

Random Token Fusion for Multi-View Medical Diagnosis

Jingyu Guo, Christos Matsoukas, Fredrik Strand, Kevin Smith



In multi-view medical diagnosis, deep learning-based models often fuse information from different imaging perspectives to improve diagnostic performance. However, existing approaches are prone to overfitting and rely heavily on view-specific features, which can lead to trivial solutions. In this work, we introduce Random Token Fusion (RTF), a novel technique designed to enhance multi-view medical image analysis using vision transformers. By integrating randomness into the feature fusion process during training, RTF addresses the issue of overfitting and enhances the robustness and accuracy of diagnostic models without incurring any additional cost at inference. We validate our approach on standard mammography and chest X-ray benchmark datasets. Through extensive experiments, we demonstrate that RTF consistently improves the performance of existing fusion methods, paving the way for a new generation of multi-view medical foundation models.

Problem description

Physicians routinely employ multi-view analysis in diagnostic procedures. Images gathered at various angles can unveil details that may be obscured in a single view, enhancing the precision of the diagnosis. It stands to reason that foundation models for medical image analysis could similarly improve their diagnostic accuracy by integrating information from multiple views. However, **multi-view models often to overfit to the most informative view** [1,2].

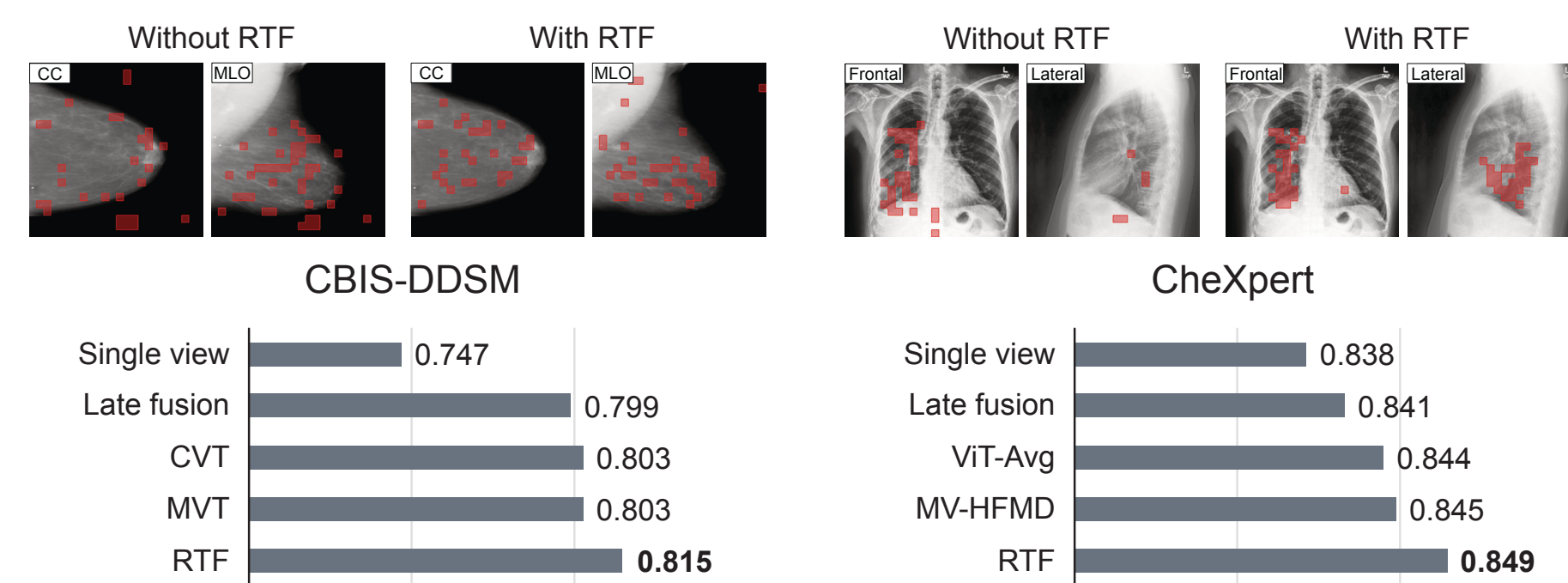


Figure 1: Illustration of the overfitting problem in multi-view medical diagnosis. The model's attention becomes overly focused on one of the two available views, resulting in an incomplete interpretation of the case. In this example (top), model attention in the MLO view dominates over the CC view in CBIS-DDSM (top left), and the frontal view over the lateral view in CheXpert (top right). RTF encourages the model to better utilize information from both views, resulting in balanced attention between both views and increased performance (bottom).

The neglect of complementary information in other views [1], can be particularly problematic in medical image analysis, where each view can provide unique diagnostic insights [3,4,5,6].

Multi-view ViTs with Random Token Fusion

We propose a solution that can **enhance existing multi-view fusion strategies**, which we call **Random Token Fusion (RTF)**.

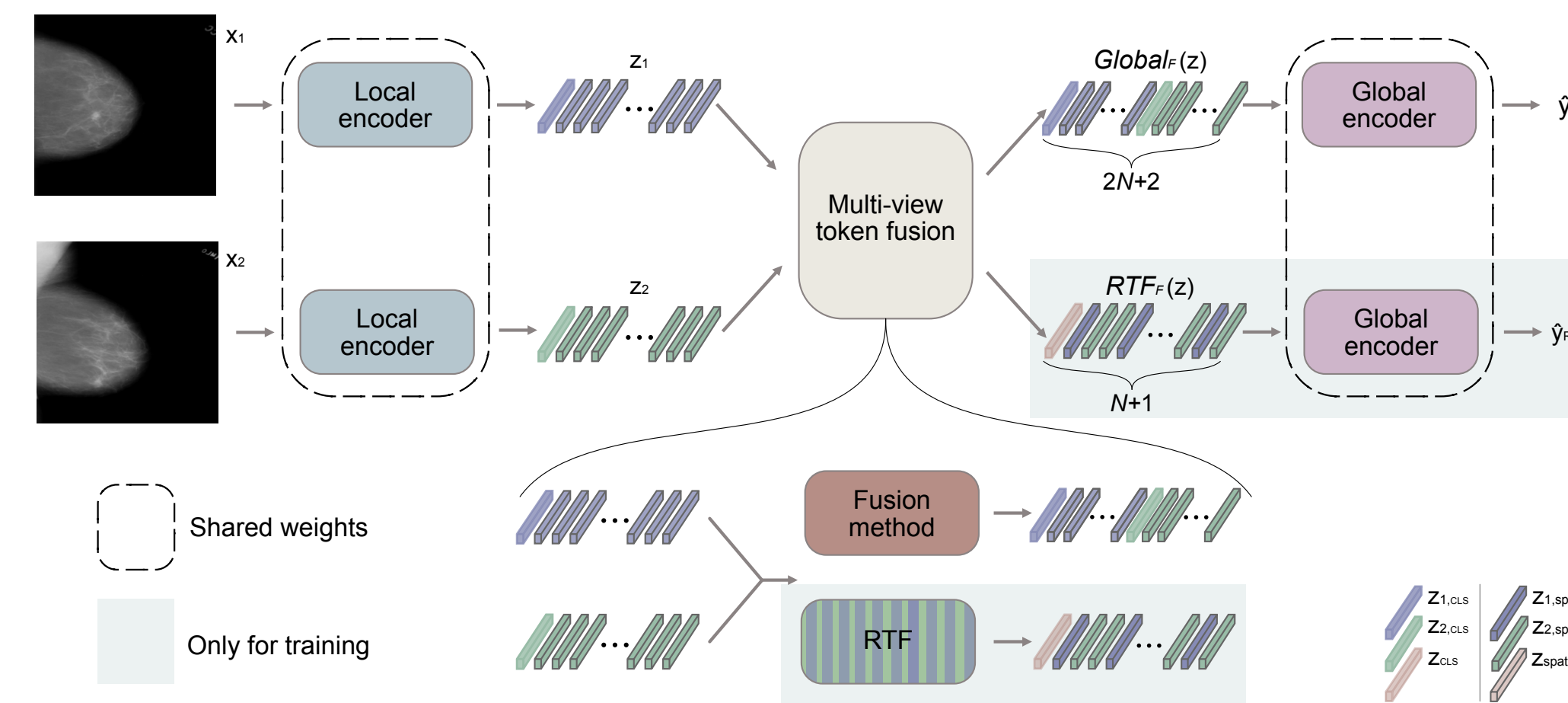


Figure 2: Multi-view ViTs with Random Token Fusion (RTF). RTF utilizes a local encoder to generate representations of different views, followed by a token fusion module. This module divides the feature fusion into two distinct branches. One branch uses some strategy to merge all tokens from both views, while the other one randomly drops spatial tokens from each view before mixing them. The fused tokens are processed by a global encoder, which produces two types of predictions: one for the global tokens and one for the RTF tokens. During training, the loss for both branches is minimized towards the same task. After training, RTF tokens are not generated, they are merged using the model's fusion method and passed to the global encoder for inference.

Integrating information occlusion and mixing into the training process has been a proven method to combat overfitting and improve robustness [7,8,9]. RTF randomly fuses tokens from different views, introducing variability in the fused representation, which acts as a regularizer. This compels the network to capture dependencies between patches originating from different views, **preventing the model from overfitting to view-specific features**.

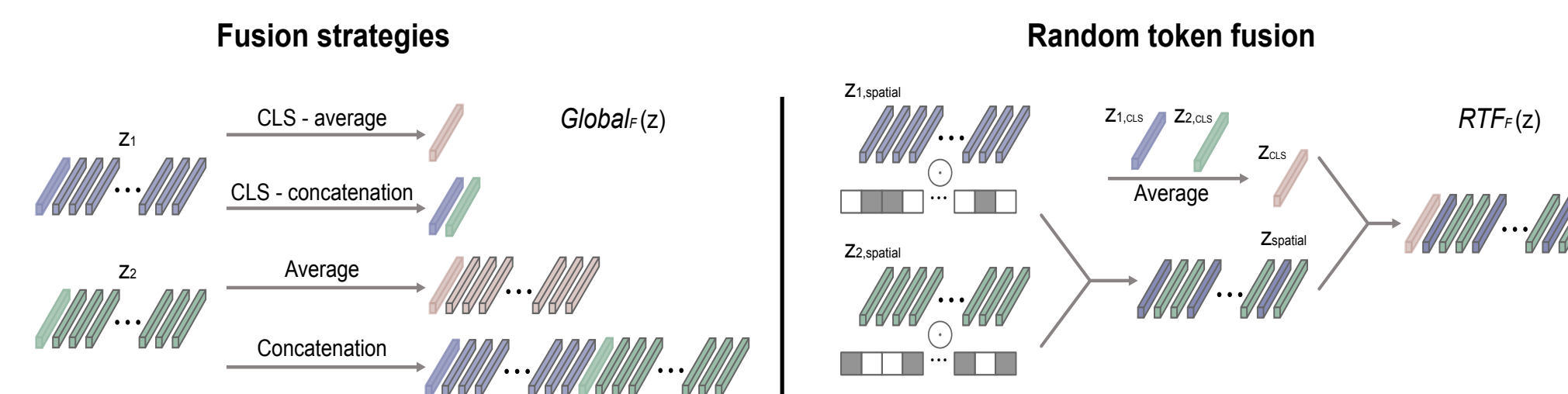


Figure 3: Illustration of different fusion strategies. (Left) Common fusion strategies to fuse the features (tokens) of different views in ViTs. (Right) The proposed random token fusion (RTF) strategy. In RTF, we randomly drop spatial tokens from both images and combine the remaining ones, augmenting the representations during training.

RTF can be **seamlessly integrated with existing multi-view fusion strategies** for vision transformers (ViTs), enriching an existing model's feature space without requiring any modification to the inference process. By **incorporating randomness** into the token fusion process, RTF also **encourages the model to learn robust and generalized features** from all views, ensuring that the fused representation captures the most informative features.

Results

Multiple views vs. single view

We assess the benefits of multi-view ViTs and train models on single and multiple views (Table 1). Using two views results in improved AUC performance compared to single-view models with the same capacity. This highlights the importance of encouraging models to utilize both perspectives in medical diagnosis

Table 1: The effect of using only a single view, multiple views with late fusion, and multiple views with RTF on CBIS-DDSM [10] (left) and CheXpert [6] (right).

View	DDSM, AUC \uparrow	View	CheXpert, AUC \uparrow
Only CC	0.730 \pm 0.004	Only Frontal	0.838 \pm 0.003
Only MLO	0.747 \pm 0.022	Only Lateral	0.832 \pm 0.002
Late fusion	0.799 \pm 0.008	Late fusion	0.841 \pm 0.001
Fusion w/ RTF	0.815 \pm 0.001	Fusion w/ RTF	0.849 \pm 0.001

RTF enhances multi-view fusion

Training with RTF consistently improves performance across all configurations (Table 2). The extent of improvement varies with the dataset and model size. CBIS-DDSM appears to gain more from RTF, particularly for larger ViT variants. We hypothesize that this is due to the regularization effects of RTF and the smaller size of the dataset, as higher-capacity models are more prone to overfitting

Table 2: AUC performance on CBIS-DDSM (left) and CheXpert (right), showing the effect of using multiple views with and without RTF for different model sizes and fusion strategies.

Method	RTF Used	ViT Tiny	ViT Small	ViT Base
Average	No	0.798 \pm 0.003	0.803 \pm 0.008	0.813 \pm 0.004
	Yes	0.802 \pm 0.001	0.809 \pm 0.002	0.825 \pm 0.005
CLS _{cut}	No	0.796 \pm 0.002	0.802 \pm 0.006	0.814 \pm 0.007
	Yes	0.801 \pm 0.001	0.811 \pm 0.008	0.826 \pm 0.004
Concat	No	0.798 \pm 0.003	0.803 \pm 0.003	0.814 \pm 0.004
	Yes	0.802 \pm 0.003	0.815 \pm 0.001	0.830 \pm 0.002

RTF outperforms SOTA fusion methods

RTF outperforms other fusion methods on both datasets, showing its efficacy in multi-view medical diagnosis. Notably, RTF can be used in conjunction with transformer-based methods, such as ViT-Average [11] and MVT [12,13], for further enhanced performance.

Table 3: Comparison vs. SOTA methods on CBIS-DDSM (left) and CheXpert (right).

Method	CBIS-DDSM	Method	CheXpert
ResNet50	0.724 \pm 0.007	MVC-NET	0.813 \pm 0.005
Shared ResNet	0.735 \pm 0.014	MVCNN	0.815 \pm 0.004
PHResNet50	0.739 \pm 0.004	CVT	0.834 \pm 0.002
MVT	0.803 \pm 0.003	MVT	0.843 \pm 0.004
CVT	0.803 \pm 0.007	ViT-Average	0.844 \pm 0.004
ViT-Average	0.803 \pm 0.008	MV-HFMD	0.845 \pm 0.002
RTF	0.815 \pm 0.001	RTF	0.849 \pm 0.001

References

- [1] Wu, N., et al. "Improving the ability of deep neural networks to use information from multiple views in breast cancer screening." Medical Imaging with Deep Learning. PMLR, 2020.
- [2] Wang, W., et al. "What makes training multi-modal classification networks hard?." Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2020.
- [3] Feigin, D. S., et al. "Lateral chest radiograph: a systematic approach." Academic radiology 17.12 (2010): 1560-1566.
- [4] Hashir, M., et al. "Quantifying the value of lateral views in deep learning for chest x-rays." Medical Imaging with Deep Learning. PMLR, 2020.
- [5] Raouf S., et al. "Interpretation of plain chest roentgenogram." Chest 141.2 (2012): 545-558.
- [6] Irvin, J., et al. "CheXpert: A large chest radiograph dataset with uncertainty labels and expert comparison." Proceedings of the AAAI conference on artificial intelligence. Vol. 33. 2019.
- [7] DeVries, T. "Improved Regularization of Convolutional Neural Networks with Cutout." arXiv preprint arXiv:1708.04552 (2017).
- [8] Zhang, H., et al. "mixup: Beyond Empirical Risk Minimization." International Conference on Learning Representations. 2018.
- [9] Liu, Y., et al. "Patchdropout: Economizing vision transformers using patch dropout." Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. 2023.
- [10] Sawyer-Lee, R., et al. "Curated Breast Imaging Subset of Digital Database for Screening Mammography (CBIS-DDSM) on The Cancer Imaging Archive. (2016).
- [11] Nguyen, H. TX., et al. "A novel multi-view deep learning approach for BI-RADS and density assessment of mammograms." 2022 44th Annual International Conference of the IEEE Engineering in Medicine & Biology Society (EMBC). IEEE, 2022.
- [12] Chen, S., et al. "Mvt: Multi-view vision transformer for 3d object recognition." arXiv preprint arXiv:2110.13083 (2021).
- [13] Chen, X., et al. "Transformers improve breast cancer diagnosis from unregistered multi-view mammograms." Diagnostics 12.7 (2022): 1549.