

Deep Sequential Segmentation of Organs in Volumetric Medical Scans

Alexey A. Novikov, David Major, Maria Wimmer, Dimitrios Lenis, Katja Bühler

Abstract—Segmentation in 3D scans is playing an increasingly important role in current clinical practice supporting diagnosis, tissue quantification, or treatment planning. The current 3D approaches based on *Convolutional Neural Networks (CNN)* usually suffer from at least three main issues caused predominantly by implementation constraints - first, they require resizing the volume to the lower-resolution reference dimensions, second, the capacity of such approaches is very limited due to memory restrictions, and third, all slices of volumes have to be available at any given training or testing time. We address these problems by a U-Net-like [1] architecture consisting of bidirectional *Convolutional Long Short-Term Memory (C-LSTM)* [2] and convolutional, pooling, upsampling and concatenation layers enclosed into time-distributed wrappers. Our network can either process the full volumes in a sequential manner, or segment slabs of slices on demand. We demonstrate performance of our architecture on vertebrae and liver segmentation tasks in 3D CT scans.

I. INTRODUCTION

Accurate segmentation of anatomical structures in volumetric medical scans is of high interest in current clinical practice as it plays an important role in many tasks involved in computer-aided diagnosis, image-guided interventions, radiotherapy and radiology. In particular, quantitative diagnostics requires accurate boundaries of anatomical organs.

Computed tomography (CT) is currently among the most used 3D imaging modalities. Despite its inability of differentiating organs with similar intensities it is widely used for diagnosis of diseases in organs. Manual segmentation in CT can be a very tedious task. Therefore, automated methods with minor or no human interaction at all, are preferable.

Automated segmentation with deep learning methods in medical images has popularized widely in the recent years, mainly due to the success of applying *Fully-Convolutional Networks (FCN)* in natural images [3] and consequently in the biomedical imaging [1]. Since then various modifications of *FCNs* have been proposed for segmentation of different anatomical organs and imaging modalities.

3D scans are generally represented as stacks of 2D images. Running a segmentation algorithm on the 2D slices directly with merging results afterwards ignores spatial inter-slice

correlations, therefore hybrid 2D/3D and direct 3D approaches gained popularity. Most of these methods are built upon 2D [1] and 3D [4] U-Net architectures. Lu et al. [5] proposed to locate and segment the liver via convolutional neural networks and graph cuts. Dou et al. [6] presented a 3D *FCN* which boosts liver segmentation accuracy by deep supervision layers. Yang et al. [7] used adversarial training in order to gain in performance for the 3D U-Net segmentation of the liver in CT scans. Sekuboyina et al. [8] proposed a pipeline approach for both localization and segmentation of the spine in CT. Here the vertebrae segmentation is performed in a blockwise manner to overcome memory limitations as well as obtain a fine-grained result. A similar blockwise approach in combination with a multi-scale two-way *CNN* was introduced by Korez et al. [9].

Other noteworthy works using variants of 2D and 3D U-Nets consider applications in cardiac MR image segmentation [10], [11], pancreas in 3D CT [12], [13] and prostate in 3D MR [14], [15]. A variety of papers have contributed to the multiple tasks in brain imaging such as segmentation of cerebrospinal fluid, gray and white matter [16], brain tumour [17], [18], multiple sclerosis lesion [19] and glioblastoma [20].

In order to overcome memory limitations modern *CNN*-based methods are usually preceded by downsampling of the input scans. This might result in a deformation of organs in the image, causing information loss.

Consequently, hybrid 2D/3D methods that process volumetric data in a slice-wise fashion followed by a 3D processing step, gained importance. For instance, Li et al. [21] applied a slice-wise densely connected variant of the 2D U-Net architecture for liver segmentation first, and refined the result by the auto-context algorithm in 3D. For the same task, Christ et al. [22] applied slice-wise 2D U-Nets to obtain a rough segmentation first, and then tuned the result in a second step with Conditional Random Fields. Relying only on intra-slice data is insufficient for proper leveraging spatial information. In order to address this issue, the above-mentioned methods applied computationally expensive 3D classical image processing refinement strategies in addition to the 2D *CNN*-based approach.

Hybrid approaches combining *FCN* with recurrent networks such as *Long Short-Term Memory (LSTM)* [23] and more recently proposed *C-LSTM* [2] are effective for processing sequential data in general. Hence, the recurrent networks have recently been introduced to the biomedical imaging context. A method proposed by Poudel et al. [24] uses a U-Net variant to get an estimate of the 2D slice-wise segmentation, which is subsequently refined by the so-called gated recurrent unit [25], a simplified version of the *LSTM*.

A. A. Novikov, D. Major, M. Wimmer, D. Lenis and K. Bühler are with the VRVis Zentrum für Virtual Reality and Visualisierung Forschungs-GmbH, 1220 Vienna, Austria, e-mail: (novikov@vrvis.at, major@vrvis.at, mwimmer@vrvis.at, lenis@vrvis.at, buehler@vrvis.at). VRVis is funded by BMVIT, BMDW, Styria, SFG and Vienna Business Agency in the scope of COMET - Competence Centers for Excellent Technologies (854174) which is managed by FFG. Thanks go to our project partner AGFA HealthCare for valuable input. Copyright (c) 2018 IEEE. Personal use of this material is permitted. However, permission to use this material for any other purposes must be obtained from the IEEE by sending a request to pubs-permissions@ieee.org.

Bates et al. [26] evaluated several architectures involving both *C-LSTM* and standard convolutional layers. In the *deep* configuration, several bidirectional *C-LSTM* units were stacked in the U-shaped architecture in which the outputs of the forward and backward *LSTM* passes were concatenated. In the *shallow* configuration a shared copy of the CNN was applied to each slice of the 3D scans separately and then the result was passed to the stacked *C-LSTM* units to produce the segmentation volume. For the purpose of designing a multi-scale architecture, Chen et al. [27] used a variant of 2D U-Net to extract features for all slices first, and then processed them with bidirectional *C-LSTM* units in order to exploit 3D context.

Though the described approaches address some issues of the deep learning based 3D segmentation algorithms such as voxel size anisotropy and intensive computations due to 3D convolutions, they still do not take into account the two following issues. First, they require that all volumes have the same fixed input reference dimensions, and, second, all slices of the volumes have to be available in order to extract 3D context at both training and testing time. The former scenario is not always applicable usually due to large variations of the number of slices in the volumes across even the same dataset, and the latter one could force reducing network capacity due to memory and timing restrictions what could potentially lead to lower accuracies.

To overcome these problems we propose to integrate bidirectional *C-LSTMs* into a U-Net-like architecture in order to extract the 3D context of slices in a sequential manner. In this way, the network is able to learn the inter-slice correlations based on the slabs of the volume. The downsampling of the input is not required anymore as only a fraction of the volume is processed at any given time. Training of this network is therefore not demanding memory-wise which is another known limitation of the current modern networks. This fully integrated sequential approach can be particularly useful for real-time applications as it enables segmentation already during data acquisition or while loading the data as both are generally performed slice-by-slice.

Furthermore, we show the invariance of our method to field-of-view and orientation by evaluating on two CT datasets depicting two different organs, namely liver and vertebrae.

For the sake of simplicity of the further explanation in the following we refer to this architecture as *Sensor3D* (acronym for "sequential segmentation of organs in 3D").

II. METHODOLOGY

A. General setup

Let $\mathcal{I} = \{I_1, \dots, I_n\}$ be a set of $n \in \mathbb{N}$ volumetric scans, where each I_i , $i = 1, \dots, n$ consists of voxels $\mathbf{x} = (x_1, x_2, x_3) \in \mathbb{R}^3$ with intensities $I_i(\mathbf{x}) \in \mathcal{J} \subset \mathbb{R}$. More specifically, each scan I_i is therefore a set of $m_i \in \mathbb{N}$ slices J_k^i , $k = 1, \dots, m_i$ within the organ area where $J_k^i(\mathbf{y}) \in \mathcal{J}$ correspond to intensities at the pixels with positions $\mathbf{y} = (y_1, y_2) \in \mathbb{R}^2$ at the k -th slice of the scan $I_i \in \mathcal{I}$.

For each slice $J_k^i \in I_i$, a set of ground truth masks $M_k^i := (M_{k,l}^i)_{l=1}^m$ is available, where l corresponds to semantic class

labels $\mathcal{L} = \{l_1, \dots, l_m\}$ and $M_{k,l}^i \in \mathcal{M}$ to the space of all 2D binary masks of the same size as the slices J_k^i .

To enforce reproducibility of the input flow shapes, we build a new training dataset in the following way.

The spatial context C_k^i of the slice J_k^i is defined as a set containing the slice J_k^i and its $(\mathbf{o} - 1)/2$ neighbouring slices above and below, selected equidistantly with a rounded stepsize \mathbf{d} and the pre-defined length of the input sequence \mathbf{o} . Rounding to the more distant slice is performed if the precise step is not possible. Training set \mathcal{I}' is defined then as follows:

$$\mathcal{I}' = \{C_k^i \mid i = 1, \dots, n, k = 1, \dots, m_i\} \quad (1)$$

For training and evaluation purposes, the dataset \mathcal{I}' is split into non-overlapping sets, namely $\mathcal{I}'_{\text{TRAIN}}$ and $\mathcal{I}'_{\text{TEST}}$. During training, the network is consecutively passed through with minibatches $\mathcal{K} \in \mathcal{N}$, where \mathcal{N} is a complete partition of the set $\mathcal{I}'_{\text{TRAIN}}$.

For each spatial context $C_k^i \in \mathcal{I}'$, i.e. $C_k^i = \{J_p^i, \dots, J_q^i\}$ for some $1 \leq p, q \leq m_i$, the multi-class output of the network is calculated: understanding the network as a function

$$\mathcal{N} : \mathcal{I}' \rightarrow \mathcal{M}, \quad (2)$$

$\mathcal{N}(C_k^i)$ derives for each pixel $\mathbf{y} \in J_t^i$ its semantic class $l \in \mathcal{L}$ in a single step with some probability, where J_t^i corresponds to the middle element of the spatial context C_k^i . In order to estimate and maximize this probability, we define a loss function

$$\Lambda : \mathcal{I}' \times \mathcal{M} \rightarrow \mathbf{R} \quad (3)$$

that estimates the deviation (error) of the network outcome from the desired ground truth. Using the formal notations derived in our work [28] we define the loss function in the following way.

For a distance function $d : \mathcal{I}' \times \mathcal{M} \rightarrow \mathbf{R}$, weighting coefficients $r_{\mathcal{K},l}$ and a spatial context $C_k^i \in \mathcal{K}$ the loss function is

$$\Lambda(C_k^i, M_k^i) := - \sum_{l \in \mathcal{L}} r_{\mathcal{K},l}^{-1} d(C_k^i, M_k^i) \quad (4)$$

over the set \mathcal{K} and the complete partition.

The distance function d_l^{dice} for the Dice coefficient for a spatial context C_k^i , a feature channel l , ground-truth mask M_k^i and sigmoid activation function $p_l(\cdot)$ can then be defined as:

$$d_l^{\text{dice}}(C_k^i, M_k^i) := 2 \frac{\sum_{\mathbf{x} \in I} \chi_{\pi_l(M_k^i)}(\mathbf{x}) p_l(\mathbf{x})}{\sum_{\mathbf{x} \in I} (\chi_{\pi_l(M_k^i)}(\mathbf{x}) + p_l(\mathbf{x}))} \quad (5)$$

where $\chi_{\pi_l(M_k^i)}(\mathbf{x})$ is a characteristic function, i.e., $\chi_{\pi_l(M_k^i)}(\mathbf{x}) = 1$ iff M_k^i is 1 at position of pixel \mathbf{x} and 0 otherwise. The definition of the loss function in this equation would allow for using multiple classes, however, this is beyond the scope of this work.

B. Building the architecture

Following the above, a 3D volumetric scan I_i can be interpreted as a time-series of 2D slices $\{J_1, \dots, J_{m_i}\}$. Such series can then be processed using methods known for successful performance on sequential data. The time-distributed

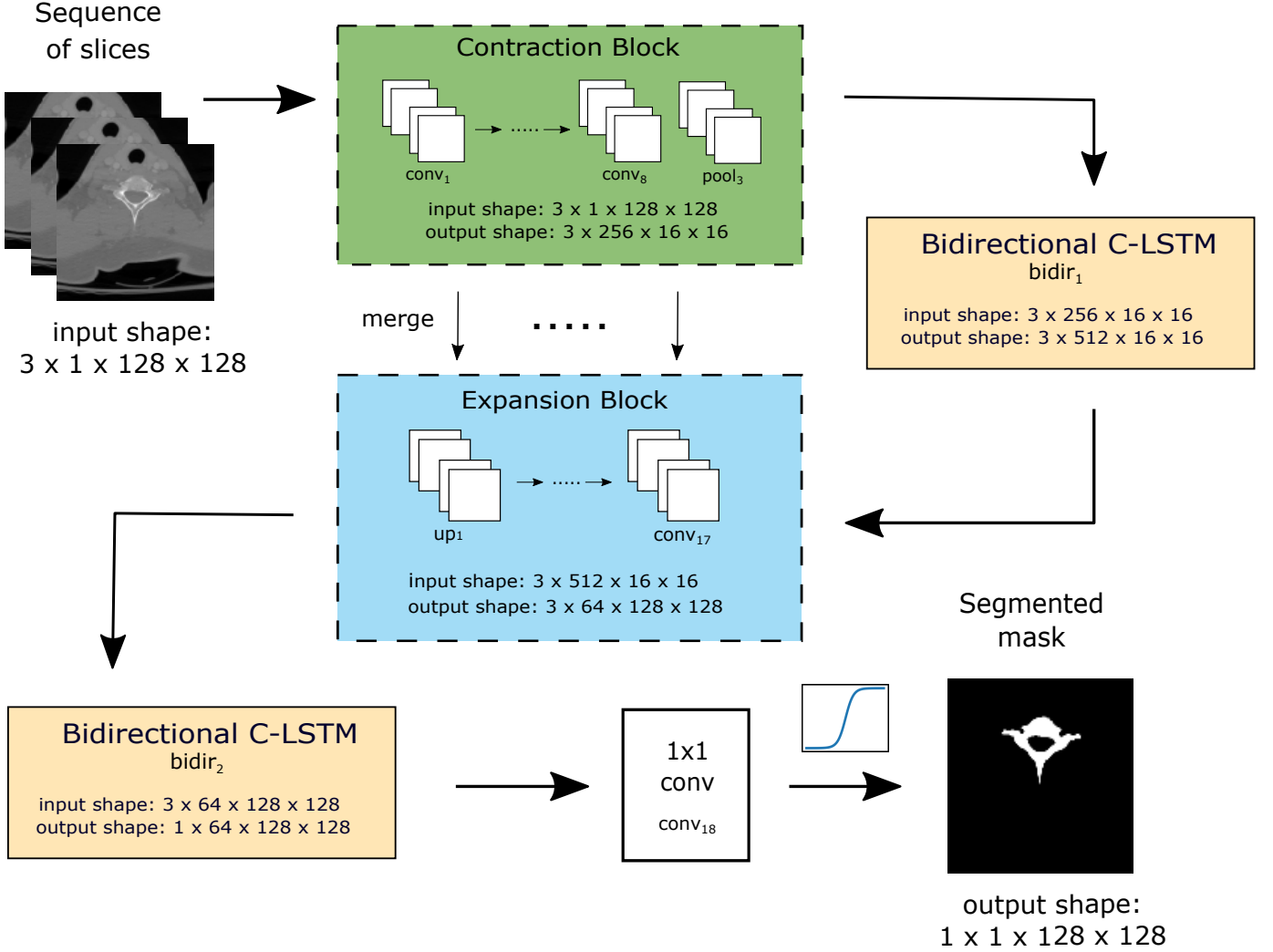


Fig. 1. Overview of the proposed *Sensor3D* architecture for a sample of three vertebrae slices and 128×128 imaging resolution used during training. Contraction and Expansion blocks are enclosed into time-distributed wrappers. Dashed merge connection corresponds to concatenations between layers of contraction and expansion blocks. The layer names in the network blocks correspond to entries in Table I.

convolutions and recurrent networks in particular are a natural choice for such 3D scans. Time-distributed convolutions are typical convolutions passed to a time-distributed wrapper that allows application of any layer to every temporal frame (or slice) of the input independently. In the context of this work such temporal frames correspond to the elements of training sequences extracted from the volumes. In our architecture the wrapper was applied to all convolutional, pooling, upsampling and concatenation layers.

In order to leverage spatio-temporal correlations of the order-preserving slices (that is elements of the C_k^i) and due to their sequential nature, we have combined the time-distributed layers and bidirectional *C-LSTMs* in an end-to-end trainable U-Net-like hybrid architecture. Main intuition for designing this architecture was that the features of the correlated slices should also be correlated. The *C-LSTMs* in our model are used to impose this correlation explicitly. To make training faster and to reduce the number of parameters, our *C-LSTM* blocks are based on the version of the *LSTM* without connections from the cell to the gates (widely known as "peephole con-

nections"). Motivation for using this variant was the research by Greff et al. [29] where it was shown that removing peephole connections in general does not hurt the overall performance.

Fig. I shows the high-level overview of the proposed architecture on a sample sequence of vertebrae slices. Table I complements the figure with tensor shapes for each layer for a particular case when the length of input sequences \mathfrak{o} is equal to three.

As mentioned previously, the network takes an odd-lengthed spatial context C_k^i as the input. This sequence is then passed to the contraction block (green in Fig. I and the corresponding layers from $conv_1$ to $pool_3$ in Table I). As all convolutional and max pooling layers are enclosed into a time-distributed wrapper, each element of the sequence is processed through the contraction block independently.

In order to capture spatio-temporal correlations between slices the features extracted for each element of the input sequence are passed into the *C-LSTM* block [2] at the end of the contraction part (layer $bidir_1$ in Table I). In order to enable the network to learn spatio-temporal correlations of the

slices in both directions, we used a bidirectional extension for the *C-LSTM* with the summation operator combining forward and backward outputs. This *C-LSTM* block aims at adding the explicit dependency of the low-dimensional high abstract features extracted for the elements of the sequence.

The sequence output of the bidirectional *C-LSTM* block is then passed to the expansion part (blue in Fig. I and the corresponding layers from up_1 to $conv_{17}$ in Table I). Similarly to the contraction part, each element of the sequence is processed independently via time-distributed convolutional as well as upsampling layers. After every upsampling layer, the features are concatenated with the corresponding features from the contraction part. When the spatial resolution of the features reaches the desired output sizes, the sequence is passed to another bidirectional *C-LSTM* block (layer $bidir_2$ in Table I). The sequence is processed in both directions and the outputs are combined by summation. At this stage this block contributes towards two goals: adding explicit dependency for the high-dimensional high-abstract features and converting the incoming sequence into a single-channelled output. The resulting features are then passed to the (1,1) convolution layer in order to map each feature vector to the desired number of classes (in the scope of this work the number of classes is equal to one). The output of the last convolutional layer (layer $conv_{18}$ in Table I) is mapped into [0,1] range via the sigmoid activation which is applied to each pixel independently. This results in the segmentation of the middle element of the spatial context C_k^i .

III. EXPERIMENTAL SETUP

To evaluate the performance and generalizability of our architecture we trained and tested it for 3D segmentation of two different anatomical organs: liver and vertebrae in CT scans. Liver segmentation is often a required step in the diagnosis of hepatic diseases while the segmentation of vertebrae is important for the identification of spine abnormalities, e.g. fractures, or image-guided spine intervention.

A. Training data and preparation

For *liver* segmentation we used two related datasets: 3Dircadb-01 and 3Dircadb-02 [30] combined together. The first consists of 20 3D CT scans with hepatic tumours in 75% cases. The second one consists of two anonymized scans with hepatic focal nodular hyperplasia. The axial in-plane resolution varied between 0.56 and 0.961 mm^2 and the slice thickness varied between 1.0 and 4.0 mm . The consecutive elements within the training sequences were generated at distances $\mathbf{d} \in \{3, 5, 7, 9\} mm$ within the liver area. These numbers were chosen based on the maximal slice thicknesses in the scans of the dataset. Unlike other existing liver datasets, 3Dircadb is more challenging due to the presence of multiple pathological cases with tumours both inside and close to the liver. The whole dataset with annotations of different organs is publicly available. Detailed per-scan information is available online [31].

We used a normalization technique similar to the one proposed by Christ et al. [22] which we applied to each

slice of the sequences independently. First, the raw slices were windowed to [-100, 400] to prevent including non-liver organs. Second, the contrast-limited adaptive histogram equalization was applied to the clipped slices. Third, the results were zero-centered by subtracting the slice-wise mean and then additionally normalized by scaling using the slice-wise standard deviation.

For *vertebrae* segmentation we used the CSI 2014 challenge train set [32]. It comprises 10 CT scans covering the entire lumbar and thoracic spine as well as full vertebrae segmentation masks for each scan. The axial in-plane resolution varies between 0.3125 and 0.3616 mm^2 . The slice thickness is 1 mm . The consecutive elements within the training sequences were generated at the distances of 1 mm within the vertebrae area.

In this work we focused on learning the 3D spatial context in a direct neighbourhood to the slices of interest only, thus in all evaluations we used sequences of three slices $\mathbf{o} = 3$. The design of the suggested architecture would allow for using larger sequences, however, this is beyond the scope of this work.

In order to prevent over-fitting for both liver and vertebrae segmentation tasks we made sure that every scan was first assigned either to the training or the testing set and only then converted into sequences. In this way, we ensured independence of the sets allowing us to estimate the generalizability of the algorithm.

All slices and their corresponding masks in the training set were downsampled to 128×128 in-plane imaging resolution. In order to compute the performance scores resulting masks were upsampled to the original 512×512 imaging resolution during testing.

B. Training strategies

We trained the networks in an end-to-end manner over the loss shown by Eq. 4 using the Adam [33] optimization algorithm with a fixed initial rate of 5×10^{-5} and the standard values of $\beta_1 = 0.9$ and $\beta_2 = 0.999$. Early stopping with the patience parameter equal to 100 epochs was used in all evaluations. Therefore, number of epochs varied between training runs.

The learning rate was chosen empirically based on the preliminary evaluations on smaller training sets. Higher learning rates caused the network training to diverge whereas lower ones slowed it down significantly.

We used zero-padding in convolutional layers and *C-LSTM* in the Sensor3D and its variants in all evaluation runs. Therefore, output channels of the layer had the same dimensions as the input.

Initialization with a random orthogonal matrix was used for the weights at the recurrent connections of the *C-LSTM* [34]. Glorot uniform [35] was utilized as an initialization for the weights at all other connections at the *C-LSTM* and at all convolutional layers.

As activation function at all convolutional layers we employed exponential linear units [36]. For the *C-LSTM* layers we used the widely used setup of hyperbolic tangent functions

TABLE I
DETAILED INFORMATION ON THE PROPOSED ARCHITECTURE WITH FILTERS AND SHAPES FOR INPUT AND OUTPUT TENSORS FOR THE CASE WHEN THE LENGTH OF INPUT SEQUENCES IS $\mathbf{o} = 3$ AND IN-PLANE IMAGING RESOLUTION IS 128×128

Layer Name	Layer Type	Input Shape	Filters	Output Shape
<i>conv</i> ₁	<i>Time-Distributed Convolutions</i>	$3 \times 1 \times 128 \times 128$	$3 \times 64 \times 3 \times 3$	$3 \times 64 \times 128 \times 128$
<i>conv</i> ₂	<i>Time-Distributed Convolutions</i>	$3 \times 64 \times 128 \times 128$	$3 \times 64 \times 3 \times 3$	$3 \times 64 \times 128 \times 128$
<i>pool</i> ₁	<i>Time-Distributed Max Pooling</i>	$3 \times 64 \times 128 \times 128$	$3 \times 2 \times 2$	$3 \times 64 \times 64 \times 64$
<i>conv</i> ₄	<i>Time-Distributed Convolutions</i>	$3 \times 64 \times 64 \times 64$	$3 \times 128 \times 3 \times 3$	$3 \times 128 \times 64 \times 64$
<i>conv</i> ₅	<i>Time-Distributed Convolutions</i>	$3 \times 128 \times 64 \times 64$	$3 \times 128 \times 3 \times 3$	$3 \times 128 \times 64 \times 64$
<i>pool</i> ₂	<i>Time-Distributed Max Pooling</i>	$3 \times 128 \times 64 \times 64$	$3 \times 2 \times 2$	$3 \times 128 \times 32 \times 32$
<i>conv</i> ₇	<i>Time-Distributed Convolutions</i>	$3 \times 256 \times 32 \times 32$	$3 \times 256 \times 3 \times 3$	$3 \times 256 \times 32 \times 32$
<i>conv</i> ₈	<i>Time-Distributed Convolutions</i>	$3 \times 256 \times 32 \times 32$	$3 \times 256 \times 3 \times 3$	$3 \times 256 \times 32 \times 32$
<i>pool</i> ₃	<i>Time-Distributed Max Pooling</i>	$3 \times 256 \times 32 \times 32$	$3 \times 2 \times 2$	$3 \times 256 \times 16 \times 16$
<i>bidir</i> ₁	<i>Bidirectional C-LSTM</i>	$3 \times 256 \times 16 \times 16$	$512 \times 3 \times 3$	$3 \times 512 \times 16 \times 16$
<i>up</i> ₁	<i>Time-Distributed Upsampling</i>	$3 \times 512 \times 16 \times 16$	$3 \times 2 \times 2$	$3 \times 512 \times 32 \times 32$
<i>concat</i> ₁	<i>Concatenation (conv₈, up₁)</i>			$3 \times 768 \times 32 \times 32$
<i>conv</i> ₁₁	<i>Time-Distributed Convolutions</i>	$3 \times 768 \times 32 \times 32$	$256 \times 3 \times 3$	$3 \times 256 \times 32 \times 32$
<i>conv</i> ₁₂	<i>Time-Distributed Convolutions</i>	$3 \times 256 \times 32 \times 32$	$256 \times 3 \times 3$	$3 \times 256 \times 32 \times 32$
<i>up</i> ₂	<i>Time-Distributed Upsampling</i>	$3 \times 256 \times 32 \times 32$	$3 \times 2 \times 2$	$3 \times 256 \times 64 \times 64$
<i>concat</i> ₂	<i>Concatenation (conv₅, up₂)</i>			$3 \times 384 \times 64 \times 64$
<i>conv</i> ₁₄	<i>Time-Distributed Convolutions</i>	$3 \times 384 \times 64 \times 64$	$128 \times 3 \times 3$	$3 \times 128 \times 64 \times 64$
<i>conv</i> ₁₅	<i>Time-Distributed Convolutions</i>	$3 \times 128 \times 64 \times 64$	$128 \times 3 \times 3$	$3 \times 128 \times 64 \times 64$
<i>up</i> ₃	<i>Time-Distributed Upsampling</i>	$3 \times 128 \times 64 \times 64$	$3 \times 2 \times 2$	$3 \times 128 \times 128 \times 128$
<i>concat</i> ₃	<i>Concatenation (conv₂, up₃)</i>			$3 \times 192 \times 128 \times 128$
<i>conv</i> ₁₇	<i>Time-Distributed Convolutions</i>	$3 \times 192 \times 128 \times 128$	$64 \times 3 \times 3$	$3 \times 64 \times 128 \times 128$
<i>bidir</i> ₂	<i>Bidirectional C-LSTM</i>	$3 \times 64 \times 128 \times 128$	$64 \times 3 \times 3$	$1 \times 64 \times 128 \times 128$
<i>conv</i> ₁₈	<i>2D Convolutions</i>	$1 \times 64 \times 128 \times 128$	$1 \times 1 \times 1$	$1 \times 1 \times 128 \times 128$

in all cases except the recurrent connections where the hard sigmoid was applied.

C. Implementation Details

All experiments were performed using Keras with TensorFlow backend in Python. The backend was used for automatic differentiation and optimization during training.

Downsampling of the ground-truth masks and upsampling of the segmentation masks were performed using the transform module of the *scikit-image* library.

D. Performance metrics

To evaluate the architectures and compare with state-of-the-art approaches, we used the Dice (D) similarity coefficient and volume overlap error (VOE), defined as follows.

Given an image I and the feature channel l , let $\pi_l(M_I)$ be a set of foreground pixels in the channel l of the ground-truth mask M_I and $P_l(I)$ be the set of pixels where the model is certain that they do not belong to the background, i.e.,

$$P_l(I) := \{\mathbf{x} : \mathbf{x} \in I \wedge |p_l(\mathbf{x}) - 1| < \epsilon\} \quad (6)$$

where $\epsilon = 0.25$ is an empirically chosen threshold value and $p_l(\mathbf{x})$ is the approximated probability of the pixel \mathbf{x} belonging to the foreground.

The coefficients D and VOE might then be computed in the following way:

$$D(I, M_I) := 2 \frac{|P_l(I) \cap \pi_l(M_I)|}{|P_l(I)| + |\pi_l(M_I)|} \quad (7)$$

$$VOE(I, M_I) = \frac{2(1 - D(I, M_I))}{2 - D(I, M_I)} \quad (8)$$

IV. RESULTS AND DISCUSSION

A. Evaluations with different inter-slice distances

Table II depicts the average Dice and volume overlap error scores for two folds of liver segmentation at different inter-slice distances \mathbf{d} . As expected, some irrelevant structures were partially segmented outside of the liver in a few cases thus lowering the scores when the full stack of volume slices is being considered.

The achieved results demonstrate that considering higher inter-slice distances is needed in order to get better segmentation performance. The lower scores for the 3 mm inter-slice distance are caused by some scans in both the training and testing data where slice thicknesses exceed 3 mm. In such scans the extracted sequences may contain direct-consecutive slices therefore adding disturbance in the training by giving the network a wrong impression that the elements in the sequences are not really different. Thus, hindering the network to learn the inter-slice context for those training sequences properly.

We additionally analysed how the segmentation results of the models with sequences generated at various distances (shown in Table II) differ in terms of statistical significance test scores. We performed pair-wise significance analysis using Wilcoxon signed-rank test for Dice scores on the test set. The results are shown in Table III where the entries with values less than 0.01 correspond to pairs of models demonstrating statistically different significance in segmentation performance. Thus, the numbers complement and confirm the detailed results provided in Table II: considering sequences of slices at the distances larger than 3 mm improves performance for the liver segmentation task significantly.

Some segmentation results at different vertebrae (top) and liver (bottom) areas are depicted in Fig. 2. The red contour

TABLE II
DETAILED SEGMENTATION RESULTS OF TWO-FOLD EVALUATIONS OF LIVER SEGMENTATION TASK FOR DIFFERENT INTER-SLICE DISTANCES

Step size	Fold 1				Fold 2			
	Organ Area		Full Volume		Organ Area		Full Volume	
	D (%)	VOE (%)	D (%)	VOE (%)	D (%)	VOE (%)	D (%)	VOE (%)
3 mm	94.8	9.8	92.8	13.4	95.1	9.4	93.7	11.8
5 mm	95.5	8.6	94.1	11.1	96.1	7.5	95.6	8.4
7 mm	95.3	8.9	94.3	10.8	96.4	6.9	96.2	7.3
9 mm	95.5	8.6	94.6	10.2	96.4	6.9	96.2	7.3

TABLE III
THE SIGNIFICANCE DIFFERENCE ANALYSIS OF SEGMENTATION RESULTS USING WILCOXON SIGNED-RANK TEST FOR DICE SCORES ON THE TEST SET FOR THE LIVER SEGMENTATION TASK. THE P-VALUES ARE GIVEN FOR FOLD 1 AND FOLD 2 (SEPARATED BY "SLASH" SIGN)

	3 mm	5 mm	7 mm	9 mm
3 mm	∞	$< 0.01 / < 0.01$	$< 0.01 / < 0.01$	$< 0.01 / < 0.01$
5 mm	$< 0.01 / < 0.01$	∞	0.17 / 0.08	0.13 / 0.07
7 mm	$< 0.01 / < 0.01$	0.17 / 0.08	∞	0.85 / 0.3
9 mm	$< 0.01 / < 0.01$	0.13 / 0.07	0.85 / 0.3	∞

TABLE IV
DETAILED SEGMENTATION RESULTS OF TWO-FOLD EVALUATIONS FOR ARCHITECTURES WITH DIFFERENT NUMBER OF FEATURES IN THE CONVOLUTIONAL LAYERS AND C -LSTM FOR THE LIVER SEGMENTATION TASK

	Fold 1				Fold 2			
	Organ Area		Full Volume		Organ Area		Full Volume	
	D (%)	VOE (%)	D (%)	VOE (%)	D (%)	VOE (%)	D (%)	VOE (%)
Original configuration	95.3	8.9	94.3	10.8	96.4	6.9	96.2	7.3
2 \times smaller	95.3	8.9	93.9	11.5	96.2	7.3	95.9	7.9
4 \times smaller	94.5	10.4	93.6	12.0	95.6	8.4	95.4	8.8
8 \times smaller	94.3	10.8	92.6	13.8	94.6	10.2	94.3	10.8

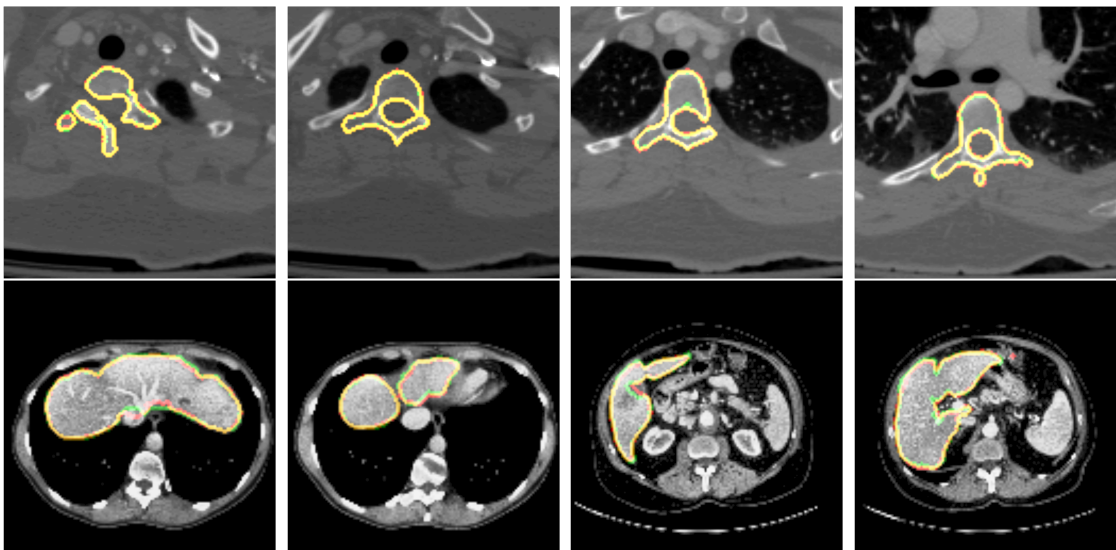


Fig. 2. Several visual examples of segmentations results in different vertebrae (top) and liver (bottom) locations. The contour in red corresponds to the outline of the prediction, green to the ground-truth and yellow to the overlap of the outlines

corresponds to the outline of the prediction, green to the ground-truth and yellow to the overlap of the outlines.

B. Evaluations on the influence of the network capacity

Table IV shows the detailed segmentation results of two-fold evaluations for architectures with different numbers of features in the convolutional layers and the C -LSTM units. The original configuration corresponds to the architecture shown in Table I. Two, four and eight times smaller configurations correspond

to the architectures where the number of feature maps in the convolutional layers and C -LSTM blocks is two, four or eight times less than in the original configuration.

Performance results demonstrate that reducing capacity of the network by two times slightly worsens the results, however, the number of parameters in this configuration is almost four times less so that training time can be reduced significantly. Making the configuration even smaller inevitably worsens results especially in a more challenging Fold 1 where the test

TABLE V

OUR METHOD COMPARED WITH STATE-OF-THE-ART METHODS ON THE LIVER SEGMENTATION ON 3DIRCADB (LEFT) AND VERTEBRAE SEGMENTATION ON CSI 2014 (RIGHT) DATASETS; "*" THE SCORE ESTIMATED USING EQ. 7 OR EQ. 8; "***" THE AREA OF VERTEBRAE AVAILABLE IN GROUND-TRUTH DATA

Method	D (%)	VOE (%)
Christ et al. [22]	94.3	10.7
Erdt et al. [37]	94.6 (*)	10.3
Li et al. [38]	94.5	10.4 (*)
Li et al. [39]	95.2 (*)	9.15
Lu et al. [5]	95.0 (*)	9.36
Sensor3D (full volume)	95.4	8.79
Sensor3D (liver area)	95.9	7.87

Method	D (%)	VOE (%)
Castro-Mateos et al. [40]	88.0	21.4 (*)
Forsberg et al. [41]	94.0	11.3 (*)
Hammernik et al. [42]	93.0	13.1 (*)
Korez et al. [43]	93.0	13.1 (*)
Seitel et al. [44]	83.0	29.1 (*)
Sensor3D (full volume)	93.1	12.9
Sensor3D (vertebrae area **)	94.9	9.7

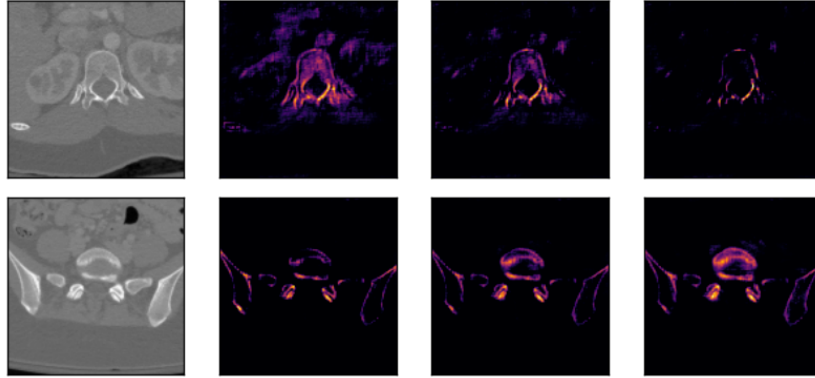


Fig. 3. Examples of features extracted after the penultimate upsampling step (after up_3 layer in Table I) for two sample contexts containing the same repeated slice

data consists of only a few challenging cases with multiple tumours inside or/and in a close proximity of the liver area.

C. Performance of variants of the Sensor3D network

1) *2D modifications*: In order to demonstrate that our Sensor3D network improves over similar 2D variants, we built and evaluated two additional architectures under the same training conditions on both folds in the liver segmentation task. *In the first* architecture, we set $\alpha = 1$, thus, changed the input in a way that the network is fed with a sequence of single slices without context. *In the second* architecture we did not change the slab size but removed the first *C-LSTM* and replaced the second one by the aggregation layer which would sum the incoming features along the time channel. Both architectures achieved similar average Dice scores of 84.3 % and 85.6 % (computed over two folds) when considering the organ area only. For the full scan scores of 73.1 % and 74.5% were achieved which are similar to the results of the U-Net performance reported by Christ et al. [22]. These scores are notably lower than the results demonstrated by Sensor3D. It shows that learning 3D context is crucial for achieving a better performance.

2) *Unidirectional modification*: We built and evaluated a unidirectional modification of the Sensor3D architecture under the same training conditions on both folds in the liver segmentation task. In this architecture we have replaced $bidir_1$ and $bidir_2$ layers in Table I with unidirectional *C-LSTM* blocks. The model achieved the Dice score of 93.53 % when considering the organ area only and 91.50 % in

the full volume. The achieved scores are significantly lower than the ones reached by both the state-of-the-art methods and Sensor3D on the same task in particular which hardens the assumption that bidirectional modification is beneficial for this architecture.

D. Comparison with state-of-the-art methods

Table V (left) compares our approach with the state-of-the-art methods trained and tested on the same 3Dircadb dataset. Though our model is trained only on the parts of the volumes where the liver is present (from 33% to 95% of slices in different scans, and in 71% of slices on average across all scans), it can still reach competitive and in many cases better results when evaluated against 2D and 3D approaches considering both the liver area and the full volume.

To demonstrate that our method generalizes on other organs as well, we have trained and evaluated the network on the CSI 2014 dataset on the vertebrae segmentation task. Table V (right) compares performance of our approach with several state-of-the-art methods. It is worth noting that some vertebrae which are not present in the ground-truth annotations are still segmented by our network thus causing lower scores in the cases when the full volume is considered.

E. Visual feature inspection

In order to visually demonstrate the sequential nature of the features learnt by our model, we performed the following test. We passed two sequences to the network (both for vertebrae), each containing three identical slices (first column in Fig. 3).

The columns show some of the features extracted after the penultimate upsampling step (after up_3 layer in Table I) before passing them to the final bidirectional *C-LSTM* block. The visualization shows that the layers respond differently to the same input element, activating different parts of the organ of interest. The brighter colour intensities correspond to higher activations. Comparing the rows, it shows that the network is able to learn spatial correlations in both directions.

V. CONCLUSIONS

In this paper we proposed Sensor3D, a general, robust, end-to-end U-Net-like hybrid architecture combining time-distributed convolutions, pooling, upsampling layers and bidirectional *C-LSTM* blocks. To demonstrate generalization of our approach, we evaluated the model on liver and additionally vertebrae segmentation task on the publicly available 3Dircadb and CSI 2014 datasets. Quantitative evaluations of the 2D variants of the Sensor3D network, statistical significance test, evaluation on the network capacity indicate that the *C-LSTM* boosts overall performance. Visual inspection of the model activation on the sequences containing the same repeated slices shows firing of different areas in the organs therefore empirically proving the sequential nature of the learnt features. Contrary to the state-of-the-art models, our network does not require full input volumes for neither training nor inference. Our network shows competitive and often superior performance on the considered liver and vertebrae segmentation tasks despite that it was trained only on slabs of the training volumes. For future work, we plan to apply our algorithm to other imaging modalities and organs in a multi-task manner.

REFERENCES

- [1] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional networks for biomedical image segmentation. In *MICCAI*, pages 234–241. Springer, 2015.
- [2] Xingjian Shi, Zhouong Chen, Hao Wang, Dit-Yan Yeung, Wai-kin Wong, and Wang-chun Woo. Convolutional LSTM network: A machine learning approach for precipitation nowcasting. In *NIPS*, pages 802–810. MIT Press, 2015.
- [3] Jonathan Long, Evan Shelhamer, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *CVPR*, 2015.
- [4] Özgün Çiçek, Ahmed Abdulkadir, Soeren S. Lienkamp, Thomas Brox, and Olaf Ronneberger. 3D U-Net: Learning dense volumetric segmentation from sparse annotation. In *MICCAI*, pages 424–432. Springer, 2016.
- [5] Fang Lu, Fa Wu, Peijun Hu, Zhiyi Peng, and Dexing Kong. Automatic 3D liver location and segmentation via convolutional neural network and graph cut. *International Journal of CARS*, 12(2):171–182, 2017.
- [6] Qi Dou, Hao Chen, Yueming Jin, Lequan Yu, Jing Qin, and Pheng-Ann Heng. 3D deeply supervised network for automatic liver segmentation from CT volumes. In *MICCAI*, pages 149–157. Springer, 2016.
- [7] Dong Yang, Daguang Xu, S Kevin Zhou, Bogdan Georgescu, Mingqing Chen, Sasa Grbic, Dimitris Metaxas, and Dorin Comaniciu. Automatic liver segmentation using an adversarial image-to-image network. In *MICCAI*, pages 507–515. Springer, 2017.
- [8] Anjany Sekuboyina, Jan Kukačka, Jan S Kirschke, Bjoern H Menze, and Alexander Valentinitich. Attention-driven deep learning for pathological spine segmentation. In *International Workshop and Challenge on Computational Methods and Clinical Applications in Musculoskeletal Imaging*, pages 108–119. Springer, 2017.
- [9] Robert Korez, Boštjan Likar, Franjo Pernuš, and Tomaž Vrtovec. Segmentation of pathological spines in CT images using a two-way CNN and a collision-based model. In *International Workshop and Challenge on Computational Methods and Clinical Applications in Musculoskeletal Imaging*, pages 95–107. Springer, 2017.
- [10] Christian F. Baumgartner, Lisa M. Koch, Marc Pollefeys, Ender Konukoglu, Maxime Sermesant, Pierre-Marc Jodoin, Alain Lalonde, Xiahai Zhuang, Guang Yang, Alistair Young, and Olivier Bernard. An exploration of 2D and 3D deep learning techniques for cardiac MR image segmentation. In *Statistical Atlases and Computational Models of the Heart. ACDC and MMWHS Challenges*, pages 111–119. Springer, 2018.
- [11] Lequan Yu, Jie-Zhi Cheng, Qi Dou, Xin Yang, Hao Chen, Jing Qin, Pheng-Ann Heng, Lena Maier-Hein, Alfred Franz, Pierre Jannin, D. Louis Collins, and Simon Duchesne. Automatic 3D cardiovascular MR segmentation with densely-connected volumetric convnets. In *MICCAI*, pages 287–295. Springer, 2017.
- [12] Yuyin Zhou, Lingxi Xie, Elliot K. Fishman, and Alan L. Yuille. Deep supervision for pancreatic cyst segmentation in abdominal CT scans. In *MICCAI*, pages 222–230. Springer, 2017.
- [13] Mattias P. Heinrich, Ozan Oktay, Lena Maier-Hein, Alfred Franz, Pierre Jannin, D. Louis Collins, and Simon Duchesne. BRIEFnet: Deep pancreas segmentation using binary sparse convolutions. In *MICCAI*, pages 329–337. Springer, 2017.
- [14] Lequan Yu, Xin Yang, Hao Chen, Jing Qin, and Pheng Ann Heng. Volumetric convnets with mixed residual connections for automated prostate segmentation from 3D MR images, 2017.
- [15] Qiuhua Liu, Min Fu, Hao Jiang, and Xinqi Gong. Volumetric densely dilated spatial pooling convnets for prostate segmentation. *CoRR*, abs/1801.10517, 2018.
- [16] Hao Chen, Qi Dou, Lequan Yu, Jing Qin, and Pheng-Ann Heng. Voxresnet: Deep voxelwise residual networks for brain segmentation from 3D MR images. *NeuroImage*, 170:446–455, 2018.
- [17] Hao Dong, Guang Yang, Fangde Liu, Yuanhan Mo, and Yike Guo. Automatic brain tumor detection and segmentation using U-Net based fully convolutional networks. In *Annual Conference on Medical Image Understanding and Analysis*, pages 506–517. Springer, 2017.
- [18] Haocheng Shen, Ruixuan Wang, Jianguo Zhang, and Stephen J McKenna. Boundary-aware fully convolutional network for brain tumor segmentation. In *MICCAI*, pages 433–441. Springer, 2017.
- [19] Tom Brosch, Lisa YW Tang, Youngjin Yoo, David K.B. Li, Anthony Trabulsee, and Roger Tam. Deep 3D convolutional encoder networks with shortcuts for multiscale feature integration applied to multiple sclerosis lesion segmentation. *IEEE Transactions on Medical Imaging*, 35(5):1229–1239, 2016.
- [20] Darvin Yi, Mu Zhou, Zhao Chen, and Olivier Gevaert. 3D convolutional neural networks for glioblastoma segmentation. *CoRR*, abs/1611.04534, 2016.
- [21] Xiaomeng Li, Hao Chen, Xiaojuan Qi, Qi Dou, Chi-Wing Fu, and Pheng Ann Heng. H-DenseUNet: Hybrid densely connected U-Net for liver and liver tumor segmentation from CT volumes. *arXiv preprint arXiv:1709.07330*, 2017.
- [22] Patrick Ferdinand Christ, Mohamed Ezzeldin A Elshaer, Florian Ertlinger, Sunil Tataavarty, Marc Bickel, Patrick Bilic, Markus Rempfler, Marco Armbruster, Felix Hofmann, Melvin DANastasi, et al. Automatic liver and lesion segmentation in CT using cascaded fully convolutional neural networks and 3D conditional random fields. In *MICCAI*, pages 415–423. Springer, 2016.
- [23] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):1735–1780, 1997.
- [24] Rudra P. K. Poudel, Pablo Lamata, Giovanni Montana, Kanwal Bhatia, Bernhard Kainz, Mehdi H. Moghari, and Danielle F. Pace. Recurrent fully convolutional neural networks for multi-slice MRI cardiac segmentation. In *Reconstruction, Segmentation, and Analysis of Medical Images*, pages 83–94. Springer, 2017.
- [25] Kyunghyun Cho, Bart van Merriënboer, Çağlar Gülçehre, Fethi Bougares, Holger Schwenk, and Yoshua Bengio. Learning phrase representations using RNN encoder-decoder for statistical machine translation. *CoRR*, abs/1406.1078, 2014.
- [26] Russell Bates, Benjamin Irving, Bostjan Markelc, Jakob Kaeppeler, Ruth Muschel, Vicente Grau, and Julia A Schnabel. Extracting 3D vascular structures from microscopy images using convolutional recurrent networks. *arXiv preprint arXiv:1705.09597*, 2017.
- [27] Jianxu Chen, Lin Yang, Yizhe Zhang, Mark Alber, and Danny Z Chen. Combining fully convolutional and recurrent neural networks for 3D biomedical image segmentation. In *NIPS*, pages 3036–3044, 2016.
- [28] A. A. Novikov, D. Lenis, D. Major, J. Hladůvka, M. Wimmer, and K. Bühler. Fully convolutional architectures for multi-class segmentation in chest radiographs. *IEEE Transactions on Medical Imaging*, PP(99):1–1, 2018.

- [29] Klaus Greff, Rupesh Kumar Srivastava, Jan Koutník, Bas R. Steunebrink, and Jürgen Schmidhuber. LSTM: A search space odyssey. *CoRR*, abs/1503.04069, 2015.
- [30] L Soler, A Hostettler, V Agnus, A Charnoz, JB Fasquel, J Moreau, A Osswald, M Bouhadjar, and J Marescaux. 3D image reconstruction for comparison of algorithm database: a patient-specific anatomical and medical image database, 2010.
- [31] 3dircadb database. Available at <https://www.ircad.fr/research/3dircadb/>.
- [32] Jianhua Yao, Joseph E. Burns, Hector Munoz, and Ronald M. Summers. Detection of vertebral body fractures based on cortical shell unwrapping. In *MICCAI*, pages 509–516. Springer, 2012.
- [33] Diederik Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- [34] Andrew M. Saxe, James L. McClelland, and Surya Ganguli. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks. *CoRR*, abs/1312.6120, 2013.
- [35] Xavier Glorot and Yoshua Bengio. Understanding the difficulty of training deep feedforward neural networks. In *Proceedings of the thirteenth international conference on artificial intelligence and statistics*, pages 249–256, 2010.
- [36] Djork-Arné Clevert, Thomas Unterthiner, and Sepp Hochreiter. Fast and accurate deep network learning by exponential linear units (ELUs). *arXiv preprint arXiv:1511.07289*, 2015.
- [37] M. Erdt and M. Kirschner. Fast automatic liver segmentation combining learned shape priors with observed shape deviation. In *2010 IEEE 23rd International Symposium on Computer-Based Medical Systems (CBMS)*, pages 249–254, 2010.
- [38] C. Li, X. Wang, S. Eberl, M. Fulham, Y. Yin, J. Chen, and D. D. Feng. A likelihood and local constraint level set model for liver tumor segmentation from CT volumes. *IEEE Transactions on Biomedical Engineering*, 60(10):2967–2977, 2013.
- [39] G. Li, X. Chen, F. Shi, W. Zhu, J. Tian, and D. Xiang. Automatic liver segmentation based on shape constraints and deformable graph cut in CT images. *IEEE Transactions on Image Processing*, 24(12):5315–5329, 2015.
- [40] Isaac Castro-Mateos, Jose M Pozo, Aron Lazary, and Alejandro Frangi. 3D vertebra segmentation by feature selection active shape model. In *Recent advances in computational methods and clinical applications for spine imaging*, pages 241–245. Springer, 2015.
- [41] Daniel Forsberg. Atlas-based segmentation of the thoracic and lumbar vertebrae. In *Recent Advances in Computational Methods and Clinical Applications for Spine Imaging*, pages 215–220. Springer, 2015.
- [42] Kerstin Hammernik, Thomas Ebner, Darko Stern, Martin Urschler, and Thomas Pock. Vertebrae segmentation in 3D CT images based on a variational framework. In *Recent advances in computational methods and clinical applications for spine imaging*, pages 227–233. Springer, 2015.
- [43] Robert Korez, Bulat Ibragimov, Boštjan Likar, Franjo Pernuš, and Tomaž Vrtovec. Interpolation-based shape-constrained deformable model approach for segmentation of vertebrae from CT spine images. In *Recent Advances in Computational Methods and Clinical Applications for Spine Imaging*, pages 235–240. Springer, 2015.
- [44] A Seitel, A Rasoulilian, R Rohling, and P Abolmaesumi. Lumbar and thoracic spine segmentation using a statistical multi-object shape+pose model. In *Recent Advances in Computational Methods and Clinical Applications for Spine Imaging*, pages 221–225. Springer, 2015.