use image similarity as a self-supervised signal to train a registration network. Most unsupervised 3D registration methods (e.g., [4,6,12]) are trained to produce a 3D deformation field that is then used to transform (or warp) the moving image. Loss (L1 or L2) between the warped moving image and the fixed image is used to train the network. With sufficient data and training time, the model learns to produce realistic deformation fields that outperform optimization-based methods in both accuracy and inference speed [6,12].

Learning-based registration methods predominantly rely on convolutional architectures (e.g., [4,14,12]), using U-Net-based architectures to generate the deformation fields. However the effectiveness of convolutional layers for registration tasks can be compromised due to their limited receptive fields that hinder capturing global context [6] and their increased susceptibility to vanishing gradients as network depth grows to enhance learning capacity [8]. Since the advent of Vision Transformers [22], transformer-based architectures have shown superior performance across various tasks, such as classification, segmentation, and registration [22,23,6,7], thanks to their modeling capabilities. In particular, they offer promising mitigations to CNN limitations. Specifically, transformers leverage global contextual information through self-attention mechanisms and provide more stable gradient flow across layers via techniques such as layer normalization and skip connections that are integral to transformers design [6].

Despite their success, transformers' advantages come at the expense of a significant increase in memory footprint and parameter count, requiring approximately 10 to 20 times more parameters than convolutional counterparts [6].

Specifically, existing transformer-based registration methods, including TransMorph [6], the current state-of-the-art transformer-based model for medical image registration, encounter *three* main significant limitations that compromise their efficiency and overall performance. *Firstly*, windowed attention approaches (e.g., the Swin transformer [17] backbone used in TransMorph [6]) optimize computational efficiency through local attention and shifted windows, enhancing interactions between adjacent windows. However, this limits global context capture, particularly in shallow layers, due to within-window constraints compared to methods that interact globally. *Secondly*, multi-head self-attention (MHSA) often learns redundant features across heads, suggesting that models could be simplified by encouraging diverse feature learning, thereby reducing computational redundancy without sacrificing accuracy. *Lastly*, the granularity of tokenization significantly impacts registration accuracy; smaller tokens capture finer details for higher accuracy but increase computational and memory requirements, potentially leading to overfitting.

In this paper, we propose EfficientMorph, a novel transformer-based framework for unsupervised 3D image registration that addresses the aforementioned challenges. We introduce a “plane attention” mechanism inspired by anatomical views (coronal, sagittal, and axial) to enhance the balance between local and global feature recognition. To reduce computational redundancy, we employ cascaded group attention [16] where each head receives only a portion of the complete feature set that is cascaded to the previous head’s representation via feature additions. Furthermore, we propose Hi-Res tokenization to reduce the model’s complexity within the encoded representation by merging neighboring tokens in a high-resolution feature space. The integration of plane cascaded group attention with Hi-Res tokenization positions EfficientMorph as a highly parameter-efficient registration architecture (see Figure 2A).

The main contributions of this paper are:

- A novel attention module for 3D registration that focuses on attention across the coronal (xy), sagittal (yz), or Axial (zx) planes within a single transformer block.
- A Hi-Res tokenization mechanism to encode high-resolution features and use cascaded group attention [16] to learn less redundant features without compromising computational efficiency.
- A new parameter-efficient architecture that achieved performance comparable to existing methods within a margin of ± 0.05 dice score, even surpassing the state-of-the-art performance on one dataset while having $\sim 16\text{-}27\times$ fewer parameters (Figure 2A) and faster convergence ($\sim 5\times$).

2 Methods

Given a 3D volume represented as $\mathbf{A} \in \mathbb{R}^{H \times W \times D}$, where H , W , and D denote the height, width, and depth dimensions, respectively. Strided convolutions are used in the patch embedding layer (with stride s) to project \mathbf{A} into a high-dimensional feature space, resulting in $\mathbf{A}' \in \mathbb{R}^{H' \times W' \times D' \times C}$, where C is the

embedding dimension, $(H', W', D') = (\frac{H}{s'}, \frac{W}{s'}, \frac{D}{s'})$. The resulting feature space is tokenized to train the downstream transformer layers. In the sequel, we describe the Hi-Res tokenization, plane attention mechanism, and cascaded group attention (CGA) of the proposed EfficientMorph and illustrate them in Figure 1.

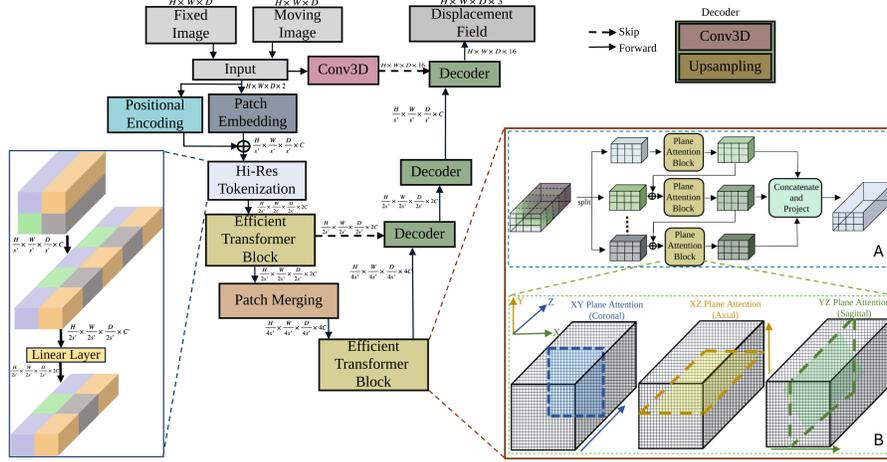


Fig. 1. EfficientMorph architecture. EfficientMorph with CGA utilizes block A whereas without CGA utilizes plane attention mechanism on the whole volume as shown in block B only. Plane Attention and CGA Architecture are highlighted from the transformer block. We use different numbers and types of plane attentions (xy, yz , or zx planes) for each block in the transformer backbone (Table 1).

2.1 Hi-Res Tokenization

For a fixed embedding dimension C , using each voxel of a 3D volume of N -voxels for tokenization would create N tokens, where $N = H \times W \times D$. Voxel-level tokenization results in attention matrices of more than a trillion parameters with a complexity of $\mathcal{O}(N^2)$. Transformer architectures often rely on s -strided convolutions (e.g., $s = 4$ [6]) for volume tokenization and patch embedding, trading off computational complexity, which is now $\mathcal{O}\left(\left(\frac{N}{s}\right)^2\right)$, at the cost of detailed features. However, fine-grained spatial information is critical for medical segmentation and registration tasks, which may be lost due to strided convolutions.

We propose a novel *Hi-Res tokenization* strategy that uses a smaller stride ($s' < s$) within the embedding layer, thereby creating high-resolution tokenized features. These tokens undergo positional encoding and are merged by grouping and concatenating the features of adjacent non-overlapping $d \times d \times d$ voxel token blocks resulting in $N' = \frac{H}{s'} \times \frac{W}{s'} \times \frac{D}{s'}$ tokens with an embedding dimension of $C' = C \times d^3$ (with $d = \frac{s}{s'}$). Then, C' is projected into a linear layer to attain

a reduced dimension of $C \times d$, as shown in Hi-Res tokenization block in Figure 1. This approach enables the use of tokens from higher resolution and reduces complexity by a factor of $\frac{d^3}{s}$.

2.2 Plane Cascaded Group Attention Mechanism

Plane Attention. Despite the use of Hi-Res tokenization, the number of tokens generated from each volume remains high. Running full attention on these tokens, while feasible, demands considerable computational resources. To address this challenge, we introduce a novel attention framework called plane attention. Instead of performing full 3D attention on all tokens, this method utilizes attention along coronal (xy), sagittal (yz), or Axial (zx) planes, as shown in Figure 1. Although attention confines focus to a specific plane, EfficientMorph achieves volume attention by sequentially employing different attention combinations xy followed by yz or zx .

$$\text{Attn}(\mathbf{A}'_{dim}) = \text{softmax} \left(\frac{\mathbf{Q}_{dim} \mathbf{K}_{dim}^T}{\sqrt{d_k}} \right) \mathbf{V}_{dim} \quad (1)$$

Here, $dim \in \{xy, yz, zx\}$, \mathbf{A}'_{dim} can be represented as $\mathbf{A}'_{xy} \in \mathbf{R}^{H' \times W' \times C}$, $\mathbf{A}'_{yz} \in \mathbf{R}^{W' \times D' \times C}$ and $\mathbf{A}'_{zx} \in \mathbf{R}^{D' \times H' \times C}$ for xy , yz , and zx planes, respectively. By decomposing the 3D attention into 2D plane attention, the proposed attention mechanism significantly reduces the parameter count while preserving the ability to capture essential volumetric features necessary for registration.

Cascaded Group Attention (CGA). MHSA is a key element of transformer models, allowing simultaneous focus on diverse input aspects. However, research has shown that MHSA modules tend to learn redundant information from the dataset [5,16]. To make our proposed architecture learn better feature representations, we incorporated CGA as shown in Figure 1B, and it has effectively minimized redundancy within the feature space without sacrificing performance.

Let \mathbf{A}' be split into h groups of tokens (i.e., $\mathbf{A}' = [\mathbf{A}'_1, \mathbf{A}'_2, \dots, \mathbf{A}'_i, \dots, \mathbf{A}'_h]$ and $1 \leq i \leq h$) where h is the number of heads. CGA can be formally expressed as:

$$\mathbf{A}'_1 = \text{Attn}(\mathbf{A}'_1) \quad (2)$$

$$\mathbf{A}'_{i+1} = \mathbf{A}'_{i+1} + \text{Attn}(\mathbf{A}'_i), \quad 1 \leq i \leq h-1 \quad (3)$$

$$\text{Attn}(\mathbf{A}') = \text{Concat}[\text{Attn}(\mathbf{A}'_i)]_{i=1}^h, \quad (4)$$

Here, each segment of the input features is denoted by \mathbf{A}'_i for the i -th segment. The enriched feature set for the subsequent head, \mathbf{A}'_{i+1} , is derived by incorporating the output from the current head, \mathbf{A}'_i , into it. The final output, denoted as $\text{Attn}(\mathbf{A}')$, is formed by concatenating the attention outputs from all heads. This layered approach ensures that the input for each head is a combination of its specific feature segment and the aggregated insights from preceding heads.

3 Results and Discussion

3.1 Datasets and Preprocessing

OASIS Brain MRI. We evaluated EfficientMorph on the publicly available dataset OASIS [18], obtained from the Learn2Reg challenge [9] for inter-patient registration and pre-processed from [10]. It has a total of 451 brain T2 MRI images. Among these, 394, 19, and 38 scans are used for training, validation, and testing, respectively. Registration accuracy is reported by performing evaluation of corresponding segmentation masks for 35 anatomical structures.

Atlas-to-Patient Brain MRI (IXI). We additionally evaluated the proposed model on IXI dataset that contains 600 MRI scans. Among these, 576 T1-weighted brain MRI images were employed as fixed images, while the moving image utilized for this task was an atlas brain MRI [14]. The dataset was partitioned into training, validation, and test sets, comprising 403, 58, and 115 volumes, respectively. Evaluation was performed on corresponding segmentation masks for 29 anatomical structures.

Implementation Details. EfficientMorph was trained on NVIDIA A100 GPUs with 40GB RAM and a batch size of 1. We used the same splits for both datasets as the existing works [6,12]. We limited training epochs to 100 to prioritize parameter efficiency and quick convergence within resource limits. We used the Adam optimizer with a learning rate of $5e-4$ for OASIS and $3e-4$ for IXI. We used a cosine annealing schedule for OASIS and stepLR for IXI. The IXI dataset was augmented by flipping in random directions while training, as done by baselines. We tested different variants: EfficientMorph-11, featuring one transformer block at each stage, and EfficientMorph-23, with two transformer blocks at stage 1 and three transformer blocks at stage 2. The corresponding plane attentions used in the variants are shown in Table 1.

Table 1. EfficientMorph Variants. EfficientMorph-AB refers to a two-stage model with A-blocks in stage 1 and B-blocks in stage 2.

Variants	Planes
EfficientMorph-11	(xy, yz)
EfficientMorph-23	(xy-yz, xy-yz-zx)

3.2 Results

The results on the OASIS dataset are shown in Table 2. We compare EfficientMorph with state-of-the-art convolutional-based methods, including VoxelMorph-H [4] and Fourier-Net [12], as well as different variants of TransMorph [6], such as TransMorph-Tiny, TransMorph, and TransMorph-L. Our proposed variants, EfficientMorph-11 and EfficientMorph-23, exhibit comparable performance albeit with fewer parameters (Figure 2A) and achieve faster convergence (Figure 2B). EfficientMorph variants utilizing the Hi-Res tokenization strategy with a stride of 2 demonstrate a parameter count similar to those employing a stride of 4, with only a marginal increase in computational overhead. Among these

Table 2. OASIS Results. Mean average dice score and standard deviation are evaluated on 35 segmented anatomies in OASIS. * indicates the performance numbers taken from TransMorph and Fourier-Net; for all others, we ran these baselines on our system for fair comparison. ‘stride’ and ‘C’ are the strides and embedding dimensions. ‘Multi-Add’ denotes the number of Multiply add operations for a forward pass.

Methods	stride	C	Epochs	Param(M)	Multi-Add	Infer(sec)	Dice Score
VoxelMorph-H*[4]	-	-	-	-	3656.2	-	0.847±0.014
Fourier-Net*[12]	-	-	-	4.19	169.07	-	0.847±0.013
TransMorph-Tiny[6]	4x4x4	6	100	0.24	11.36	0.0161	0.80±0.056
TransMorph[6]	4x4x4	96	100	46.5	251.50	0.0992	0.8486±0.0137
TransMorph[6]	4x4x4	96	500	46.5	251.50	0.0998	0.858±0.0143
TransMorph-L* [6]	4x4x4	128	500	108.11	416.30	-	0.862±0.014
EfficientMorph-11	4x4x4	96	100	1.8	171.14	0.0610	0.8408±0.0127
EfficientMorph-23	4x4x4	96	100	2.8	171.14	0.0810	0.8458±0.0127
EfficientMorph-11	2x2x2	96	100	1.7	1359.85	0.5585	0.8623±0.0133
EfficientMorph-23	2x2x2	96	100	2.8	1359.85	0.9179	0.8671±0.0135
EfficientMorph-11 (CGA)	2x2x2	96	100	1.6	1359.85	0.5525	0.8506±0.0136
EfficientMorph-11	2x2x2	24	100	1.2	92.12	0.0834	0.8403±0.0114
EfficientMorph-23	2x2x2	24	100	2.25	92.12	0.0959	0.843±0.01360

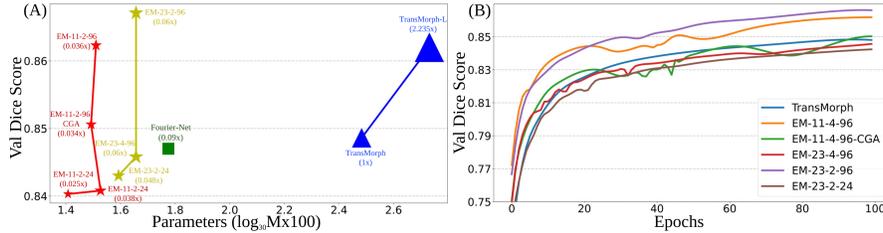


Fig. 2. OASIS quantitative results. The proposed variants are formatted as EfficientMorph-11-stride-C and EfficientMorph-23-stride-C. (A) Comparison of parameter count in millions(M) and Dice scores between the proposed Variants and baselines. (B) Dice score curves of EfficientMorph variants as a function of epochs.

variants, EfficientMorph-23 achieves the highest Dice score with only 2.8M parameters (16 times fewer parameters than TransMorph) and outperforms all the compared baselines, even TransMorph-L, which has more than 100M parameters. With even fewer parameters, Efficient-11 and its CGA variant outperform all other baselines. The results indicate that leveraging the plane attentions, with Hi-Res tokenization and CGA, leads to fewer parameters and better performance. It is observed that the ability to learn diverse features is more beneficial when we have fewer parameters because the model attempts to maximize its learning capacity with the parameter count available. Figure 2A shows that reducing 75% of embedding dimensions for “EfficientMorph-23” and “EfficientMorph-11” results in a significant parameter reduction (19.6% and 29.4%, respectively) with a small performance decrease (2.8% and 2.55%, respectively). This indicates that altering embedding dimensions in Hi-Res tokenization can substantially reduce parameters with a limited impact on performance. Supplementary Figure 5 illustrates the performance comparison between the EfficientMorph variants and the baseline on different brain MR substructures, demonstrating notable enhancements for the proposed models.

The accuracy curves in Figure 2B clearly show that TransMorph learns quickly in a few initial epochs but then slowly saturates to the final performance, whereas all EfficientMorph variants slowly and steadily converge to higher dice scores. EfficientMorph starts to outperform TransMorph by a significant margin as early as 10 epochs. Supplementary Figure 3 shows that EfficientMorph’s warped segmentations are qualitatively better than those of TransMorph.

Results of the IXI dataset are presented in Table 3. EfficientMorph outperforms traditional optimization-based methods such as SyN, NiftiReg, and various convolutional-based approaches such as VoxelMorph-H [4] and CycleMorph [14] by a significant margin. EfficientMorph variants EfficientMorph-11 and EfficientMorph-23 with 4x4x4 strides achieve comparable performance (within ± 0.003) with less than 3 million parameters compared to TransMorph’s 46 million parameters and 5 \times fewer epochs. Variants employing the Hi-Res tokenization technique with a stride 2 do not perform well for IXI. However, the ablations experiment with fewer embedding dimensions (C=24) improved the performance of 0.7317 to TransMorph’s 0.7293 at 100 epochs, achieving similar accuracy as Fourier-Net-s and has better inference speed than all other baselines. If trained for a longer period (> 100 epochs), EfficientMorph will probably be as accurate as TransMorph (maybe even higher), but this is left for future experiments. Accuracy vs epochs curves shown in supplementary Figure 4 indicate that most EfficientMorph variants outperform TransMorph in initial epochs, but then performance tends to saturate. Qualitative segmentations for IXI dataset, shown in supplementary Figure 7, show that EfficientMorph produces results of similar quality to TransMorph. For different substructures, EfficientMorph performs on par with the baseline, as shown in supplementary Figure 6.

Table 3. IXI Results. Mean average dice score and standard deviation are evaluated on 29 segmented anatomies in IXI. * indicates the performance numbers taken from TransMorph and Fourier-Net; for all others, we ran these baselines on our system for fair comparison. ‘stride’ and ‘C’ are the strides and channel layer for initial embedding layer. ‘Multi-Add’ is the number of Multiply add operations needed for a forward pass.

Methods	stride	C	Epochs	Param(M)	Multi-Add	Infer(sec)	Dice Score	
							Val	Test
SyN*	-	-	-	-	-	-	-	0.645 \pm 0.152
NiftiReg*	-	-	-	-	-	-	-	0.645 \pm 0.167
voxelMorph-1*[4]	-	-	-	0.3	-	-	-	0.548 \pm 0.317
cycleMorph*[14]	-	-	-	-	966.9	-	-	0.528 \pm 0.321
Fourier-Net-s[12]	-	-	200	1.05	43.82	0.318	0.729 \pm 0.024	0.730 \pm 0.025
Fourier-Net-s[12]	-	-	1000	1.05	43.82	0.318	0.735 \pm 0.026	0.736 \pm 0.027
Fourier-Net[12]*	-	-	1000	4.19	169.07	0.342	-	0.763\pm0.129
TransMorph-Tiny*[6]	4x4x4	6	500	0.24	122.3	-	0.545 \pm 0.180	0.543 \pm 0.180
TransMorph[6]	4x4x4	96	100	46.7	686.80	0.2044	0.7293 \pm 0.029	0.7324 \pm 0.0314
TransMorph[6]	4x4x4	96	500	46.7	686.80	0.2044	0.7405 \pm 0.0283	0.7408 \pm 0.0299
TransMorph-L[6]*	4x4x4	128	500	108.34	1084.9	-	0.753 \pm 0.130	0.754\pm0.128
EfficientMorph-11	4x4x4	96	100	2.01	577.29	0.1567	0.7233 \pm 0.0305	0.7224 \pm 0.0324
EfficientMorph-23	4x4x4	96	100	3.04	577.29	0.1749	0.7233 \pm 0.0303	0.7298 \pm 0.0322
EfficientMorph-11	2x2x2	96	100	1.7	1359.86	0.5716	0.6739 \pm 0.0322	0.6749 \pm 0.0323
EfficientMorph-23	2x2x2	96	100	2.8	1359.86	0.9365	0.7159 \pm 0.0307	0.7174 \pm 0.0330
EfficientMorph-11(CGA)	2x2x2	96	100	1.6	1359.86	0.5490	0.6843 \pm 0.03332	0.6859 \pm 0.0330
EfficientMorph-11	2x2x2	24	100	2.02	576.3	0.1715	0.7206 \pm 0.0315	0.7210 \pm 0.0337
EfficientMorph-23	2x2x2	24	100	3.0	576.3	0.1906	0.7312\pm0.0298	0.7317\pm0.0320

4 Conclusion and Future Work

We propose EfficientMorph, a parameter-efficient transformer-based architecture for unsupervised 3D deformable image registration. EfficientMorph uses a novel attention plane attention mechanism. EfficientMorph attends to 3D volumetric features by sequentially placing different plane attention blocks xy followed by yz or zx , thus attending to features along all three axes. Additionally, we propose a Hi-Res tokenization strategy to increase the feature resolution while maintaining computational complexity. To mitigate the redundant feature learned by transformer layers, we use cascaded group attention (CGA). Evaluations of two datasets demonstrate that EfficientMorph can achieve state-of-the-art results with a considerably lower parameter count ($\sim 16\text{-}27\times$). As future work, we plan to explore other attention mechanisms that can be paired with EfficientMorph to further reduce the computational overhead. Furthermore, reducing decoder complexity can further improve the efficiency and efficacy of EfficientMorph.

Acknowledgements

The National Science Foundation supported this work under grant number NSF2217154. The content is solely the authors' responsibility and does not necessarily represent the official views of the National Science Foundation.

References

1. Alam, F., Rahman, S.U., Ullah, S., Gulati, K.: Medical image registration in image guided surgery: Issues, challenges and research opportunities. *Biocybernetics and Biomedical Engineering* **38**(1), 71–89 (2018)
2. Avants, B.B., Tustison, N., Song, G., et al.: Advanced normalization tools (ants). *Insight j* **2**(365), 1–35 (2009)
3. Azam, M.A., Khan, K.B., Ahmad, M., Mazzara, M.: Multimodal medical image registration and fusion for quality enhancement. *Computers, Materials & Continua* **68**(1), 821–840 (2021)
4. Balakrishnan, G., Zhao, A., Sabuncu, M.R., Guttag, J., Dalca, A.V.: Voxelmorph: a learning framework for deformable medical image registration. *IEEE transactions on medical imaging* **38**(8), 1788–1800 (2019)
5. Bian, Y., Huang, J., Cai, X., Yuan, J., Church, K.: On attention redundancy: A comprehensive study. In: *Proceedings of the 2021 conference of the north american chapter of the association for computational linguistics: human language technologies*. pp. 930–945 (2021)
6. Chen, J., Frey, E.C., He, Y., Segars, W.P., Li, Y., Du, Y.: Transmorph: Transformer for unsupervised medical image registration. *Medical image analysis* **82**, 102615 (2022)
7. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929* (2020)

8. Hanin, B.: Which neural net architectures give rise to exploding and vanishing gradients? *Advances in neural information processing systems* **31** (2018)
9. Hering, A., Hansen, L., Mok, T.C., Chung, A.C., Siebert, H., Häger, S., Lange, A., Kuckertz, S., Heldmann, S., Shao, W., et al.: Learn2reg: comprehensive multi-task medical image registration challenge, dataset and evaluation in the era of deep learning. *IEEE Transactions on Medical Imaging* **42**(3), 697–712 (2022)
10. Hoopes, A., Hoffmann, M., Fischl, B., Guttag, J., Dalca, A.V.: Hypermorph: Amortized hyperparameter learning for image registration. In: *Information Processing in Medical Imaging: 27th International Conference, IPMI 2021, Virtual Event, June 28–June 30, 2021, Proceedings 27*. pp. 3–17. Springer (2021)
11. Huang, Z., Liu, J., Fan, X., Liu, R., Zhong, W., Luo, Z.: Reconet: Recurrent correction network for fast and efficient multi-modality image fusion. In: *European Conference on Computer Vision*. pp. 539–555. Springer (2022)
12. Jia, X., Bartlett, J., Chen, W., Song, S., Zhang, T., Cheng, X., Lu, W., Qiu, Z., Duan, J.: Fourier-net: Fast image registration with band-limited deformation. In: *Proceedings of the AAAI Conference on Artificial Intelligence*. vol. 37, pp. 1015–1023 (2023)
13. Kataria, T., Rajamani, S., Ayubi, A.B., Bronner, M., Jedrzekiewicz, J., Knudsen, B.S., Elhabian, S.Y.: Automating ground truth annotations for gland segmentation through immunohistochemistry. *Modern Pathology* **36**(12), 100331 (2023)
14. Kim, B., Kim, D.H., Park, S.H., Kim, J., Lee, J.G., Ye, J.C.: Cyclemorph: cycle consistent unsupervised deformable image registration. *Medical image analysis* **71**, 102036 (2021)
15. Klein, S., Staring, M., Murphy, K., Viergever, M.A., Pluim, J.P.: Elastix: a toolbox for intensity-based medical image registration. *IEEE transactions on medical imaging* **29**(1), 196–205 (2009)
16. Liu, X., Peng, H., Zheng, N., Yang, Y., Hu, H., Yuan, Y.: Efficientvit: Memory efficient vision transformer with cascaded group attention. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 14420–14430 (2023)
17. Liu, Z., Lin, Y., Cao, Y., Hu, H., Wei, Y., Zhang, Z., Lin, S., Guo, B.: Swin transformer: Hierarchical vision transformer using shifted windows. In: *Proceedings of the IEEE/CVF international conference on computer vision*. pp. 10012–10022 (2021)
18. Marcus, D.S., Wang, T.H., Parker, J., Csernansky, J.G., Morris, J.C., Buckner, R.L.: Open Access Series of Imaging Studies (OASIS): Cross-sectional MRI Data in Young, Middle Aged, Nondemented, and Demented Older Adults. *Journal of Cognitive Neuroscience* **19**(9), 1498–1507 (09 2007). <https://doi.org/10.1162/jocn.2007.19.9.1498>, <https://doi.org/10.1162/jocn.2007.19.9.1498>
19. Modat, M., Ridgway, G.R., Taylor, Z.A., Lehmann, M., Barnes, J., Hawkes, D.J., Fox, N.C., Ourselin, S.: Fast free-form deformation using graphics processing units. *Computer methods and programs in biomedicine* **98**(3), 278–284 (2010)
20. Rigaud, B., Simon, A., Castelli, J., Lafond, C., Acosta, O., Haigron, P., Cazoulat, G., de Crevoisier, R.: Deformable image registration for radiation therapy: principle, methods, applications and evaluation. *Acta Oncologica* **58**(9), 1225–1237 (2019)
21. Sokooti, H., De Vos, B., Berendsen, F., Lelieveldt, B.P., Išgum, I., Staring, M.: Nonrigid image registration using multi-scale 3d convolutional neural networks. In: *Medical Image Computing and Computer Assisted Intervention- MICCAI 2017*:

- 20th International Conference, Quebec City, QC, Canada, September 11-13, 2017, Proceedings, Part I 20. pp. 232–239. Springer (2017)
22. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, L., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017)
 23. Xie, E., Wang, W., Yu, Z., Anandkumar, A., Alvarez, J.M., Luo, P.: Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems* **34**, 12077–12090 (2021)
 24. Yang, X., Kwitt, R., Niethammer, M.: Fast predictive image registration. In: *Deep Learning and Data Labeling for Medical Applications: First International Workshop, LABELS 2016, and Second International Workshop, DLMIA 2016, Held in Conjunction with MICCAI 2016, Athens, Greece, October 21, 2016, Proceedings 1*. pp. 48–57. Springer (2016)

5 Supplementary

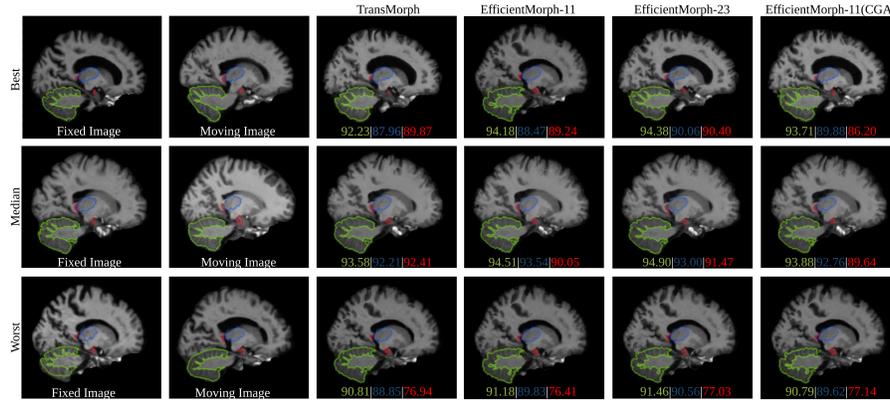


Fig. 3. OASIS qualitative results. Comparison among the best, median, and worst output of TransMorph with the variants of the proposed method. Here, EfficientMorph-23 and EfficientMorph-11 are the different variants with $2 \times 2 \times 2$ stride size and 96 embedded dimension; CGA means variants with cascaded group attention.

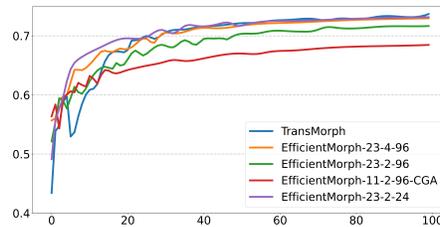


Fig. 4. Dice scores as a function of number of epochs (IXI).

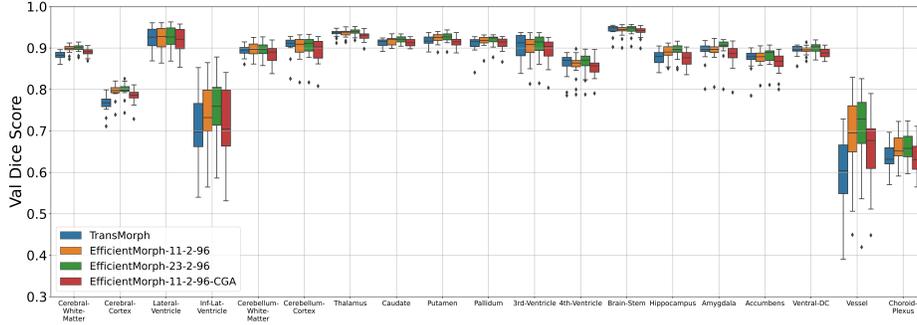


Fig. 5. OASIS boxplot. Quantitative comparison of the proposed models with TransMorph showing dice scores for 19 anatomical substructures.

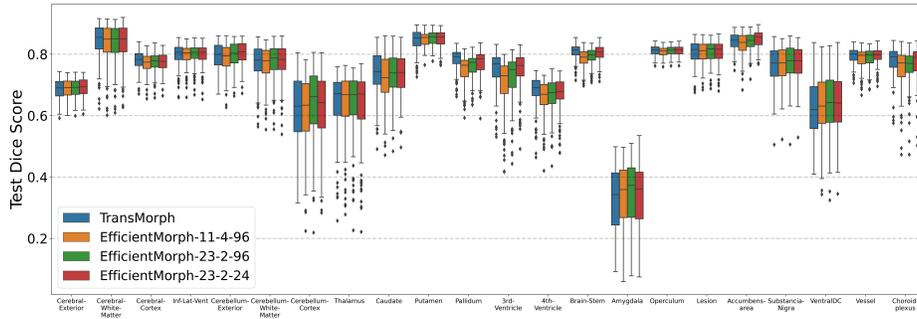


Fig. 6. IXI boxplot. Quantitative comparison of the proposed models with TransMorph showing dice scores for 22 anatomical substructures.

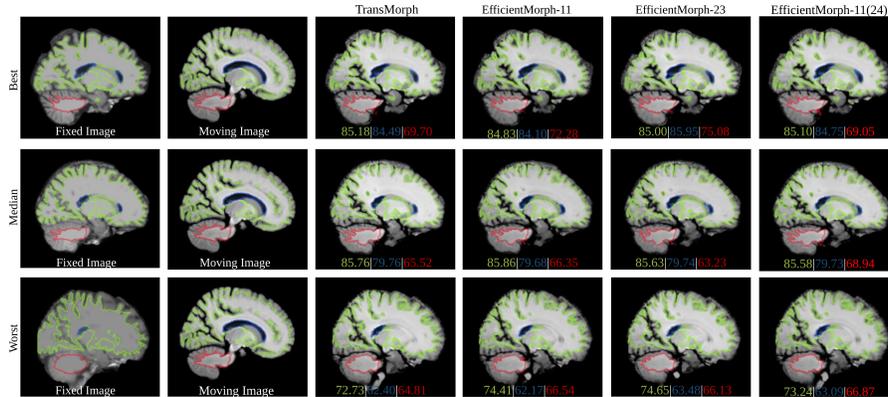


Fig. 7. IXI qualitative results. Comparison among the best, median, and worst output of TransMorph with the variants of the proposed method. EfficientMorph-23 and EfficientMorph-11 are the different variants with 4x4x4 stride size and 96 embedded dimensions; EfficientMorph-11(24) has 24 embedding dimensions.