

Asymmetric 3D Context Fusion for Universal Lesion Detection

Jiancheng Yang^{1,2,*}, Yi He^{2,*}, Kaiming Kuang², Zudi Lin³,
Hanspeter Pfister³, and Bingbing Ni^{1**}

¹ Shanghai Jiao Tong University, Shanghai, China
nibingbing@sjtu.edu.cn

² Dianei Technology, Shanghai, China

³ Harvard University, Cambridge MA, USA

Abstract. Modeling 3D context is essential for high-performance 3D medical image analysis. Although 2D networks benefit from large-scale 2D supervised pretraining, it is weak in capturing 3D context. 3D networks are strong in 3D context yet lack supervised pretraining. As an emerging technique, *3D context fusion operator*, which enables conversion from 2D pretrained networks, leverages the advantages of both and has achieved great success. Existing 3D context fusion operators are designed to be spatially symmetric, *i.e.*, performing identical operations on each 2D slice like convolutions. However, these operators are not truly equivariant to translation, especially when only a few 3D slices are used as inputs. In this paper, we propose a novel asymmetric 3D context fusion operator (A3D), which uses different weights to fuse 3D context from different 2D slices. Notably, A3D is NOT translation-equivariant while it significantly outperforms existing symmetric context fusion operators without introducing large computational overhead. We validate the effectiveness of the proposed method by extensive experiments on DeepLesion benchmark, a large-scale public dataset for universal lesion detection from computed tomography (CT). The proposed A3D consistently outperforms symmetric context fusion operators by considerable margins, and establishes a new *state of the art* on DeepLesion. To facilitate open research, our code and model in PyTorch is available at <https://github.com/M3DV/AlignShift>.

Keywords: 3D context · universal lesion detection · DeepLesion · A3D.

1 Introduction

Computer vision for medical image analysis has been dominated by deep learning [11, 15], thanks to the availability of large-scale open datasets [1, 24, 18, 6] and powerful infrastructure. In this study, we focus on 3D medical image analysis, *e.g.*, computed tomography (CT) and magnetic resonance imaging (MRI).

* These authors have contributed equally: Jiancheng Yang and Yi He.

** Corresponding author: Bingbing Ni (nibingbing@sjtu.edu.cn).

Spatial information from 3D voxel grids can be effectively learned by convolutional neural networks (CNNs), while 3D context modeling is still essential for high-performance models. There have been considerable debates over 2D and 3D representation learning on 3D medical images; 2D networks benefit from large-scale 2D pretraining [3], whereas the 2D representation is fundamentally weak in large 3D context. 3D networks learn 3D representations; However, few publicly available 3D medical datasets are large enough for 3D pretraining.

Recently, there have been a family of techniques that enable building 3D networks with 2D pretraining [2,13,23,9,22], we refer to it as *3D context fusion operators*. See Sec. 2.1 for a review of existing techniques. These operators learn 3D representations while their (partial) learnable weights can be initialized from 2D convolutional kernels. Existing 3D context fusion operators are convolution-like, *i.e.*, either axial convolutions to fuse slice-wise information [2,13,23] or shifting adjacent slices [9,22]. Therefore, these operators are designed to be spatially **symmetric**: each 2D slice is operated identically. However, convolution-like operations are not truly translation-equivariant [12], due to padding and limited effective receptive fields. In many 3D medical image applications, only a few 2D slices are used as inputs to models due to the memory constraints. It may be meaningless to pursue translation-equivariance in these cases.

In this study, we propose a novel **asymmetric** 3D context fusion operator (A3D). See Sec. 2.2 for the methodology details. Basically, given D slices of 3D input features, A3D uses different weights to fuse the input D slices for each output slice. Therefore, the A3D is **NOT** translation-equivariant. However, it significantly outperforms existing symmetric context fusion operators without introducing large computational overhead in terms of both parameters and FLOPs. We validate the effectiveness of the proposed method by extensive experiments on DeepLesion benchmark [21], a large-scale public dataset for universal lesion detection from computed tomography (CT). As described in Sec. 3, the proposed A3D consistently outperforms symmetric context fusion operators by considerable margins, and establishes a new *state of the art* on DeepLesion.

2 Methods

2.1 Preliminary: 3D Context Fusion Operators with 2D Pretraining

In this section, we briefly review the 3D context fusion operators that enable 2D pretraining, including (a) no fusion, (b) I3D [2], (c) P3D [13], (d) ACS [23] and (e) Shift [9,22]. As an emerging technique, 3D context fusion operator leverages advantages of both 2D pretraining and 3D context modeling.

Given a 3D input feature $\mathbf{X}_i \in \mathbb{R}^{C_i \times D \times H \times W}$, we would like to obtain a transformed 3D output $\mathbf{X}_o \in \mathbb{R}^{C_o \times D \times H \times W}$ with a (pretrained) 2D convolutional kernel $\mathbf{W}_{2D} \in \mathbb{R}^{C_i \times C_o \times K \times K}$, where $D \times H \times W$ is the spatial size of 3D features, C_i and C_o are the input and output channels, and K denotes the kernel size. For simplicity, only cases with same padding are considered here. Apart from convolutions, we simply convert 2D pooling and normalization into 3D [22]. We then introduce each operator as follows:

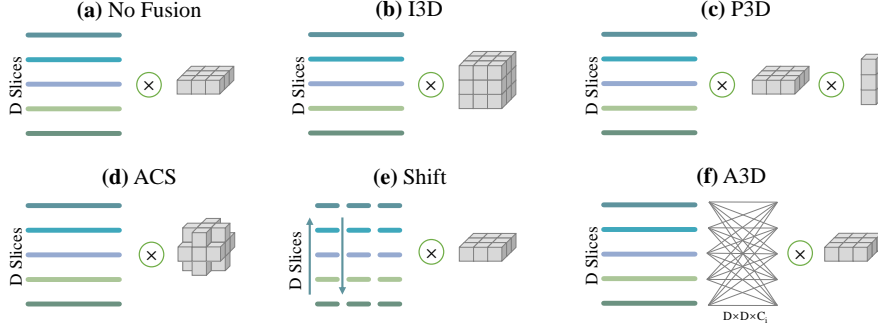


Fig. 1: Illustration of various 3D context fusion operators: (a) no fusion, (b) I3D [2], (c) P3D [13], (d) ACS [23], (e) Shift [9,22] and (f) the proposed A3D. In each sub-figure, left: D slices of C_i -channel 3D features as inputs; middle: \otimes means convolution; right: illustration of convolutional kernels.

(a) *no fusion*. We run 2D convolutions on each 2D slice, which is equivalent to 3D convolutions with $\mathbf{W}_{3D} \in \mathbb{R}^{C_i \times C_o \times 1 \times K \times K}$ converted from the 2D kernel.

(b) *I3D* [2]. I3D is basically an initialization technique for 3D convolution, $\mathbf{W}_{I3D} \in \mathbb{R}^{C_i \times C_o \times K \times K \times K}$ is initialized with K repeats of \mathbf{W}_{2D}/K , so that the distribution expectation of 3D features is the same as that of 2D features.

(c) *P3D* [13]. P3D convolution is a $1 \times K \times K$ 3D convolution followed by a $K \times 1 \times 1$ 3D convolution, where the first convolutional kernel is converted from a 2D kernel (same as *no fusion*), and the second is initialized as $[0, \dots, 1, \dots, 0]$ (e.g., $[0, 1, 0]$ if $K = 3$) to make it as no fusion before training.

(d) *ACS* [23]. ACS runs 2D-like (3D) convolutions in three views of 3D volumes, by splitting the 2D kernel into three 3D kernels: $\mathbf{W}_a \in \mathbb{R}^{C_i \times C_o^{(a)} \times 1 \times K \times K}$, $\mathbf{W}_c \in \mathbb{R}^{C_i \times C_o^{(c)} \times 1 \times K \times K}$ and $\mathbf{W}_s \in \mathbb{R}^{C_i \times C_o^{(s)} \times 1 \times K \times K}$ ($C_o^{(a)} + C_o^{(c)} + C_o^{(s)} = C_o$). 3D context is fused with layer-by-layer ACS transformation without introducing computational cost compared to no fusion.

(e) *Shift* [9,22]. Shift is a family of techniques that fuse 3D context by shifting adjacent 2D slices. Take TSM [9] as an example. It first splits the input feature $\mathbf{X}_i \in \mathbb{R}^{C_i \times D \times H \times W}$ by channel into 3 parts: $\mathbf{X}_i^+ \in \mathbb{R}^{C_i^+ \times D \times H \times W}$, $\mathbf{X}_i^- \in \mathbb{R}^{C_i^- \times D \times H \times W}$ and $\mathbf{X}_i^\pm \in \mathbb{R}^{C_i^\pm \times D \times H \times W}$ ($C_i^+ + C_i^- + C_i^\pm = C_i$). \mathbf{X}_i^+ , \mathbf{X}_i^- and \mathbf{X}_i^\pm are then shifted up, shifted down and kept among the axial axis (D dimension), respectively. Finally, a 3D convolution with $\mathbf{W}_{3D} \in \mathbb{R}^{C_i \times C_o \times 1 \times K \times K}$ (as in no fusion) can fuse 3D context with a single slice. AlignShift [22] is a shift operator adaptive to medical imaging thickness, thus improves the performance of TSM on mixed-thickness data (e.g., a mix of thin- and thick-slice CT scans).

Table 1: Parameters and theoretical (theo.) FLOPs analysis for 3D context fusion operators, in terms of overhead over no fusion, whose parameters and FLOPs are $C_o C_i K^2$ and $\mathcal{O}(DHW C_o C_i K^2)$, respectively. D denotes the number of slices, $D \times W$ denotes the spatial size of each slice, C_i and C_o denote the input and output channel, and K denotes the kernel size. We also provide the numeric FLOPs of the 3D backbone part for 3/7-slice inputs, *i.e.*, GFLOPs (3/7). Additional FLOPs introduced by A3D are marginal given a two-decimal precision.

Operators	No Fusion	I3D [2]	P3D [13]	ACS [23]	Shift [9,22]	A3D (Ours)
Parameters	1	K	$1 + C_o/(C_i K)$	1	1	$1 + D^2/(C_o K^2)$
Theo. FLOPs	1	K	$1 + C_o/(C_i K)$	1	1	$1 + D/(C_o K^2)$
GFLOPs (3)	40.64	78.69	67.79	40.64	40.64	40.64
GFLOPs (7)	94.83	183.61	158.18	94.83	94.83	94.83

Algorithm 1: Asymmetric 3D Context Fusion (A3D)

- Input:** 3D input feature $\mathbf{X}_i \in \mathbb{R}^{C_i \times D \times H \times W}$.
Parameter: asymmetric fusion weight $\mathbf{P} \in \mathbb{R}^{D \times D \times C_i}$,
 2D (pretrained) convolutional kernel $\mathbf{W}_{2D} \in \mathbb{R}^{C_i \times C_o \times K \times K}$.
Output: 3D output feature $\mathbf{X}_o \in \mathbb{R}^{C_o \times D \times H \times W}$.
 1 $\mathbf{W}_{3D} = \text{unsqueeze}(\mathbf{W}_{2D}, \text{dim} = 2) \in \mathbb{R}^{C_i \times C_o \times 1 \times K \times K}$,
 2 $\mathbf{X} = \text{einsum}(\text{"cdhw, dkc} \rightarrow \text{ckhw"} , [\mathbf{X}_i, \mathbf{P}]) \in \mathbb{R}^{C_i \times D \times H \times W}$,
 3 $\mathbf{X}_o = \text{Conv3D}(\mathbf{X}, \text{kernel} = \mathbf{W}_{3D})$.
-

Fig. 1 illustrates these operators, and Table 1 summarizes the computational overhead over no fusion, in terms of parameters and FLOPs. Apart from theoretical FLOPs, we also provide the numeric FLOPs for 3/7-slice inputs to better understand the algorithm complexity in practice. To fairly compare these methods, only FLOPs in 3D backbone are counted, those in 3D-to-2D feature layer and detection heads on 2D feature maps are ignored. Interestingly, additional FLOPs introduced by A3D are marginal given a two-decimal precision.

2.2 Asymmetric 3D Context Fusion (A3D)

The 3D context fusion operators above are designed to be spatially symmetric, *i.e.*, each 2D slice is transformed identically to ensure these convolution-like operations to be translation-equivariant. However, in many medical imaging applications, only a few slices are used as model inputs because of memory constraints ($D = 3$ or 7 in this study). In this case, padding (zero or others) on the axial axis induces a significant distribution shift near top and bottom slices. Moreover, convolution-like operations are not truly translation-equivariant [12] due to limited effective receptive fields. It is not necessary to use spatially symmetric operators in pursuit of translation-equivariance for 3D context fusion.

To address this issue, we propose a novel asymmetric 3D context fusion operator (A3D), which uses different weights to fuse 3D context for each slice.

Mathematically, given a 3D input feature $\mathbf{X}_i \in \mathbb{R}^{C_i \times D \times H \times W}$, A3D fuses features from different slices by creating dense linear connections within the slice dimension for each channel separately. We introduce a trainable asymmetric fusion weight $P \in \mathbb{R}^{D \times D \times C}$, then

$$X^{(c)} = P^{(c)} \cdot X_i^{(c)} \in \mathbb{R}^{D \times H \times W}, c \in \{1, \dots, C\}, \quad (1)$$

where $P^{(c)} \in \mathbb{R}^{D \times D}$ and $X_i^{(c)} \in \mathbb{R}^{D \times H \times W}$ denotes the channel c of P and X_i , respectively, \cdot denotes matrix multiplication. The output X denotes the 3D features after 3D context fusion, it is then transformed by a 3D convolution with $\mathbf{W}_{3D} \in \mathbb{R}^{C_i \times C_o \times 1 \times K \times K}$ (as in no fusion). A3D can be implemented using Einstein summation and 3D convolution in lines of code. Einstein summation saves up extra memories occupied by intermediate results of operations such as transposing, therefore makes A3D faster and more memory-efficient. We depict a PyTorch-fashion pseudo-code of A3D in Algorithm 1. Batch dimension is ignored for simplicity, while the algorithm is easily batched by changing “ $cdhw, dkc \rightarrow ckhw$ ” into “ $bcdhw, dkc \rightarrow bckhw$ ”. A3D is a simple operator that can be plugged into any 3D image model with ease.

To facilitate stable training and faster convergence, the convolution kernels in A3D operation can be initialized with ImageNet [3] pretrained weights to take advantage of supervised pretraining. Furthermore, we initialize each channel of asymmetric fusion weight $P^{(c)}$ with a identity matrix $I \in \mathbb{R}^{D \times D}$ added with a random perturbation following uniform distribution in $[-0.1, 0.1]$, *i.e.*, the A3D is initialized to be like no fusion before training.

Compared to symmetric 3D context fusion operators, A3D uses dense linear connections to gather global contextual information along the axial axis (illustrated in Fig. 1 (f)), thus avoids the padding issue around the top and bottom slices. Besides, as depicted in Table 1, A3D introduces negligible computational overhead in terms of both parameters and FLOPs compared with no fusion. Since D is typically much smaller than C_o , A3D is more lightweight than I3D [2] and P3D [13]. Moreover, as A3D can be implemented with natively supported *einsum*, it is faster than ACS [23] and Shift [9,22] with channel splitting in actual running time. Note that the A3D is NOT translation-equivariant, as it uses different weights for each output slice to fuse the 3D context from input D slices. However, it significantly outperforms existing symmetric context fusion operators with negligible computational overhead.

2.3 Network Structure for Universal Lesion Detection

We develop a universal lesion detection model following Mask R-CNN [4]. An overview of our network is depicted in Fig. 2. The network consists of a DenseNet-121 [5] based 3D backbone with 3D context fusion operators (the proposed A3D or others) plugged in and 2D detection heads. The network backbone takes a gray-scale 3D tensor in shape of $1 \times D \times H \times W$ as input, where D is the number of slices included in each sample ($D \in \{3, 7\}$ in this study). Three dense blocks gradually downsample feature maps and increase number of channels

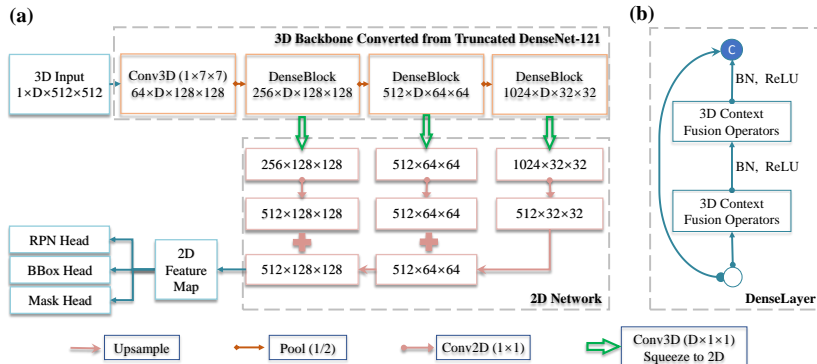


Fig. 2: Universal lesion detection model on DeepLesion [21]. The 3D backbone derived from DenseNet-121 [5, 20] takes a grey-scale 3D input of $D \times 512 \times 512$, where D is the number of slices ($D \in \{3, 7\}$ in this study). Features from different scales are collected and fused together in a feature pyramid [10]. Detection is based on instance segmentation framework using Mask R-CNN [4, 22].

while the depth dimension stays at D . After spatial and channel-wise unification by upsampling and $D \times 1 \times 1$ convolution, 3D features output by three dense blocks are added together and squeezed to 2D by a $D \times 1 \times 1$ convolution. Finally, the 2D feature maps are used for lesion detection on key slices.

3 Experiments

3.1 Dataset and Experiment Settings

DeepLesion dataset [21] includes 32,120 axial CT slices extracted from 10,594 studies of 4,427 patients. There are 32,735 lesions labelled in various organs in total. Each slice contains 1 to 3 lesions, sizes of which range from 0.21 to 342.5mm. RECIST diameter coordinates and bounding boxes are annotated in key slices. Adjacent slices within the range of ± 15 mm from the key slice are given as contextual information.

Our experiments are based on the official code of AlignShift [22], and A3D code is merged into the same code repository. Since DeepLesion does not contain pixel-wise segmentation labels, we use GrabCut [14] to generate weak segmentation labels from RECIST annotations following [22, 26]. Input CT Hounsfield units are clipped to $[-1024, 2050]$ and then normalized to $[-50, 205]$. For AlignShift [22], we process the inputs as in its official code since it uses imaging thickness as inputs. For A3D and other counterparts, we normalize the axial thickness of all data to 2mm and resize each slice to 512×512 . In terms of data augmentation, we apply random horizontal flip, shift, rescaling and rotation during the training stage. No test-time augmentation is adopted. We follow the official data split of 70%/15%/15% for training, validation and test, respectively.

Table 2: Performance evaluated on the large-scale DeepLesion benchmark [21] of the proposed A3D versus other 3D context fusion operators, in terms of sensitivities (%) at various false positives (FPs) per image.

Methods	Slices	0.5	1	2	4	8	16	Avg.[0.5,1,2,4]
No Fusion	×3	72.57	79.89	86.80	91.04	94.24	96.32	82.58
I3D [2]	×3	72.01	80.09	86.54	91.29	93.91	95.68	82.48
P3D [13]	×3	62.13	73.21	82.14	88.6	92.37	94.95	76.52
ACS [23]	×3	72.82	81.15	87.40	91.35	94.69	96.42	83.18
TSM [9]	×3	71.80	80.11	86.97	91.10	93.75	95.56	82.50
AlignShift [22]	×3	73.00	81.17	87.05	91.78	94.63	95.48	83.25
A3D (Ours)	×3	74.10	81.81	87.87	92.13	94.60	96.50	83.98
No Fusion	×7	73.66	82.15	87.72	91.38	93.86	95.98	83.73
I3D [2]	×7	75.37	83.43	88.68	92.20	94.52	96.07	84.92
P3D [13]	×7	74.84	82.17	87.57	91.72	94.90	96.23	84.07
ACS [23]	×7	78.38	85.39	90.07	93.19	95.18	96.75	86.76
TSM [9]	×7	75.98	83.65	88.44	92.14	94.89	96.50	85.05
AlignShift [22]	×7	78.68	85.69	90.37	93.49	95.48	97.05	87.06
A3D (Ours)	×7	80.27	86.73	91.33	94.12	96.15	97.33	88.11

As per [20,26,8], the proposed method and its counterparts are evaluated on the test set using sensitivities at various false positive levels (*i.e.*, FROC analysis). We also implement the mentioned 3D context fusion operators to validate the effectiveness of the proposed A3D.

3.2 Performance Analysis

We compare A3D with a variety of 3D context fusion operators (see Sec. 2.1) on the DeepLesion dataset. Table 2 gives the detailed performances of A3D and all its counterparts on 3 and 7 slices. A3D delivers superior performances compared with all counterparts on both 3 slices and 7 slices. We attribute this performance boost to A3D’s ability of gathering information among globally along the axial axis by creating dense connections among slices, which can be empirically validated by the observation that A3D has a higher performance boost on 7 slices than on 3 slices when compared with the previous *state-of-the-art* AlignShift [22] (+1.05 vs. +0.73) since 7 slices provide more contextual information. Moreover, A3D introduces no padding along the axial axis, this advantage also leads to the performance boost compared to other operators. Note that AlignShift-based model is adaptive to imaging thickness, which is an orthogonal contribution to this study. The asymmetric operation-based methods could be potentially improved by adapting imaging thickness.

Table 3 shows a performance comparison of A3D and previous *State of the Art*. Without heavy engineering and data augmentations, our proposed method outperforms the previous *state-of-the-art* AlignShift [22] on both 3 slices and 7 slices by considerable margin. It is worth noting that A3D with image only

Table 3: Performance evaluated on the large-scale DeepLesion benchmark [21] of the proposed A3D versus previous *state-of-the-art*, in terms of sensitivities (%) at various false positives (FPs) per image.

Methods	Venue	Slices	0.5	1	2	4	8	16	Avg.[0.5,1,2,4]
3DCE [19]	MICCAI'18	×27	62.48	73.37	80.70	85.65	89.09	91.06	75.55
ULDor [16]	ISBI'19	×1	52.86	64.80	74.84	84.38	87.17	91.80	69.22
V.Attn [17]	MICCAI'19	×3	69.10	77.90	83.80	-	-	-	-
Retina. [26]	MICCAI'19	×3	72.15	80.07	86.40	90.77	94.09	96.32	82.35
MVP [8]	MICCAI'19	×3	70.01	78.77	84.71	89.03	-	-	80.63
MVP [8]	MICCAI'19	×9	73.83	81.82	87.60	91.30	-	-	83.64
MULAN [20]	MICCAI'19	×9	76.12	83.69	88.76	92.30	94.71	95.64	85.22
Bou.Maps [7]	MICCAI'20	×3	73.32	81.24	86.75	90.71	-	-	83.01
MP3D [25]	MICCAI'20	×9	79.60	85.29	89.61	92.45	-	-	86.74
AlignShift [22]	MICCAI'20	×3	73.00	81.17	87.05	91.78	94.63	95.48	83.25
AlignShift [22]	MICCAI'20	×7	78.68	85.69	90.37	93.49	95.48	97.05	87.06
ACS [23]	JBHI'21	×3	72.82	81.15	87.40	91.35	94.69	96.42	83.18
ACS [23]	JBHI'21	×7	78.38	85.39	90.07	93.19	95.18	96.75	86.76
A3D	Ours	×3	74.10	81.81	87.87	92.13	94.60	96.50	83.98
A3D	Ours	×7	80.27	86.73	91.33	94.12	96.15	97.33	88.11

surpasses MULAN [20] by nearly 3% even though it takes less slices and no additional information apart from CT images such as medical report tags and demographic information as inputs.

4 Conclusion

In this study, we focus on 3D context fusion operators that enable 2D pre-training, which is an emerging technique that leverages advantages of both 2D pre-training and 3D context modeling. We analyze the unnecessary pursuit of translation-equivariance in existing spatially symmetric 3D context fusion operators especially when only a few 2D slices are used as model inputs. To this end, we further propose a novel asymmetric 3D context fusion operator (A3D) that is translation-equivariant. The A3D significantly outperforms existing symmetric context fusion operators without introducing large computational overhead. Extensive experiments on DeepLesion benchmark validate the effectiveness of the proposed method, and we establish a new *state of the art* that surpasses prior arts by considerable margins.

Acknowledgment. This work was supported by National Science Foundation of China (U20B2072, 61976137).

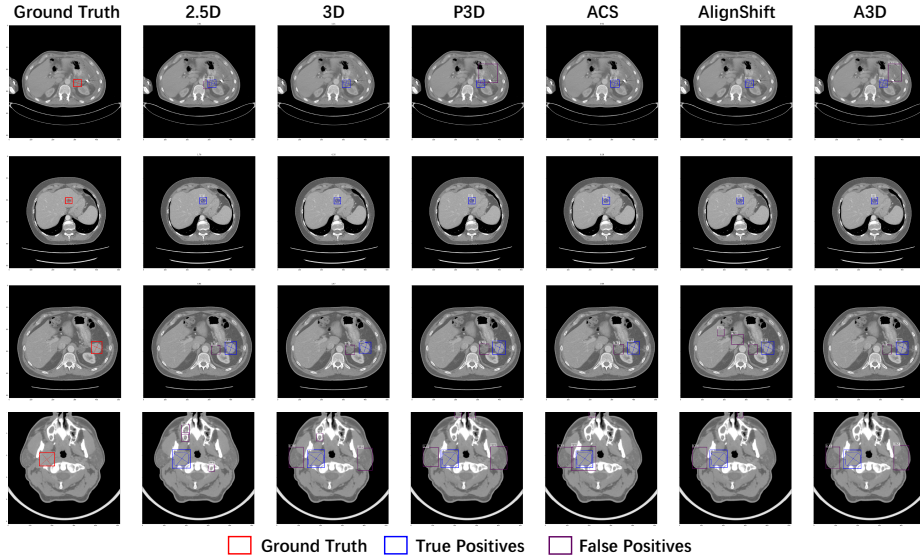


Fig. 3: Visualization of DeepLesion slices highlighted with ground truth and predictions generated by different 3D context fusion operators.

References

1. Antonelli, M., Reinke, A., Bakas, S., et al.: The medical segmentation decathlon. arXiv preprint arXiv:2106.05735 (2021)
2. Carreira, J., Zisserman, A.: Quo vadis, action recognition? a new model and the kinetics dataset. In: CVPR. pp. 6299–6308 (2017)
3. Deng, J., Dong, W., Socher, R., Li, L.J., Li, K., Fei-Fei, L.: Imagenet: A large-scale hierarchical image database. In: CVPR. pp. 248–255 (2009)
4. He, K., Gkioxari, G., Dollár, P., Girshick, R.B.: Mask r-cnn. ICCV pp. 2980–2988 (2017)
5. Huang, G., Liu, Z., Van Der Maaten, L., Weinberger, K.Q.: Densely connected convolutional networks. In: CVPR. vol. 1, p. 3 (2017)
6. Jin, L., Yang, J., Kuang, K., Ni, B., Gao, Y., Sun, Y., Gao, P., Ma, W., Tan, M., Kang, H., et al.: Deep-learning-assisted detection and segmentation of rib fractures from ct scans: Development and validation of fracnet. *EBioMedicine* **62**, 103106 (2020)
7. Li, H., Han, H., Zhou, S.K.: Bounding maps for universal lesion detection. In: MICCAI. pp. 417–428. Springer (2020)
8. Li, Z., Zhang, S., Zhang, J., Huang, K., Wang, Y., Yu, Y.: Mvp-net: Multi-view fpn with position-aware attention for deep universal lesion detection. In: MICCAI. pp. 13–21. Springer (2019)
9. Lin, J., Gan, C., Han, S.: Tsm: Temporal shift module for efficient video understanding. In: ICCV. pp. 7083–7093 (2019)
10. Lin, T.Y., Dollár, P., Girshick, R.B., He, K., Hariharan, B., Belongie, S.J.: Feature pyramid networks for object detection. CVPR pp. 936–944 (2016)

11. Litjens, G., Kooi, T., Bejnordi, B.E., Setio, A.A.A., Ciampi, F., Ghafoorian, M., Van Der Laak, J.A., Van Ginneken, B., Sánchez, C.I.: A survey on deep learning in medical image analysis. *Medical image analysis* **42**, 60–88 (2017)
12. Luo, W., Li, Y., Urtasun, R., Zemel, R.S.: Understanding the effective receptive field in deep convolutional neural networks. In: *NIPS* (2016)
13. Qiu, Z., Yao, T., Mei, T.: Learning spatio-temporal representation with pseudo-3d residual networks. In: *ICCV*. pp. 5533–5541 (2017)
14. Rother, C., Kolmogorov, V., Blake, A.: "grabcut" interactive foreground extraction using iterated graph cuts. *ACM transactions on graphics (TOG)* **23**(3), 309–314 (2004)
15. Shen, D., Wu, G., Suk, H.I.: Deep learning in medical image analysis. *Annual review of biomedical engineering* **19**, 221–248 (2017)
16. Tang, Y.B., Yan, K., Tang, Y.X., Liu, J., Xiao, J., Summers, R.M.: Uldor: a universal lesion detector for ct scans with pseudo masks and hard negative example mining. In: *ISBI*. pp. 833–836. *IEEE* (2019)
17. Wang, X., Han, S., Chen, Y., Gao, D., Vasconcelos, N.: Volumetric attention for 3d medical image segmentation and detection. In: *MICCAI*. pp. 175–184. Springer (2019)
18. Wei, D., Lin, Z., Franco-Barranco, D., Wendt, N., Liu, X., Yin, W., Huang, X., Gupta, A., Jang, W.D., Wang, X., et al.: Mitoem dataset: Large-scale 3d mitochondria instance segmentation from em images. In: *MICCAI*. pp. 66–76. Springer (2020)
19. Yan, K., Bagheri, M., Summers, R.M.: 3d context enhanced region-based convolutional neural network for end-to-end lesion detection. In: *MICCAI*. pp. 511–519. Springer (2018)
20. Yan, K., Tang, Y., Peng, Y., Sandfort, V., Bagheri, M., Lu, Z., Summers, R.M.: Mulan: Multitask universal lesion analysis network for joint lesion detection, tagging, and segmentation. In: *MICCAI* (2019)
21. Yan, K., Wang, X., Lu, L., Zhang, L., Harrison, A.P., Bagheri, M., Summers, R.M.: Deep lesion graphs in the wild: relationship learning and organization of significant radiology image findings in a diverse large-scale lesion database. In: *CVPR*. pp. 9261–9270 (2018)
22. Yang, J., He, Y., Huang, X., Xu, J., Ye, X., Tao, G., Ni, B.: Alignshift: bridging the gap of imaging thickness in 3d anisotropic volumes. In: *MICCAI*. pp. 562–572. Springer (2020)
23. Yang, J., Huang, X., He, Y., Xu, J., Yang, C., Xu, G., Ni, B.: Reinventing 2d convolutions for 3d images. *IEEE Journal of Biomedical and Health Informatics* (2021)
24. Yang, J., Shi, R., Ni, B.: Medmnist classification decathlon: A lightweight automl benchmark for medical image analysis. In: *ISBI* (2021)
25. Zhang, S., Xu, J., Chen, Y.C., Ma, J., Li, Z., Wang, Y., Yu, Y.: Revisiting 3d context modeling with supervised pre-training for universal lesion detection in ct slices. In: *MICCAI*. pp. 542–551. Springer (2020)
26. Zlocha, M., Dou, Q., Glocker, B.: Improving retinanet for ct lesion detection with dense masks from weak recist labels. In: *MICCAI*. pp. 402–410. Springer (2019)