

Reinventing 2D Convolutions for 3D Images

Jiancheng Yang¹, Xiaoyang Huang¹, Yi He, Jingwei Xu, Canqian Yang, Guozheng Xu, and Bingbing Ni¹

Abstract—There have been considerable debates over 2D and 3D representation learning on 3D medical images. 2D approaches could benefit from large-scale 2D pretraining, whereas they are generally weak in capturing large 3D contexts. 3D approaches are natively strong in 3D contexts, however few publicly available 3D medical dataset is large and diverse enough for universal 3D pretraining. Even for hybrid (2D + 3D) approaches, the intrinsic disadvantages within the 2D/3D parts still exist. In this study, we bridge the gap between 2D and 3D convolutions by reinventing the 2D convolutions. We propose ACS (axial-coronal-sagittal) convolutions to perform natively 3D representation learning, while utilizing the pretrained weights on 2D datasets. In ACS convolutions, 2D convolution kernels are split by channel into three parts, and convoluted separately on the three views (axial, coronal and sagittal) of 3D representations. Theoretically, ANY 2D CNN (ResNet, DenseNet, or DeepLab) is able to be converted into a 3D ACS CNN, with pretrained weight of a same parameter size. Extensive experiments validate the consistent superiority of the pretrained ACS CNNs, over the 2D/3D CNN counterparts with/without pretraining. Even without pretraining, the ACS convolution can be used as a plug-and-play replacement of standard 3D convolution, with smaller model size and less computation.

Index Terms—2D-to-3D transfer learning, 3D medical images, ACS convolutions, deep learning.

Manuscript received October 19, 2020; revised December 30, 2020; accepted January 3, 2021. Date of publication January 6, 2021; date of current version August 5, 2021. This work was supported by the National Science Foundation of China (U20B2072, 61976137, U1611461). Authors would like to appreciate the Student Innovation Center of SJTU for providing GPUs. (J. Yang and X. Huang contributed equally to this article.) (Corresponding author: Bingbing Ni.)

Jiancheng Yang is with the Shanghai Jiao Tong University, Shanghai 200240, China, and with MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai 200040, China, and also with Dianei Technology, Shanghai 200050, China (e-mail: jekyll4168@sjtu.edu.cn).

Xiaoyang Huang, Jingwei Xu, Canqian Yang, and Guozheng Xu are with the Shanghai Jiao Tong University, Shanghai 200240, China, and also with MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai 200240, China (e-mail: huangxiaoyang@sjtu.edu.cn; xjwxjw@sjtu.edu.cn; charles.young@sjtu.edu.cn; guozhengxu@sjtu.edu.cn).

Yi He is with the Dianei Technology, Shanghai 200250, China (e-mail: yi.he@dianei-ai.com).

Bingbing Ni is with the Shanghai Jiao Tong University, Shanghai 200240, China, and with MoE Key Lab of Artificial Intelligence, AI Institute, Shanghai Jiao Tong University, Shanghai 200240, China, and also with Huawei Hisilicon, Shanghai 200120, China (e-mail: nibingbing@sjtu.edu.cn).

This article has supplementary downloadable material available at <https://doi.org/10.1109/JBHI.2021.3049452>, provided by the authors.

Digital Object Identifier 10.1109/JBHI.2021.3049452

I. INTRODUCTION

EMERGING deep learning technology has been dominating the medical image analysis research [1], in a wide range of data modalities (e.g., ultrasound [2], CT [3], MRI [4], X-Ray [5]) and tasks (e.g., classification [6], segmentation [7], detection [8], registration [9]). Thanks to contributions from dedicated researchers from academia and industry, there have been much larger medical image datasets than ever before. With large-scale datasets, strong infrastructures and powerful algorithms, numerous challenging problems in medical images seem solvable. However, the data-hungry nature of deep learning limits its applicability in various real-world scenarios with limited annotations. Compared to millions (or even billions) of annotations in natural image datasets, the medical image datasets are not large enough. Especially for 3D medical images, datasets with thousands of supervised training annotations [10], [11] are “large” due to imperfect medical annotations [12]: hardly-accessible and high dimensional medical data, expensive expert annotators (radiologists/clinicians), and severe class-imbalance issues [13].

Transfer learning, with pretrained weights from large-scale datasets (e.g., ImageNet [14], MS-COCO [15]), is one of the most important solutions for annotation-efficient deep learning with insufficient data [12]. Unfortunately, widely-used pretrained CNNs are developed on 2D datasets, which are non-trivial to transfer to 3D medical images. Prior art on 3D medical images follows either 2D-based approaches or 3D-based approaches (compared in Fig. 1). 2D-based approaches [16]–[18] benefit from large-scale pretraining on 2D natural images, while the 2D representation learning are fundamentally weak in large 3D contexts. 3D-based approaches [19]–[21] learn natively 3D representations. However, few publicly available 3D medical dataset is large and diverse enough for universal 3D pretraining. Therefore, compact network design and sufficient training data are essential for training the 3D networks from scratch. Hybrid (2D + 3D) approaches [22]–[24] seem to combine the best of both worlds, nevertheless these ensemble-like approaches do not fundamentally overcome the intrinsic issues of 2D-based and 3D-based approaches. Please refer to Section II for in-depth discussion on these related methods.

There has been considerable debates over 2D and 3D representation learning on 3D medical images: prior studies choose either large-scale 2D pretraining or natively 3D representation learning. This paper presents an alternative to bridge the gap between the 2D and 3D approaches. To solve the intrinsic disadvantages from the 2D convolutions and 3D convolutions in modeling 3D images, we argue that an ideal method should adhere to the following principles:

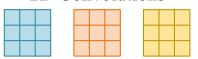

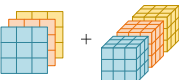
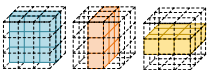
	Pros	Cons
2D Convolutions 	2D pretrained weights on large 2D datasets	Natively 2D representations
3D Convolutions 	Natively 3D representations	Lack of 3D pretrained weights on large datasets
Hybrid (2D + 3D) 	2D + 3D representations	a. 2D representation within 2D parts b. Lack of 3D pretrained weights c. Redundant multi-stage / multi-stream models
ACS Convolutions 	a. Natively 3D representations b. 3D pretrained weights on large 2D datasets c. Converting ANY 2D model into a 3D model seamlessly without extra computation costs	

Fig. 1. A comparison between the proposed ACS convolutions and prior art on modeling the 3D medical images: pure 2D/2.5D approaches with 2D convolution kernels, pure 3D approaches with 3D convolution kernels, and hybrid approaches with both 2D and 3D convolution kernels. The ACS convolutions run multiple 2D convolution kernels among the three views (axial, coronal and sagittal).

- 1) **Natively 3D representation**;
- 2) **2D weight transferable**: it benefits from 2D pretraining;
- 3) **ANY model convertible**: it enables any 2D model, including classification, detection and segmentation backbones, to be converted into a 3D one.

These principles cannot be achieved simultaneously with standard 2D convolutions or standard 3D convolutions, which directs us to develop a novel convolution operator. Inspired from the widely-used tri-planar representations of 3D medical images [16], we propose *ACS convolutions* satisfying these principles. Instead of explicitly treating the input 3D volumes as three orthogonal 2D planar images [16] (axial, coronal and sagittal), we operate on the convolution kernels to perform view-based 3D convolutions, via splitting the 2D convolution kernels into three parts by channel. Notably, **no additional 3D fusion layer** is required to fuse the three-view representations from the 3D convolutions, since they will be seamlessly fused by the subsequent ACS convolution layers (Section III).

The ACS convolution aims at a generic and plug-and-play replacement of standard 3D convolutions for 3D medical images. Even without pretraining, the ACS convolution is comparable to 3D convolution with a smaller model size and less computation. When pretrained on large 2D datasets, it consistently outperforms 2D/3D convolution by a large margin. To improve research reproducibility, a PyTorch [25] implementation of ACS convolution is open-source.¹ Using the provided functions, standard 2D CNNs (e.g., those from PyTorch torchvision package) could be converted into ACS CNNs for 3D images with a single line of code, where 2D pretrained weights could be directly loaded. Compared with 2D models, it introduces no extra computation costs, in terms of FLOPs, memory and model size. The proposed

ACS convolutions could be used in various neural networks for diverse tasks; Extensive experiments on several medical benchmarks (including classification, segmentation and detection tasks) validate the consistent effectiveness of the proposed method.

II. RELATED WORK

In this section, we first review 2D/2.5D, 3D and hybrid approaches for 3D medical images, including their advantages and disadvantages. We then discuss the pretraining for 3D medical images by transfer learning and self-supervised learning techniques. Compared with existing 2D/2.5D/3D/hybrid approaches, ACS convolution focuses on how to use the existing pretrained 2D networks in a 3D way. Note that the contribution of this study is also orthogonal to pretraining methods. It is possible to pretrain ACS CNNs on 2D images, videos and 3D medical images with supervised or self-supervised learning. This paper uses ACS convolutions with supervised pretraining on 2D natural images.

A. 2D/2.5D Approaches

Transfer learning from 2D CNNs, trained on large-scale datasets (e.g., ImageNet [14]), is a widely-used approach in 3D medical image analysis. To mimic the 3-channel image representation (*i.e.*, RGB), prior studies follow either multi-planar or multi-slice representation of 3D images as 2D inputs. In these studies, pretrained 2D CNNs are usually fine-tuned on the target medical dataset.

Early study [16], [26] proposes tri-planar representation of 3D medical images, where three views (axial, coronal and sagittal) from a voxel are regarded as the three channels of 2D input. Although this method is empirically effective, there is a fundamental flaw that the channels are not spatially aligned. More studies follow tri-slice representations [17], [18], [27], where a center slice together with its two neighbor slices are treated as the three channels. In these representations, the channels are spatially aligned, which conforms to the inductive biases in convolution. There are also studies [17], [28] combining both multi-slice and multi-planar approaches, using multi-slice 2D representations in multiple views. The multi-view representations are averaged [17] or fused by additional networks [28]. Recent work [29] extracts multi-view information by applying 2D CNNs on rotated and permuted data. Song *et al.* [30] projects the 3D object boundary surface into a 2D matrix to allow 2D CNN for segmentation.

Even though these approaches benefit from large-scale 2D pretraining, which is empirically effective in numerous studies [31]–[34], both multi-slice and multi-planar representation with 2D convolutions are fundamentally weak in capturing large 3D contexts.

B. 3D Approaches

Instead of regarding the 3D spatial information as input channels in 2D approaches, there are numbers of studies using pure 3D convolutions for 3D medical image analysis [19]–[21], [35],

¹Code is open-source at: <https://github.com/M3DV/ACSCnv/>

[36]. Compared to limited 3D contexts along certain axis in 2D approaches, the 3D approaches are theoretically capable of capturing arbitrarily large 3D contexts in any axis. Therefore, the 3D approaches are generally better at tasks requiring large 3D contexts, e.g., distinguishing small organs, vessels, and lesions.

However, there are also drawbacks for pure 3D approaches. One of the most important is the lack of large-scale universal 3D pretraining. For this reason, efficient training of 3D networks is a pain point for 3D approaches. Several techniques are introduced to (partially) solve this issue, e.g., deep supervision [36], compact network design [21], [37], [38]. Nevertheless, these techniques are not directly targeting the issue of 3D pretraining.

A related study of our method is Parallel Separable Convolution (PSC) [39] and Long-Range Asymmetric Branch (LRAB) [40], which extend pseudo 3D convolution (P3D) [41] in multiple parallel streams of various directions. Both introduce additional layers apart from 2D convolutions, thereby not all weights could be pretrained. As a comparison, our approach focusing on the use of pretrained 2D weights, converts whole pretrained networks seamlessly, while keeping same computation as the 2D variants, in terms of FLOPs, memory and parameters (Table III). In video analysis, there is work [42] similar to the proposed Soft-ACS variant.

C. Hybrid Approaches

Hybrid approaches are proposed to combine the advantages of both 2D and 3D approaches [22]–[24], [28]. In these studies, 2D pretrained networks with multi-slice inputs, and 3D randomly-initialized networks with volumetric inputs are (jointly or separately) trained for the target tasks.

The hybrid approaches could be mainly categorized into multi-stream and multi-stage approaches. In multi-stream approaches [22], [24], 2D networks and 3D networks are designed to perform a same task (e.g., segmentation) in parallel. In multi-stage (*i.e.*, cascade) approaches [23], [24], [28], several 2D networks (and 3D networks) are developed to extract representations from multiple views, and a 3D fusion network is then used to fuse the multi-view representations into 3D representations to perform the target tasks.

Although empirically effective, the hybrid approaches do not solve the intrinsic disadvantages of 2D and 3D approaches: the 2D parts are still not able to capture large 3D contexts, and the 3D parts still lacks large-scale pretraining. Besides, these ensemble-like methods are generally redundant to deploy.

D. Transfer Learning and Self-Supervised Learning

Medical annotations require expertise in medicine and radiology, which are thereby expensive to be scalable. For certain rare diseases or novel applications (e.g., predicting response for novel treatment [43]), the data scale is naturally very small. Transfer learning from large-scale datasets to small-scale datasets is a de-facto paradigm in this case.

Human without any radiological experience could recognize basic anatomy and lesions on 2D and 3D images with limited demonstration. Based on this observation, we believe that transfer learning from universal vision datasets (e.g., ImageNet [14],

TABLE I
COMPARISON OF TRANSFER LEARNING FOR 3D MEDICAL IMAGES FROM VARIOUS SOURCES

Source	Data Scale	Data Diversity	Supervised	Medical
2D Image	Very Large	Very Diverse	Y	N
Video [49]	Large	Diverse	Y	N
Med3D [50]	Moderate	Moderate	Y	Y
MG [51]	Large	Moderate	N	Y

MS-COCO [15]) should be beneficial for 3D medical image analysis. Although there is literature reporting that universal pretraining is useless for target tasks [44], [45], this phenomenon is usually observed when target datasets are large enough. Apart from boosting task performance, the universal pretraining is expected to improve model robustness and uncertainty quantification [46].

Unfortunately, 2D-to-3D transfer learning has not been adequately studied. Research efforts [35], [47] have been paid to pretrain natively 3D CNNs on 3D datasets, however few publicly available 3D medical dataset is large and diverse enough for universal pretraining. Prior research explores the transfer learning of 3D CNNs trained on spatio-temporal video datasets [48]. However, there are two kinds of domain shift between video and 3D medical images: 1) natural images vs. medical images, and 2) spatio-temporal data vs. 3D spatial data. The domain shift makes video pretraining [49] less applicable for 3D medical images. To reduce domain shift, there is research (Med3D [50]) building pretrained 3D models on numbers of 3D medical image datasets. Despite the tremendous effort on collecting data from multiple sources, the data scale of involved 1000+ training samples is still too much small compared to 1000000+ training samples in natural image datasets.

In addition to supervised pretraining, Models Genesis [51] explores unsupervised (self-supervised) learning to obtain the pretrained 3D models. Though very impressive, the model performance of up-to-date unsupervised learning is generally not comparable to that of fully supervised learning; even *state-of-the-art* unsupervised/semi-supervised learning techniques [52], [53] could not reproduce the model performance using fully supervised training data.

Table I compares the sources of transfer learning for 3D medical images. Compared to transfer learning from video [49]/Med3D [50]/Models Genesis [51], the key advantage of 2D image pretraining is the *overwhelming* data scale and diversity of datasets. With the ACS convolutions proposed in this study, we are able to develop natively 3D CNNs using 2D pretrained weights. We compare these pretraining approaches in our experiments, and empirically prove the superiority of the proposed ACS convolutions.

III. ACS CONVOLUTIONAL NEURAL NETWORKS

A. ACS Convolutions

Convolution layers capture spatial correlation. Intuitively, the formal difference between 2D and 3D convolutions is the kernel size: the 2D convolutions use 2D kernels ($C_o \times C_i \times K \times K$)

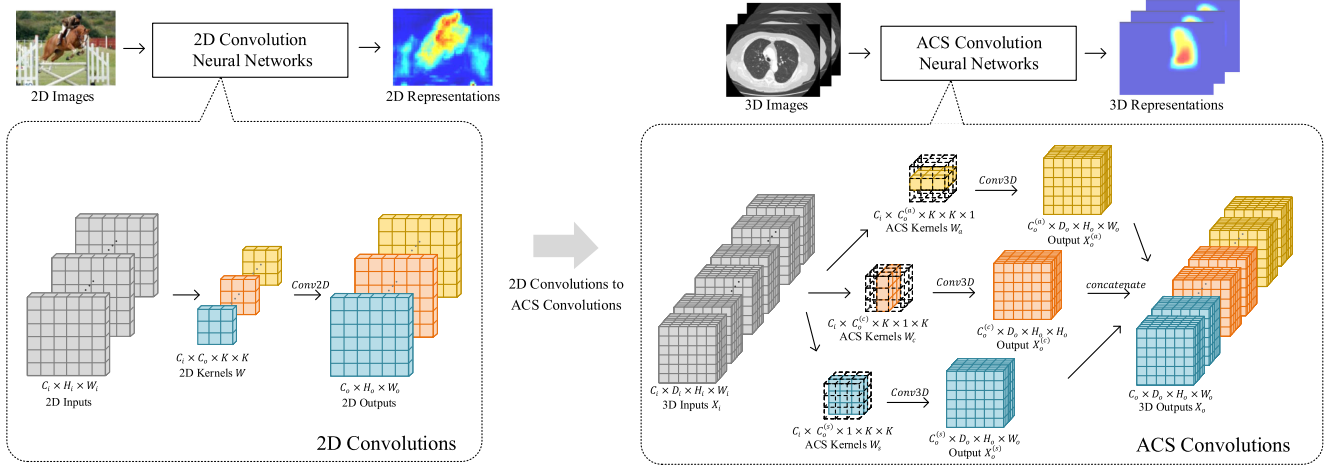


Fig. 2. Illustration of ACS convolutions and 2D-to-ACS model conversion. With a kernel-splitting design, a 2D convolution kernel could be seamlessly transferred into ACS convolution kernels to perform natively 3D representation learning. The ACS convolutions enable **ANY** 2D model (ResNet [54], DenseNet [55], or DeepLab [34]) to be converted into a 3D model.

for 2D inputs ($C_i \times H_i \times W_i$), whereas the 3D convolutions use 3D kernels ($C_o \times C_i \times K \times K \times K$) for 3D inputs ($C_i \times D_i \times H_i \times W_i$), where C_i, C_o denote the channels of inputs and outputs, K denotes the kernel size, and $(D_i \times H_i \times W_i)$ denotes the input size. To transfer the 2D kernels to 3D kernels, there are basically two prior approaches: 1) “inflate” the pretrained 2D kernels into 3D kernels size ($K \times K \rightarrow K \times K \times K$), *i.e.*, Inflated 3D (I3D [56]), where the 2D kernels are repeated along an axis and then normalized; 2) unsqueeze the 2D kernels into pseudo 3D kernels on an axis ($K \times K \rightarrow 1 \times K \times K$), *i.e.*, AH-Net-like [57], which could not effectively capture 3D contexts. Note that in both cases, the existing methods assume a specific axis to transfer the 2D kernels. It is meaningful to assign a special axis for spatio-temporal videos, while controversial for 3D medical images. Even for anisotropic medical images, *any view of a 3D image is still a 2D spatial image*.

Based on this observation, we develop ACS (axial-coronal-sagittal) convolutions to learn spatial representations from the axial, coronal and sagittal views. Instead of treating channels of 2D kernels identically [56], [57], we split the kernels into three parts for extracting 3D spatial information from the axial, coronal and sagittal views. The detailed algorithm of ACS convolutions is shown in the supplementary materials. For simplicity, we introduce ACS convolutions with same padding as follows (Fig. 2).

Given a 3D input $\mathbf{X}_i \in \mathbb{R}^{C_i \times D_i \times H_i \times W_i}$, we would like to obtain a 3D output $\mathbf{X}_o \in \mathbb{R}^{C_o \times D_o \times H_o \times W_o}$, with pretrained/non-pretrained 2D kernels $\mathbf{W} \in \mathbb{R}^{C_o \times C_i \times K \times K}$. Here, C_i and C_o denote the input and output channels, $D_i \times H_i \times W_i$ and $D_o \times H_o \times W_o$ denote the input and output sizes, K denotes the kernel size. Instead of presenting 3D images into tri-planar 2D images [16], we split and reshape the kernels into three parts (named ACS kernels) by the output channel, to obtain the view-based 3D representations for each volume: $\mathbf{W}_a \in \mathbb{R}^{C_o^{(a)} \times C_i \times K \times K \times 1}$, $\mathbf{W}_c \in \mathbb{R}^{C_o^{(c)} \times C_i \times K \times 1 \times K}$, $\mathbf{W}_s \in \mathbb{R}^{C_o^{(s)} \times C_i \times 1 \times K \times K}$, where $C_o^{(a)} + C_o^{(c)} + C_o^{(s)} = C_o$. It

is theoretically possible to assign an “optimal axis” for a 2D kernel; However, considering the feature redundancy in CNNs [58], in practice we simply set $C_o^{(a)} \approx C_o^{(c)} \approx C_o^{(s)} \approx \lfloor C_o/3 \rfloor$. We then compute the axial, coronal and sagittal view-based 3D features via 3D convolutions:

$$\mathbf{X}_o^{(a)} = \text{Conv3D}(\mathbf{X}_i, \mathbf{W}_a) \in \mathbb{R}^{C_o^{(a)} \times D_o \times H_o \times W_o}, \quad (1)$$

$$\mathbf{X}_o^{(c)} = \text{Conv3D}(\mathbf{X}_i, \mathbf{W}_c) \in \mathbb{R}^{C_o^{(c)} \times D_o \times H_o \times W_o}, \quad (2)$$

$$\mathbf{X}_o^{(s)} = \text{Conv3D}(\mathbf{X}_i, \mathbf{W}_s) \in \mathbb{R}^{C_o^{(s)} \times D_o \times H_o \times W_o}. \quad (3)$$

The output feature \mathbf{X}_o is obtained by concatenating $\mathbf{X}_o^{(a)}$, $\mathbf{X}_o^{(c)}$ and $\mathbf{X}_o^{(s)}$ by the channel axis. It is noteworthy that, **no 3D fusion layer is required additionally**. The view-based output features will be automatically fused by subsequent convolution layers, without any additional operation, since the convolution kernels are **not** split by input channel. Thanks to linearity of convolutions, expectation of features from converted ACS convolutions keeps the same as that of 2D ones, thereby no weight rescaling [56] is needed. It is also the prerequisite for the usefulness of 2D pretraining in the converted ACS convolutions. The ACS convolution could be regarded as a special case of 3D convolutions, whose kernels are block sparse.

The parameter size of ACS convolutions is exactly same as that of 2D convolutions, as the ACS kernels: $\mathbf{W}_a, \mathbf{W}_c$ and \mathbf{W}_s are directly split and reshaped from the 2D kernels \mathbf{W} , therefore the proposed method enables ANY 2D model to be converted into a 3D model. Table II lists how operators in 2D CNNs are converted to those in ACS CNNs. Note that the converted models could load the 2D weights directly.

B. Counterparts and Related Methods

1) **2D Convolutions**: We include a simple AH-Net-like [57] 2D counterpart, by replacing all ACS convolutions in ACS CNNs with Conv3D $1 \times K \times K$. We name this pseudo 3D

TABLE II

MAIN OPERATOR CONVERSION FROM 2D CNNs INTO ACS CNNs.
 $\{1, K\} \times K \times K$ AND $(1 \times) 1 \times 1$ DENOTE THE KERNEL SIZES

2D CNNs	ACS CNNs
Conv2D $K \times K$	ACSCov $K \times K$
Conv2D 1×1	Conv3D $1 \times 1 \times 1$
{Batch,Group}Norm2D	{Batch,Group}Norm3D
{Max,Avg}Pool2D $K \times K$	{Max,Avg}Pool3D $\{1, K\} \times K \times K$

TABLE III

SPACE AND TIME COMPLEXITY ANALYSIS, FOR 2D (2.5D), 3D, ACS, MEAN-ACS, AND SOFT-ACS CONVOLUTIONS. $D \times H \times W$ DENOTES SPATIAL SIZE OF A 3D INPUT, K DENOTES KERNEL SIZE (IDENTICAL KERNEL FOR SIMPLICITY), AND C_i, C_o DENOTE THE INPUT AND OUTPUT CHANNELS. BIAS TERMS ARE NOT COUNTED IN PARAMETERS

Kernels	FLOPs	Memory	Parameters
Conv2D	$\mathcal{O}(DHW C_o C_i K^2)$	$DHWC_o$	$C_o C_i K^2$
Conv3D	$\mathcal{O}(DHW C_o C_i K^3)$	$DHWC_o$	$C_o C_i K^3$
ACSCov	$\mathcal{O}(DHW C_o C_i K^2)$	$DHWC_o$	$C_o C_i K^2$
M-ACSCov	$\mathcal{O}(3DHW C_o C_i K^2)$	$3DHW C_o$	$C_o C_i K^2$
S-ACSCov	$\mathcal{O}(3DHW C_o C_i K^2)$	$3DHW C_o$	$C_o(C_i K^2 + 3)$

counterpart as “2.5D” in our experiments, which enables 2D pretrained weight transferring with ease. The 3D pooling and normalization layers enable 3D context fusion in this case; although insufficient in 3D, we would rather call this variant as “2.5D”.

2) 3D Convolutions: For the 3D counterparts, we replace all convolutions in ACS CNNs with standard 3D convolutions. Various pretraining sources (I3D [56] with 2D images, Med3D [50], Video [49]) are included for fair comparison. If there is any difference between the converted 3D models and the pretrained 3D models, we keep the pretrained 3D network architectures to load the pretrained weights. Models Genesis [51] uses 3D UNet-based [19], [20] network architecture. We train the same network from scratch/with its self-supervised pretraining to compare with our models. Moreover, we implement $P_3 SC_1$ to compare Parallel Separable Convolutions [39] with ACS convolutions. The $P_3 SC_1$ models are trained from scratch due to the lack of pretraining with this method.

Table III compares the time and space complexity of 2D (2.5D), 3D and ACS convolutions. The proposed ACS convolution could be used as a generic and plug-and-play replacement of 3D convolution, with less computation and smaller size. Besides, the ACS convolution enables 2D pretraining. We demonstrate its superiority over the counterparts with extensive experiments (Section IV).

C. ACS Convolution Variants

Apart from the kernel splitting approach used in the proposed ACS convolutions, there are possible variants to implement the 2D-transferable ACS convolutions.

1) Mean-ACS Convolutions: Instead of splitting the 2D convolution kernels, we replicate and reshape \mathbf{W} into $\mathbf{W}'_a \in \mathbb{R}^{C_o \times C_i \times K \times K \times 1}$, $\mathbf{W}'_c \in \mathbb{R}^{C_o \times C_i \times K \times 1 \times K}$, $\mathbf{W}'_s \in \mathbb{R}^{C_o \times C_i \times 1 \times K \times K}$, and obtain the 3D features $\mathbf{X}'_o(a) = \text{Conv3D}(\mathbf{X}_i, \mathbf{W}'_a)$, $\mathbf{X}'_o(c) = \text{Conv3D}(\mathbf{X}_i, \mathbf{W}'_c)$,

$\mathbf{X}'_o(s) = \text{Conv3D}(\mathbf{X}_i, \mathbf{W}'_s)$. The output features is

$$\mathbf{X}_o^M = (\mathbf{X}'_o(a) + \mathbf{X}'_o(c) + \mathbf{X}'_o(s))/3. \quad (4)$$

Note that the replication and reshaping operations for \mathbf{W} are conducted on the fly, thereby only one copy of the \mathbf{W} are required to be stored for the model.

2) Soft-ACS Convolutions: Note that the Mean-ACS convolution uses a symmetric aggregation, thereby it could not distinguish any view-based information. To this regard, we introduce weighted sum of Mean-ACS, *i.e.*, Soft-ACS,

$$\mathbf{X}_o^S = \alpha^{(a)} \mathbf{X}'_o(a) + \alpha^{(c)} \mathbf{X}'_o(c) + \alpha^{(s)} \mathbf{X}'_o(s), \quad (5)$$

where $\alpha^{(a)}, \alpha^{(c)}, \alpha^{(s)} \in \mathbb{R}$ are learnable weights.

In Table III, we compare the time and space complexity. The two variants are more computationally intensive in terms of FLOPs and memory. Unfortunately, they do not provide significant performance boost empirically. Therefore, we only report the model performance of ACS convolutions in Section IV, and analyze these variants in the supplementary materials.

IV. EXPERIMENTS

We experiment with the proposed method on a proof-of-concept dataset and medical benchmarks. To fairly compare model performance, we include several counterparts (2.5D/3D/ACS {Network} **r/p.**) with the same experiment setting, where **r.** denotes random initialization, and **p.** denotes pretraining on various sources. We use separate network architectures in different experiments to demonstrate the flexibility and versatility of the proposed method.

A. Proof of Concept: How Does the Pretraining Work?

1) Motivation: We would like to intuitively understand whether the 2D pretraining could be useful for the converted ACS models. For this reason, we design a proof-of-concept experiments with a synthetic dataset, where the “knowledge” from the 2D space is guaranteed to be useful in the 3D space. UNet-based models [19], [59] segment the foregrounds, on the dataset consisting of sufficient 2D samples (foreground: circle and square) for pretraining the 2D networks and limited 3D samples (foreground: sphere, cube, cylinder, cone and pyramid) for evaluating the converted ACS networks. Note that the shapes of 2D dataset are exactly the projected single views of 3D volumes (except for triangle), thereby the 2D pretraining is expected to be useful in the 3D segmentation. We quantitatively analyze feature discriminative ability on the 3D dataset **without** training on it, using a mAUC metric based on Receiver Operating Characteristic (ROC) analysis of final layer features to discriminate the foreground. We then train the UNets and compare the ACS, 2.5D and 3D counterparts. Dice and mIoU averaged on the 5 foreground classes are reported on the 3D dataset. Details of the synthetic dataset, network, training and mAUC metric are provided in the supplementary materials.

2) Result Analysis: We first illustrate two examples of feature maps from ACS **r.** and ACS **p.** on the 3D dataset in Fig. 3. Even without any training on the target dataset, the features from ACS **p.** are well aligned with the foreground. As shown in

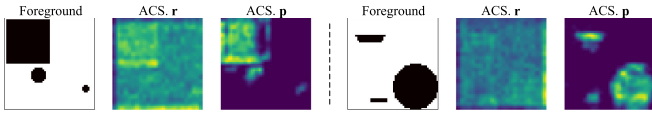


Fig. 3. Examples of features from ACS **r.** and **p.** on the 3D dataset **without** any training on it. ACS **p.** is pretrained on the 2D dataset.

TABLE IV
SEGMENTATION PERFORMANCE ON THE PROOF-OF-CONCEPT DATASET

Models	Feature mAUC w/o training	Dice	mIoU	Size
2.5D UNet r.	69.0	82.2	72.5	1.6 Mb
2.5D UNet p.	85.1	82.7	73.3	1.6 Mb
3D UNet r.	72.1	94.6	90.8	4.7 Mb
ACS UNet r.	68.7	94.7	90.7	1.6 Mb
ACS UNet p.	88.1	95.4	92.0	1.6 Mb

Table IV, the mAUC metric of ACS **p.** is significantly higher than ACS **r.**, which empirically proves the usefulness of pretraining for ACS convolutions. Notably, features from ACS **p.** are even more discriminative than 2.5D **p.** After training on the 3D dataset, the performance of ACS UNet **r.** is comparable to 3D UNet **r.**, and the ACS UNet **p.** achieves the best performance. The results indicate that ACS convolution is an alternative to 3D convolution with comparable or even better performance, and a smaller model size.

B. Lung Nodule Classification and Segmentation

1) *Dataset*: We then validate the effectiveness of the proposed method on a large medical data LIDC-IDRI [60], the largest public lung nodule dataset, for both lung nodule segmentation and malignancy classification task. There are 2635 lung nodules annotated by at most 4 experts, from 1018 CT scans. The annotations include pixel-level labelling of the nodules and 5-level classification of the malignancy, from “1” (highly benign) to “5” (highly malignant). For segmentation, we choose one of the up to 4 annotations for all cases. For classification, we take the mode of the annotations as its category. In order to reduce ambiguity, we ignore nodules with level-“3” (uncertain labelling) and perform binary classification by categorizing the cases with level “1/2,” “4/5” into class 0, 1. It results in a total of 1633 nodules for classification. We randomly divide the dataset into 4 : 1 for training and evaluation, respectively. At training stage we perform data augmentation including random-center cropping, random-axis rotation and flipping.

2) *Experiment Setting*: We compare the ACS models with 2.5D and 3D counterparts with or without pretraining. The pretrained 2.5D/ACS weights are from PyTorch torchvision package [25], trained on ImageNet [14]. For 3D pretraining, we use the official pretrained models by Med3D [50] and Video [49], while I3D [56] weights are transformed from the 2D ImageNet-pretrained weights as well. For PSC [39], we use the configuration of P₃ SC₁, which resembles our ACS methods. To take advantage of the pretrained weights from Med3D [50] and video [49] for comparison, all models are adopted a ResNet-18 [54] architecture, except for Model Genesis [51], since the

TABLE V
LIDC-IDRI [60] LUNG NODULE SEGMENTATION (DICE GLOBAL) AND CLASSIFICATION (AUC) PERFORMANCE

Models	Segmentation (Dice)	Classification (AUC)
Models Genesis [51] r.	75.5	94.3
Models Genesis [51] p.	75.9	94.1
P ₃ SC ₁ [39] r.	74.3	90.9
2.5D ResNet-18 r.	68.8	89.4
2.5D ResNet-18 p.	69.8	92.0
3D ResNet-18 r.	74.7	90.3
3D ResNet-18 p. I3D [56]	75.7	91.5
3D ResNet-18 p. Med3D [50]	74.9	90.6
3D ResNet-18 p. Video [49]	75.7	91.0
ACS ResNet-18 r.	75.1	92.5
ACS ResNet-18 p.	76.5	94.9

official pretrained model is based on a 3D UNet [19] architecture. For all model training, we use an Adam optimizer [61] with an initial learning rate of 0.001 and train the model for 100 epochs, and delay the learning rate by 0.1 after 50 and 75 epochs. For ResNet-18 backbone, in order to keep higher resolution for output feature maps, we modify the stride of first layer (7×7 stride-2 convolution) into 1, and remove the first max-pooling. Note that this modification still enables pretraining. A FCN-like [31] decoder is applied with progressive upsampling twice. Dice loss with a batch of 8 is used for segmentation, and binary cross-entropy loss with a batch of 24 for classification. Dice global and AUC are reported for these two tasks. To demonstrate the flexibility and versatility of ACS convolutions, we also report the results of VGG [62] and DenseNet [55] in the supplementary materials, which is consistent with the ResNet-18.

To further demonstrate the effectiveness of pretraining, we evaluate models trained with various percentages of training data (25%, 50%, 75% and 100%) on both segmentation and classification. To maintain the same number of training iterations, the numbers of epoch are increased to 100/(0.25, 0.5, 0.75, 1). The best results among all epochs are reported. We plot the results in Fig. 4.

3) *Result Analysis*: Experiment results are depicted in Table V. The ACS models consistently outperform all the counterparts by large margins, including 2.5D and 3D models in both random initialization or pretraining settings. P₃ SC₁ [39] performance is similar to standard 3D convolution. We observe that the 3D models (ACS, PSC and 3D) generally outperform the 2.5D models, indicating that the usefulness of 3D contexts. Except for the pretrained 2.5D model on classification task, its superior performance over 3D counterparts explain the prior art [63], [64] with 2D networks on this dataset. As for pretraining, the ImageNet [14] provides significant performance boost (see 2.5D **p.**, 3D **p.** I3D [56] and ACS **p.**), while Med3D [50] brings limited performance boost. We conjecture that it owes to the overwhelming data scale and diversity of 2D image dataset.

Due to the difference on network architecture (ResNet-based FCN vs. UNet), we experiment with the official code of self-supervised pretrained Models Genesis [51] with exactly same setting. Even without pretraining, the segmentation and classification performance of the UNet-based models are strong on this dataset. Despite this, the pretrained ACS model is still

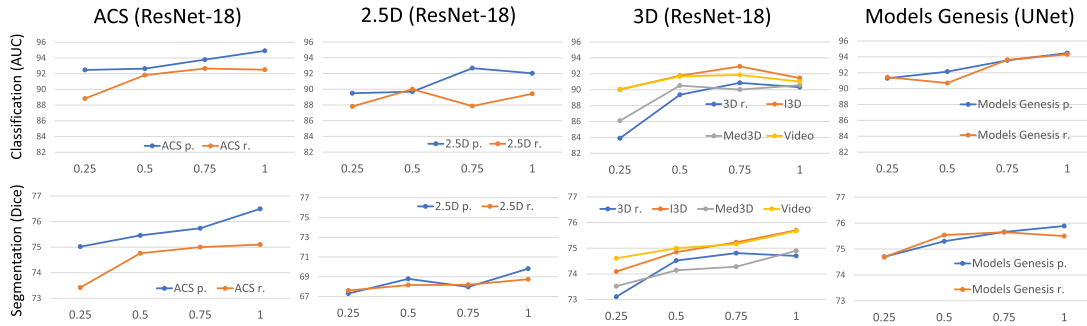


Fig. 4. Performance of ACS, 2.5D, 3D and MG [51] on LIDC-IDRI [60] lung nodule classification and segmentation vs. the training data scale.

better performing. Besides, negative transferring is observed for classification experiments by the Models Genesis [51] encoder-only transferring, whereas the ImageNet pretraining consistently improves the performance. Apart from the superior model performance, the ACS model achieves the best parameter efficiency in our experiments. Take the segmentation task for example, the size of ACS model is 49.8 Mb, compared to 49.8 Mb (2.5D), 142.5 Mb (3D) and 65.4 Mb (MG [51]).

As showed in Fig 4, models with pretraining consistently outperform those without pretraining under various training data scales. Moreover, when trained with 25% data, the performance gap between *p.* and *n.* reaches the highest, which implies the efficiency of pretraining for limited annotated data. Note that ACS *p.* consistently outperforms all counterparts, no matter how much training data are leveraged.

C. Liver Tumor Segmentation (LiTS) Benchmark

1) *Dataset*: We further experiment with our approach on LiTS [65], a challenging 3D medical image segmentation dataset. It consists of 131 and 70 enhanced abdominal CT scans for training and testing respectively, to segment the liver and liver tumors. The training annotations are open to public while the test ones are only accessible by online evaluation. The sizes of *x, y* axis are 512, while the sizes of *z* axis vary in the range of [50, 1000]. We transpose the axes into *z, y, x* to keep the concept consistent as previously mentioned. For pre-processing, we clip the Hounsfield Unit to $[-200, 250]$ and then normalize to $[0, 1]$. Training data augmentation includes random-center cropping, random-axis flipping and rotation, and random-scale resampling.

2) *Experiment Setting*: A DeepLabv3+ [34] with a backbone of ResNet-101 [54] is used in this experiment. The pre-trained 2D model is directly obtained from PyTorch torchvision package [25]. The compared baselines are similar to those in the above LIDC experiment (Section IV-C). We train all the models for 6000 epochs. An Adam optimizer [61] is used with an initial learning rate of 0.001, and we decay the learning rate by 0.1 after 3000 and 4500 epochs. At training stage, we crop the volumes to the size of $64 \times 224 \times 224$. As for testing stage, we crop the volumes to the size of $64 \times 512 \times 512$ and adopt window sliding at a step of 24 at *z* axis. Dice global and Dice per case of lesion and liver are reported as standard evaluation on this dataset.

TABLE VI

LiTS [65] SEGMENTATION PERFORMANCE ON LESION AND LIVER. DG: DICE GLOBAL. DPC: DICE PER CASE

Models	Lesion		Liver	
	DG	DPC	DG	DPC
H-DenseUNet [22]	82.4	72.2	96.5	96.1
Models Genesis [51] ²	-	-	-	91.13 \pm 1.51
P ₃ SC ₁ [39] r.	72.6	59.1	93.9	94.2
2.5D DeepLab r.	72.9	56.9	93.1	92.7
2.5D DeepLab p.	73.4	60.4	92.9	92.0
3D DeepLab r.	75.5	62.6	94.8	94.8
3D DeepLab p. I3D [56]	76.5	58.2	94.1	93.4
3D DeepLab p. Med3D [50]	67.1	53.9	92.0	93.6
3D DeepLab p. Video [49]	66.3	56.9	92.5	93.2
ACS DeepLab r.	75.3	62.4	95.0	94.9
ACS DeepLab p.	79.1	65.8	96.7	96.2

3) *Result Analysis*: As shown in Table VI,² consistent model performance as LIDC experiment (Section IV-C) can be observed. The pretrained ACS DeepLab achieves better performance than the 2D and 3D counterparts (including self-supervised pretraining [51]) by large margins; without pretraining, ACS DeepLab achieves comparable or better performance than 3D DeepLab. According to pretraining results on I3D [56], Med3D [50] and Video [49] for 3D DeepLab, negative transferring is observed, probably due to severe domain shift and anisotropy on LiTS dataset. We also report a *state-of-the-art* performance on LiTS dataset using H-DenseUNet [22] as a reference. Note that it adopts a completely different training strategy and network architecture (a cascade of 2D and 3D DenseNet [55] models), thereby it is insuitable to compare to our models directly. It is feasible to integrate these orthogonal contributions into ours to improve the model performance.

A key advantage of the proposed ACS convolution is that it enables flexible whole-network conversion together with the pretrained weights, including the task head. In the supplementary materials, we validate the superiority of whole-network weight transferring over encoder-only weight transferring. The previous pretraining methods, e.g., I3D [56], Med3D [50] and Video [49], hardly take care of the scenarios.

²The author only releases the pretrained model on chest CTs, thereby we simply report the evaluation metric provided by the paper.

TABLE VII

PERFORMANCE ON THE DEEPLesion BENCHMARK [3], IN TERMS OF DETECTION SENSITIVITY (SENS, %) AT VARIOUS FALSE POSITIVES (FPs) PER IMAGE. NOTE THAT MULAN [8] USES EXTRA CLASSIFICATION TAG SUPERVISION AND AN ADDITION SCORE REFINEMENT LAYER (SRL) WITH TAG INPUTS; WE REPORT ITS PERFORMANC UNDER PUBLIC 171-TAG SUPERVISION AS WELL AS THAT WITHOUT SRL

Methods	Venue	Slices	Sens@0.5	Sens@1	Sens@2	Sens@4	Sens@8	Sens@16	Avg.[0.5,1,2,4]
3DCE [67]	MICCAI'18	$\times 27$	62.48	73.37	80.70	85.65	89.09	91.06	75.55
ULDor [68]	ISBI'19	$\times 1$	52.86	64.8	74.84	84.38	87.17	91.8	69.22
Volumetric Attention [69]	MICCAI'19	$\times 3$	69.10	77.90	83.80	-	-	-	-
Improved RetinaNet [70]	MICCAI'19	$\times 3$	72.15	80.07	86.40	90.77	94.09	96.32	82.35
MVP-Net [71]	MICCAI'19	$\times 3$	70.01	78.77	84.71	89.03	-	-	80.63
MVP-Net [71]	MICCAI'19	$\times 9$	73.83	81.82	87.60	91.30	-	-	83.64
MULAN (Mask R-CNN) [8]	MICCAI'19	$\times 9$	76.12	83.69	88.76	92.30	94.71	95.64	85.22
MULAN (Mask R-CNN) w/o SRL [8]	MICCAI'19	$\times 9$	-	-	-	-	-	-	84.22
2.5D Mask R-CNN r.	Ours	$\times 3$	70.34	79.11	86.16	90.94	94.04	96.01	81.64
2.5D Mask R-CNN p.	Ours	$\times 3$	72.57	79.89	86.80	91.04	94.24	96.32	82.58
3D Mask R-CNN r.	Ours	$\times 3$	63.58	74.63	82.88	88.03	91.39	93.99	77.28
3D Mask R-CNN p. I3D [56]	Ours	$\times 3$	72.01	80.09	86.54	91.29	93.91	95.68	82.48
ACS Mask R-CNN r.	Ours	$\times 3$	72.52	80.85	87.10	91.05	94.39	96.12	82.88
ACS Mask R-CNN p.	Ours	$\times 3$	73.00	81.17	87.05	91.78	94.63	95.48	83.25
2.5D Mask R-CNN r.	Ours	$\times 7$	73.37	81.13	86.73	90.96	93.99	95.79	83.05
2.5D Mask R-CNN p.	Ours	$\times 7$	73.66	82.15	87.72	91.38	93.86	95.98	83.73
3D Mask R-CNN r.	Ours	$\times 7$	64.10	75.14	82.13	87.44	91.53	94.22	77.20
3D Mask R-CNN p. I3D [56]	Ours	$\times 7$	75.37	83.43	88.68	92.20	94.52	96.07	84.92
ACS Mask R-CNN r.	Ours	$\times 7$	78.01	84.75	88.97	91.76	93.79	95.26	85.87
ACS Mask R-CNN p.	Ours	$\times 7$	78.38	85.39	90.07	93.19	95.18	96.75	86.76

D. Universal Lesion Detection on DeepLesion

1) **Dataset:** DeepLesion dataset [3] consists of 32 120 axial CT slices from 10 594 studies of unique patients. There are 1 to 3 lesions in each slice, with totally 32735 lesions from several organs, whose sizes vary from 0.21 to 342.5 mm. RECIST diameter coordinates and bounding boxes were labeled on the key slices, with adjacent slices (above and below 30 mm) provided as contextual information. We use GrabCut algorithm to generate weak segmentation “ground truth” from weak RECIST labels [70]. Hounsfield units of the input are clipped into $[-1024, 2050]$ and normalized. The thickness of all slices is normalized into 2 mm.

Data augmentation including horizontal flip, shift, rescaling and rotation is applied during training stage, no test-time augmentation (TTA) is applied. We resize each input slice to 512×512 before feeding into the networks. We use official data split (training/validation/test: 70%/15%/15%). To evaluate detection performance, sensitivity at various false positives levels (*i.e.*, FROC analysis) is tuned on the validation set and evaluated on the test set.

2) **Experiment Setting:** As illustrated in Fig. 5, a detection and instance segmentation network based on Mask R-CNN [66] with a same backbone of previous *state-of-the-art* MULAN [8] is developed for universal lesion detection. Since the inputs are 3D slices whereas only 2D key-slice annotations are available on DeepLesion dataset [3], we use a 3D backbone (truncated DenseNet-121 [55] converted with 2.5D/3D/ACS convolutions) with a 2D outputs. All $2D \times 2$ pooling operators are converted into $3D \times 1 \times 2$ pooling. The encoder takes a grey-scale 3D tensor of $1 \times D \times 512 \times 512$ as input, where D is the length of key slices ($D = 3, 7$ in this study), and extracts 3D features through three dense blocks. The feature output of each dense blocks is processed by a $D \times 1 \times 1$ 3D convolution and then squeezed into a 2D shape. A 2D decoder then combines these features under different resolutions and upsamples the features

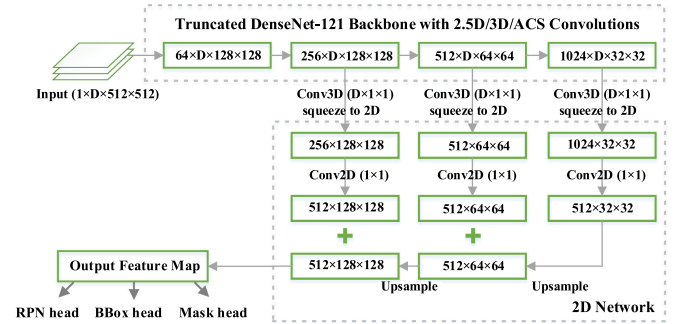


Fig. 5. Illustration of the Mask R-CNN [66] for universal lesion detection on DeepLesion [3]. The 3D backbone is converted from a truncated DenseNet-121 [8], [55] with 2.5D/3D/ACS convolutions.

step by step. The final feature map is fed into RPN head, BBox head and Mask head for detection and instance segmentation supervised by weak RECIST labels. We implement the Mask-RCNN with PyTorch [25] and MMDetection [72].

We adopt cross entropy loss for classification and smooth L1 loss for bounding box regression in the RPN head and BBox head, while in the Mask head we adopt the dice loss. These losses weigh equally to form the final loss function. We use an momentum SGD optimizer to train the models for 20 epochs. The learning rate is initialized as 0.02 and multiplied by 0.1 at epoch 10 and 13.

3) **Result Analysis:** As depicted in Table VII, the proposed ACS Mask R-CNN with pretraining significantly outperforms previous *state-of-the-art* MULAN [8]. Notably, we use only the detection and RECIST supervision, without additional information beyond the CT images such as tags from medical reports and demographic information. 3D context modeling with 3D and ACS convolutions is proven effective, especially with more input slices. Pretraining consistently improves the model performance for 2.5D, 3D and ACS convolutions. Large performance gap is observed for 3D **r.** and 3D **p.**, perhaps due to the large model

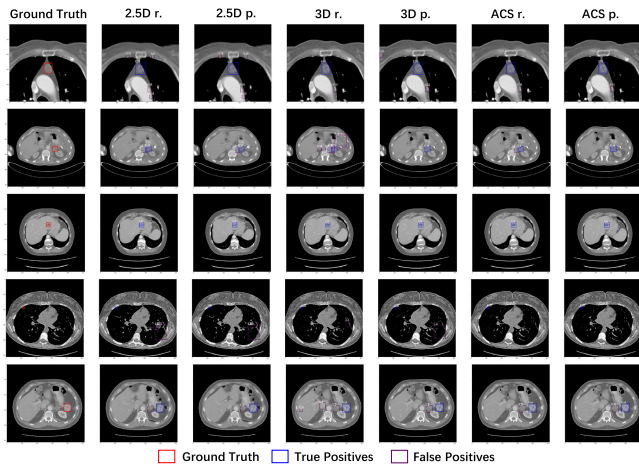


Fig. 6. Illustration of universal lesion detection on the DeepLesion Benchmark. Predicted results with 7-slice inputs are depicted. We also demonstrate the predicted segmentation contours and RECIST diameters in the figure.

size of 3D convolutions. We also visualize several examples of detection results in Fig. 6. Both 3D context modeling and pretraining reduce the predicted false positives.

V. CONCLUSION

We propose ACS convolution for 3D medical images, as a generic and plug-and-play replacement of standard 3D convolution. It enables pretraining from 2D images, which consistently provides significant performance boost in our experiments. Even without pretraining, the ACS convolution is comparable or even better than 3D convolution, with smaller model size and less computation. In further study, we will focus on optimal ACS kernel axis assignment and integration with other 2D-to-3D transfer learning operators.

REFERENCES

- [1] G. Litjens *et al.*, “A survey on deep learning in medical image analysis,” *Med. Image Analysis*, vol. 42, pp. 60–88, 2017.
- [2] R. Droste *et al.*, “Ultrasound image representation learning by modeling sonographer visual attention,” in *Proc. Int. Conf. Inf. Proc. Med. Imag.*, 2019, pp. 592–604.
- [3] K. Yan *et al.*, “Deep lesion graphs in the wild: Relationship learning and organization of significant radiology image findings in a diverse large-scale lesion database,” in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 9261–9270.
- [4] B. H. Menze *et al.*, “The multimodal brain tumor image segmentation benchmark (brats),” *IEEE Trans. Med. Imag.*, vol. 34, no. 10, pp. 1993–2004, Oct. 2014.
- [5] X. Wang, Y. Peng, L. Lu, Z. Lu, M. Bagheri, and R. M. Summers, “Chestx-ray8: Hospital-scale chest X-ray database and benchmarks on weakly-supervised classification and localization of common thorax diseases,” in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 2097–2106.
- [6] V. Gulshan *et al.*, “Development and validation of a deep learning algorithm for detection of diabetic retinopathy in retinal fundus photographs,” *JAMA*, vol. 316, no. 22, pp. 2402–2410, 2016.
- [7] F. Isensee *et al.*, “NNU-net: Self-adapting framework for u-net-based medical image segmentation,” *Nat. Methods*, 2020.
- [8] K. Yan *et al.*, “MULAN: Multitask universal lesion analysis network for joint lesion detection, tagging, and segmentation,” in *Int. Conf. Medical Image Comput. Comput.-Assisted Intervention*, Springer, 2019, pp. 194–202.
- [9] G. Balakrishnan, A. Zhao, M. R. Sabuncu, J. Guttag, and A. V. Dalca, “Voxelmorph: A learning framework for deformable medical image registration,” *IEEE Trans. Med. Imag.*, vol. 38, no. 8, pp. 1788–1800, Aug. 2019.
- [10] A. A. Setio *et al.*, “Validation, comparison, and combination of algorithms for automatic detection of pulmonary nodules in computed tomography images: The LUNA16 challenge,” *Med. Image Anal.*, vol. 42, pp. 1–13, 2017.
- [11] A. L. Simpson *et al.*, “A large annotated medical image dataset for the development and evaluation of segmentation algorithms,” 2019, *arXiv:1902.09063*.
- [12] N. Tajbakhsh, L. Jeyaseelan, Q. Li, J. Chiang, Z. Wu, and X. Ding, “Embracing imperfect datasets: A review of deep learning solutions for medical image segmentation,” *Med. Image Analysis*, vol. 63, 2020, Art. no. 101693.
- [13] K. Yan, Y. Peng, V. Sandfort, M. Bagheri, Z. Lu, and R. M. Summers, “Holistic and comprehensive annotation of clinically significant findings on diverse ct images: Learning from radiology reports and label ontology,” in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 8523–8532.
- [14] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2009, pp. 248–255.
- [15] T.-Y. Lin *et al.*, “Microsoft COCO: Common objects in context,” in *Eur. Conf. Comput. Vis.*, 2014, pp. 740–755.
- [16] H. R. Roth *et al.*, “New 2.5 D representation for lymph node detection using random sets of deep convolutional neural network observations,” in *Med. Image Comput. Comput.-Assist. Interv.*, 2014, pp. 520–527.
- [17] Q. Yu, L. Xie, Y. Wang, Y. Zhou, E. K. Fishman, and A. L. Yuille, “Recurrent saliency transformation network: Incorporating multi-stage visual cues for small organ segmentation,” in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 8280–8289.
- [18] T. Ni, L. Xie, H. Zheng, E. K. Fishman, and A. L. Yuille, “Elastic boundary projection for 3D medical image segmentation,” in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 2109–2118.
- [19] Ö. Çiçek, A. Abdulkadir, S. S. Lienkamp, T. Brox, and O. Ronneberger, “3D U-Net: Learning dense volumetric segmentation from sparse annotation,” in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, 2016, pp. 424–432.
- [20] F. Milletari, N. Navab, and S.-A. Ahmadi, “V-Net: Fully convolutional neural networks for volumetric medical image segmentation,” in *Proc. 4th Int. Conf. 3DVis. (3DV)*, 2016, pp. 565–571.
- [21] W. Zhao *et al.*, “3D deep learning from CT scans predicts tumor invasiveness of subcentimeter pulmonary adenocarcinomas,” *Cancer Res.*, vol. 78, no. 24, pp. 6881–6889, 2018.
- [22] X. Li, H. Chen, X. Qi, Q. Dou, C.-W. Fu, and P.-A. Heng, “H-DenseUNet: Hybrid densely connected UNet for liver and tumor segmentation from CT volumes,” *IEEE Trans. Med. Imaging*, vol. 37, no. 12, pp. 2663–2674, Dec. 2018.
- [23] Y. Xia, L. Xie, F. Liu, Z. Zhu, E. K. Fishman, and A. L. Yuille, “Bridging the gap between 2 d and 3 d organ segmentation with volumetric fusion net,” in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, 2018, pp. 445–453.
- [24] H. Zheng *et al.*, “A new ensemble learning framework for 3d biomedical image segmentation,” *Assoc. Adv. Artif. Intell.*, vol. 33, 2019, pp. 5909–5916.
- [25] A. Paszke *et al.*, “Automatic differentiation in PyTorch,” in *Neural Inf. Process. Syst. Autodiff Workshop*, 2017.
- [26] A. Prason, K. Petersen, C. Igel, F. Lauze, E. Dam, and M. Nielsen, “Deep feature learning for knee cartilage segmentation using a triplanar convolutional neural network,” *Med. Image Comput. Comput.-Assist. Interv.*, vol. 16 Pt, no. 2, pp. 246–53, 2013.
- [27] J. Ding, A. Li, Z. Hu, and L. Wang, “Accurate pulmonary nodule detection in computed tomography images using deep convolutional neural networks,” in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, Springer, 2017, pp. 559–567.
- [28] M. Perslev, E. B. Dam, A. Pai, and C. Igel, “One network to segment them all: A general, lightweight system for accurate 3d medical image segmentation,” in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, 2019, pp. 30–38.
- [29] Y. Xia *et al.*, “3d semi-supervised learning with uncertainty-aware multi-view co-training,” in *Proc. IEEE Winter Conf. Appl. Comp. Vis.*, 2020, pp. 3646–3655.
- [30] Y. Song *et al.*, “Learning 3d features with 2d cnns via surface projection for ct volume segmentation,” in *Proc. Int. Conf. Med. Imag. Comput. Comput.-Assist. Interv.*, 2020, pp. 176–186.

- [31] J. Long, E. Shelhamer, and T. Darrell, "Fully convolutional networks for semantic segmentation," *Proc. Conf. Comput. Vis. Pattern Recognit.*, pp. 3431–3440, 2015.
- [32] A. Esteva *et al.*, "Dermatologist-level classification of skin cancer with deep neural networks," *Nature*, vol. 542, no. 7639, pp. 115–118, 2017.
- [33] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal loss for dense object detection," in *Proc. Int. Conf. Comput. Vis.*, 2017, pp. 2980–2988.
- [34] L.-C. Chen, Y. Zhu, G. Papandreou, F. Schroff, and H. Adam, "Encoder-decoder with atrous separable convolution for semantic image segmentation," in *Proc. Eur. Conf. Comput. Vis.*, 2018, pp. 801–818.
- [35] K. Kamnitsas *et al.*, "Efficient multi-scale 3d CNN with fully connected CRF for accurate brain lesion segmentation," *Med. Image Analysis*, vol. 36, pp. 61–78, 2017.
- [36] Q. Dou *et al.*, "3d deeply supervised network for automated segmentation of volumetric medical images," *Med. Image Analysis*, vol. 41, pp. 40–54, 2017.
- [37] Z. Zhou, M. M. R. Siddiquee, N. Tajbakhsh, and J. Liang, "UNET: A nested U-net architecture for medical image segmentation," in *Deep Learn. Med. Imag. Anal. Multimodal Learn. Clin. Decis. Support*, Springer, 2018, pp. 3–11.
- [38] J. Zhang, Y. Xie, P. Zhang, H. Chen, Y. Xia, and C. Shen, "Light-weight hybrid convolutional network for liver tumor segmentation," in *Proc. Int. Joint Conf. Artif. Intell.*, 2019, pp. 4271–4277.
- [39] F. Gonda, D. Wei, T. Parag, and H. Pfister, "Parallel separable 3D convolution for video and volumetric data understanding," in *Proc. Brit. Mach. Vis. Conf.*, 2018.
- [40] H. Zheng *et al.*, "HFA-NET: 3d cardiovascular image segmentation with asymmetrical pooling and content-aware fusion," in *Int. Conf. Med. Image Comput. Comput.-Assist. Interv.*, Springer, 2019, pp. 759–767.
- [41] Z. Qiu, T. Yao, and T. Mei, "Learning spatio-temporal representation with pseudo-3 d residual networks," *Proc. Int. Conf. Comput. Vis.*, pp. 5534–5542, 2017.
- [42] C. Li, Q. Zhong, D. Xie, and S. Pu, "Collaborative spatiotemporal feature learning for video action recognition," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 7872–7881.
- [43] R. Sun *et al.*, "A radiomics approach to assess tumour-infiltrating cd8 cells and response to anti-pd-1 or anti-pd-11 immunotherapy: An imaging biomarker, retrospective multicohort study," *Lancet Oncol.*, vol. 19, no. 9, pp. 1180–1191, 2018.
- [44] K. He, R. Girshick, and P. Dollár, "Rethinking imagenet pre-training," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2019, pp. 4918–4927.
- [45] M. Raghu, C. Zhang, J. Kleinberg, and S. Bengio, "Transfusion: Understanding transfer learning with applications to medical imaging," in *Adv. Neural Inf. Process. Syst.*, 2019, pp. 3347–3357.
- [46] D. Hendrycks, K. Lee, and M. Mazeika, "Using pre-training can improve model robustness and uncertainty," in *Proc. ICML, ser. Proc. Mach. Learn. Res.*, K. Chaudhuri and R. Salakhutdinov, Eds., vol. 97. Long Beach, California, USA: PMLR, Jun. 09–15, 2019, pp. 2712–2721.
- [47] E. Gibson *et al.*, "NIFTYNET: A deep-learning platform for medical imaging," *Comput. Methods Programs Biomed.*, vol. 158, pp. 113–122, 2018.
- [48] S. Hussein, K. Cao, Q. Song, and U. Bagci, "Risk stratification of lung nodules using 3 d CNN-based multi-task learning," in *Proc. Int. Conf. Inf. Process. Med. Imag.*, Springer, 2017, pp. 249–260.
- [49] K. Hara, H. Kataoka, and Y. Satoh, "Can spatiotemporal 3d CNNs retrace the history of 2d CNNs and imagenet?," in *Proc. IEEE Conf. Comput. Vis. Pattern Recognit.*, 2018, pp. 6546–6555.
- [50] S. Chen, K. Ma, and Y. Zheng, "MED3D: Transfer learning for 3d medical image analysis," 2019, *arXiv:1904.00625*.
- [51] Z. Zhou *et al.*, "Models genesis: Generic autodidactic models for 3d medical image analysis," in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, 2019, pp. 384–393.
- [52] D. Berthelot, N. Carlini, I. Goodfellow, N. Papernot, A. Oliver, and C. Raffel, "MIXMATCH: A holistic approach to semi-supervised learning," in *Adv. Neural Inf. Process. Syst.*, 2019, pp. 5049–5059.
- [53] O. J. Hénaff, "Data-efficient image recognition with contrastive predictive coding," in *Int. Conf. Machine Learn. PMLR*, 2020, pp. 4182–4192.
- [54] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2016, pp. 770–778.
- [55] G. Huang, Z. Liu, L. Van Der Maaten, and K. Q. Weinberger, "Densely connected Convolutional networks," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 4700–4708.
- [56] J. Carreira and A. Zisserman, "Quo vadis, action recognition? A new model and the kinetics dataset," in *Proc. Conf. Comput. Vis. Pattern Recognit.*, 2017, pp. 6299–6308.
- [57] S. Liu *et al.*, "3d anisotropic hybrid network: Transferring convolutional features from 2d images to 3d anisotropic volumes," in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, Springer, 2018, pp. 851–858.
- [58] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," 2015, *arXiv:1510.00149*.
- [59] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Proc. Int. Conf. Med. Imag. Comput. Comput.-Assist. Interv.*, Springer, 2015, pp. 234–241.
- [60] S. G. Armato *et al.*, "The lung image database consortium (LIDC) and image database resource initiative (IDRI): A completed reference database of lung nodules on CT scans," *Med. Physics*, vol. 38, no. 2, pp. 915–931, 2011.
- [61] D. P. Kingma and J. Ba, "ADAM: A method for stochastic optimization," in *Proc. Int. Conf. Learn. Represent.*, 2014.
- [62] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *Proc. Int. Conf. Learn. Representations*, 2015.
- [63] Y. Xie, Y. Xia, J. Zhang, D. D. Feng, M. Fulham, and W. Cai, "Transferable multi-model ensemble for benign-malignant lung nodule classification on chest ct," in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, 2017, pp. 656–664.
- [64] L. Liu, Q. Dou, H. Chen, J. Qin, and P.-A. Heng, "Multi-task deep model with margin ranking loss for lung nodule analysis," *IEEE Trans. Med. Imag.*, vol. 39, no. 3, pp. 718–728, 2019.
- [65] P. Bilic *et al.*, "The liver tumor segmentation benchmark (LITS)," 2019, *arXiv:1901.04056*.
- [66] K. He, G. Gkioxari, P. Dollár, and R. B. Girshick, "Mask R-CNN," *Proc. Int. Conf. Comput. Vis.*, pp. 2980–2988, 2017.
- [67] K. Yan, M. Bagheri, and R. M. Summers, "3d context enhanced region-based convolutional neural network for end-to-end lesion detection," in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, 2018, pp. 511–519.
- [68] Y.-B. Tang, K. Yan, Y.-X. Tang, J. Liu, J. Xiao, and R. M. Summers, "ULDOR: A universal lesion detector for CT scans with pseudo masks and hard negative example mining," in *Proc. IEEE Int. Symp. Biomed. Imag.*, 2019, pp. 833–836.
- [69] X. Wang, S. Han, Y. Chen, D. Gao, and N. Vasconcelos, "Volumetric attention for 3d medical image segmentation and detection," in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, 2019, pp. 175–184.
- [70] M. Zlocha, Q. Dou, and B. Glocker, "Improving retinanet for CT lesion detection with dense masks from weak recist labels," in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, 2019, pp. 402–410.
- [71] Z. Li, S. Zhang, J. Zhang, K. Huang, Y. Wang, and Y. Yu, "MVP-NET: Multi-view FPN with position-aware attention for deep universal lesion detection," in *Proc. Med. Image Comput. Comput.-Assist. Interv.*, 2019, pp. 13–21.
- [72] K. Chen *et al.*, "MMDETECTION: Open mmlab detection toolbox and benchmark," 2019, *arXiv:1906.07155*.