# CATS: COMPLEMENTARY CNN AND TRANSFORMER ENCODERS FOR SEGMENTATION

*Hao Li, Dewei Hu, Han Liu, Jiacheng Wang, Ipek Oguz*

Department of Electrical Engineering and Computer Science, Vanderbilt University

## ABSTRACT

Recently, deep learning methods have achieved state-of-the-art performance in many medical image segmentation tasks. Many of these are based on convolutional neural networks (CNNs). For such methods, the encoder is the key part for global and local information extraction from input images; the extracted features are then passed to the decoder for predicting the segmentations. In contrast, several recent works show a superior performance with the use of transformers, which can better model long-range spatial dependencies and capture low-level details. However, transformer as sole encoder underperforms for some tasks where it cannot efficiently replace the convolution based encoder. In this paper, we propose a model with double encoders for 3D biomedical image segmentation. Our model is a U-shaped CNN augmented with an independent transformer encoder. We fuse the information from the convolutional encoder and the transformer, and pass it to the decoder to obtain the results. We evaluate our methods on three public datasets from three different challenges: BTCV, MoDA and Decathlon. Compared to the state-of-the-art models with and without transformers on each task, our proposed method obtains higher Dice scores across the board.

***Index Terms***— Convolutional neural network, Transformer, Medical image segmentation

## 1. INTRODUCTION

In recent years, convolutional neural networks (CNNs) with U-shaped structures have dominated the medical image segmentation field [1, 2, 3]. The U-shaped networks consist of an encoder and a decoder, with skip connections in between. The encoder extracts information by consecutive convolution and down-sampling operations. The encoded information is sent to the decoder via skip connections to obtain the segmentation result. The U-Net [1] and its many variants have shown great performance on many medical image segmentation tasks [4, 5, 6, 7].

However, convolutional encoders are somewhat limited for modeling long-range dependencies, due to the local receptive field of the convolution kernels. A potential solution is the transformer, which was originally proposed in the context of nature language processing, and has been used for image segmentation and classification [8, 9, 10, 11, 12] due

to its ability to better capture global information and modeling long-range context. For medical images, Chen et al. proposed the transUnet [13], which is the first segmentation framework with transformers. They added a transformer layer to the CNN-based encoder to better leverage global information. The UNEt TRansformers (UNETR) [14] is an architecture that completely replaces the CNN-based encoder with a transformer. The UNETR passes the multi-level scale information from this transformer to CNN-based decoder, and it is a state-of-the-art model for multi-organ segmentations on computed tomography (CT) images. Similar to U-Net, Cao et al. proposed Swin-Unet [15] which is a pure transformer-based U-shaped architecture for medical image segmentation. Since the transformer has limited ability to encode high-level information, these networks with transformers may not work well for some medical image segmentation tasks. To overcome this, Zhang et al. proposed a 2D architecture that combines a CNN-based encoder and a transformer-based segmentation network in parallel [16].

In this paper, we propose CATS (Complementary CNN and Transformer Encoders for Segmentation), a U-shaped architecture with double encoders. Inspired by UNETR [14] and TransFuse [16], we use a transformer as an independent encoder in addition to the CNN encoder. However, unlike the TransFuse, we use multi-scale features from the transformer rather than only the highest level features (i.e., the transformer output). The proposed CATS is a straightforward way to combine CNN and transformer without requiring a complex network architecture, such as the attention blocks and the BiFusion module in TransFuse. We use a 3D architecture and train from scratch instead of pre-training. Unlike the UNETR, we include both a transformer-based encoder and a CNN-based encoder. Multi-scale features extracted from the transformer are added with CNN features. The fused information is delivered to the CNN-based decoder for segmentation. We compare our model to state-of-the-art models with and without transformers on three public datasets.

## 2. METHODS

### 2.1. Framework overview

The proposed CATS framework is shown as Fig. 1. Our model contains two encoder paths, a CNN path and a trans-
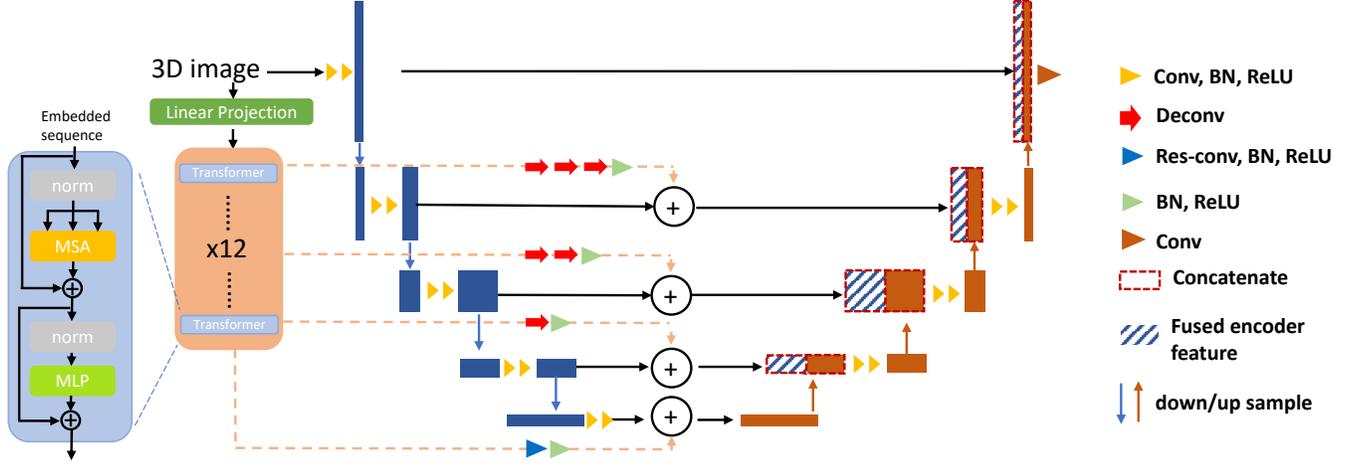
**Fig. 1**. Proposed network architecture with two independent encoder paths: a CNN-based encoder and a transformer encoder.

former path. For the CNN-based encoder, the information is gradually coded by the convolution and down-sampling operations. For the transformer path, to make sure the low-level details are well-preserved, we directly send the raw input to the transformer. Then, the information from the two paths are fused by addition operations at each level, and delivered to the CNN-based decoder to predict the final segmentation.

## 2.2. Transformer

We begin with an input 3D image $x \in \mathbb{R}^{C \times W \times H \times D}$, where $W$, $H$ and $D$ are the image dimensions and $C$ is the number of channels. We construct $x_p = \{x_p^i | i \in [1, N]\}$, which is a set consisting of $N$ non-overlapping patches $x_p^i \in \mathbb{R}^{(P^3 \times C)}$, where $N = \frac{W \times H \times D}{P^3}$ and each $x_p^i$ is a patch of $P^3$ voxels with $C$ channels. Next, we send the patches to the linear projection layer to obtain the embedded projection $\boldsymbol{E}$ with dimension $M = P^3 C$. To keep the position information, we add the position embedding $\boldsymbol{E_p}$ to form the transformer input:

$$z_0 = [x_p^1 \boldsymbol{E}; x_p^2 \boldsymbol{E}; ...; x_p^N \boldsymbol{E}] + \boldsymbol{E_p} \tag{1}$$

where $\boldsymbol{E} \in \mathbb{R}^{(P^3 \times C) \times M}$ and $\boldsymbol{E_p}, z_0 \in \mathbb{R}^{N \times M}$.

The encoder path of transformer (Fig. 1) has $L$ layers of Multihead Self-Attention (MSA) and Multi-Layer Perceptron (MLP) blocks:

$$z_l^{MSA} = MSA(norm(z_{l-1})) + z_{l-1} \tag{2}$$

$$z_l = z_l^{MLP} = MLP(norm(z_l^{MSA})) + z_l^{MSA} \tag{3}$$

where $z_l^{MSA}$ and $z_l^{MLP}$ are the outputs from MSA and MLP blocks, $norm(\cdot)$ denotes the layer normalization and $l$ is the layer index. The MLP blocks contain two linear layers followed by the GELU activation functions.

There are $n$ Self-Attention heads (SAs) in the MSA block to extract global information from the embedded sequence:

$$SA(z_i) = softmax(\frac{qk^T}{\sqrt{M_n}})v \tag{4}$$

where $z_i \in \mathbb{R}^{N \times M}$, $q$, $k$ and $v$ are the query, key and value of $z_i$ respectively, $q = W_q z_i$, $k = W_k z_i$ and $v = W_v z_i$. $W$s are the three weight matrices and $\sqrt{M_n} = \frac{M}{n}$ is the scaling factor. The output of the $softmax(\cdot)$ function is the similarity weight between $q$ and $k$. Then the MSA is defined as:

$$MSA(z) = [SA_1(z); SA_2(z); ...; SA_n(z)]W_{msa} \tag{5}$$

where $W_{msa}$ are the trainable weights.

Inspired by the UNETR [14], we use the same strategy for visualizing the multi-scale features $z_3$, $z_6$, $z_9$, $z_{12}$ from transformer encoder path. The feature size of each $z_i$ is $N \times M = \frac{H}{P} \times \frac{W}{P} \times \frac{D}{P} \times M$. We upsample the $z_i$ into the same size as the corresponding outputs from CNN-based encoder by deconvolution (deconv) operations, batch normalization (BN) and RELU activation function. The details can be viewed in Fig. 1. For the last level of the transformer features ($z_{12}$), we directly apply a residual convolution (res-conv) [17], BN and RELU for reshaping.

## 2.3. Convolutional neural network architecture

Our CNN (Fig. 1) is adapted from the 3D U-Net [1]. There are two parts of the CNN model, the encoder and the decoder. Four max-pooling and deconvolution operations are used for down-sampling and upsampling respectively. The feature maps from the top level are directly forwarded to the decoder, and the rest of the feature maps from the lower levels are fused with the encoded information from transformer path by addition. Then, the fused information is delivered to the decoder by skip connections that follows the way of 3D U-Net.

**Table 1**. Mean Dice scores in BTCV dataset. Bold numbers denote the highest Dice scores. The results of TransUNet are directly copied from [13]. The experiments of UNETR and proposed method use the public pipeline of UNETR. The organs from left to right are: spleen, right and left kidney, gallbladder, esophagus, liver, stomach, aorta, inferior vena cava, portal vein and splenic vein, pancreas, right and left adrenal gland, and overall average.

| Method | Spl | RKid | LKid | Gall | Eso | Liv | Sto | Aor | IVC | Veins | Pan | RAG | LAG | Avg. |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| TransUNet [13] | 85.1 | 77.0 | 81.9 | 63.1 | - | 94.1 | 75.6 | 87.2 | - | - | 55.9 | - | - | 77.5 |
| UNETR [14] | 93.4 | 85.5 | 87.6 | 61.9 | 74.7 | 95.7 | 76.8 | 85.2 | 77.2 | 69.8 | 61.5 | 64.4 | 59.4 | 76.9 |
| CATS | **95.8** | **90.2** | **93.4** | **65.9** | **77.1** | **96.8** | **83.0** | **88.6** | **83.1** | **76.9** | **73.8** | **70.2** | **62.6** | **81.4** |

## 3. EXPERIMENTS AND RESULTS

### 3.1. Datasets and Implementation Details

We used three public datasets for our experiments, and followed the evaluation metrics of each challenge.

**Beyond the Cranial Vault (BTCV)**[1] dataset contains 30/20 subjects with abdominal CT images for training/testing, with 13 different organs labeled by experts. To preprocess, we resampled the images and clipped HU values to range [-175, 250]. Three data augmentations were used: random flip, rotation and intensity shift. UNETR [14] is the winner for this challenge, and we compared our results to the publicly available UNETR implementation [2]. We also compare to the TransUNet model [13]. We note that both of these methods have been shown to be superior to CNN-only models such as 3D U-Net for this challenge [13, 14].

**Cross-Modality Domain Adaptation for Medical Image Segmentation (MoDA)**[3] has 105 contrast-enhanced T1-weighted MRIs with manual labels for vestibular schwannomas (VS). We split the dataset into 55/20/30 for training/validation/ testing. The preprocessing pipeline consists of rigid registration to MNI space and cropping. Adam optimizer and Dice loss are used. We again compare our results to TransUNet[13] and UNETR [14], as well as to a 2.5D CNN [4] which was provided as the baseline for the challenge. We also report the 95% Hausdorff distance in this dataset in addition to the metrics used in the challenge.

**Task 5 of Decathlon (Decathlon-5)**[4] consists of 32 MRIs with manual prostate labels. 2 MRIs in validation were excluded due to the wrong labels being provided in the public dataset. We compare to the TransFuse [16] model and the nn-Unet [3], which was the top-performing approach of the challenge at the time of the initial competition for this task. For a direct comparison with these two models, we use this
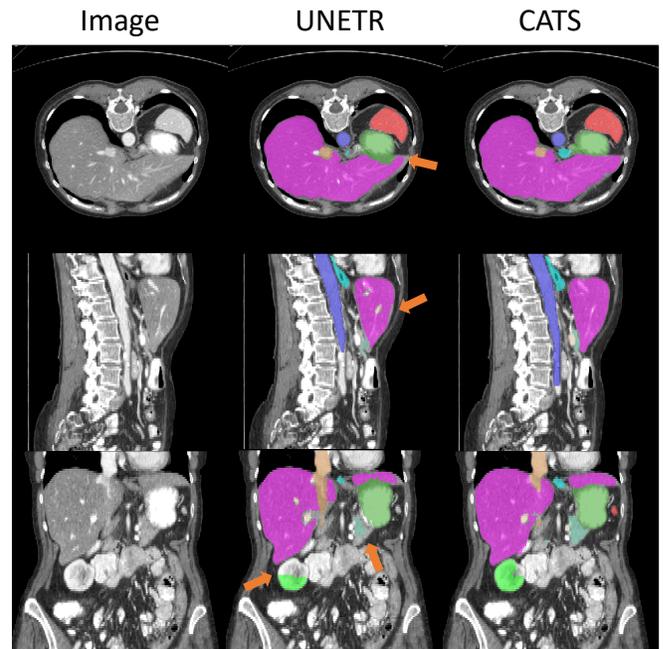


Image UNETR CATS

**Fig. 2**. Qualitative results in BTCV test set. Some major differences are highlighted by orange arrows.

dataset in a 5-fold cross-validation framework, and follow the setting in [3].

**Implementation details.** The training batch size was 2 for all three experiments, and constant learning rate was 0.0001. All intensities were normalized to range [0, 1]. We used Pytorch, MONAI and an Nvidia Titan RTX GPU.

### 3.2. BTCV Results

The Dice score is used to evaluate the BTCV experiment; the results can be viewed in Tab. 1. Our proposed method outperformed the state-of-the-art transformer-based models for each organ in this dataset. The most dramatic improvements (Dice improvement > 5%) between UNETR and our proposed method are in left kidney, stomach, inferior vena cava, portal vein and splenic vein, pancreas and right adrenal gland.
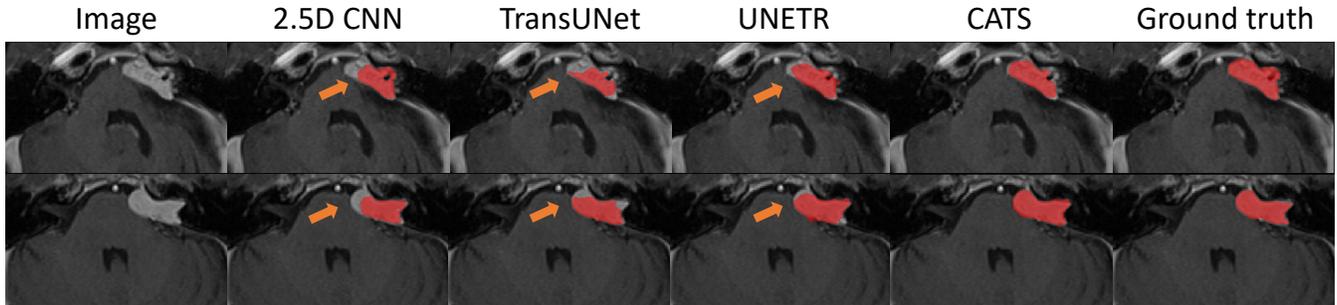
---

[1] https://www.synapse.org/#!Synapse:syn3193805/wiki/217753
[2] https://github.com/Project-MONAI/tutorials/blob/master/3d_segmentation/unetr_btcv_segmentation_3d.ipynb
[3] https://crossmoda.grand-challenge.org/
[4] http://medicaldecathlon.com/

**Fig. 3**. Quantitative results in MoDA. Local segmentation errors are highlighted with arrows.

It is noteworthy that our method improves the segmentation accuracy not only on larger organs such as stomach and kidney, but also on small organs such as the right adrenal gland.

Fig. 2 shows qualitative results. Compared to the UNETR, our proposed model produces smoother segmentation for the stomach and liver (axial view, arrow). We can also see (coronal view, arrows) that UNETR undersegments the right kidney and liver, unlike our proposed model.

### 3.3. MoDA results

The quantitative results of MoDA dataset are shown in Tab. 2. We report the Dice score, average surface distance (ASD) and 95-percent Hausdorff distance (HD95) as metrics. It is easy to observe that the CNN-only network has better performance than the transformer-only encoders for this task. However, our proposed CATS model outperformed the 2.5D CNN [4], which was specifically designed for segmenting VS from MRIs with large difference between in-plane resolution and slice thickness, as is the case for this dataset.

Fig. 3 shows the qualitative results of VS segmentation. While all methods appear to undersegment the VS, our proposed model most closely resembles the ground truth segmentations.

**Table 2**. Quantitative results in MoDA dataset, presented as $mean(std.dev.)$. Bold numbers indicate the best performance.

| Method | Dice | ASD | HD95 |
|---|---|---|---|
| 2.5D CNN [4] | 0.856 (1.000) | 0.69 (1.20) | 3.5 (5.2) |
| TransUNet [13] | 0.792 (0.234) | 7.86 (27.6) | 12 (31) |
| UNETR [14] | 0.772 (0.139) | 7.95 (14.2) | 26 (43) |
| CATS | **0.873 (0.088)** | **0.48 (0.63)** | **2.6 (3.6)** |

**Table 3**. Mean Dice scores in Decathlon-5 dataset. PZ and TZ denote the peripheral zone and the transition zone of the prostate, respectively.

| Method | PZ | TZ | Avg. |
|---|---|---|---|
| 2D nnUnet [3] | 0.6285 | 0.8380 | 0.7333 |
| 3D nnUnet full [3] | 0.6663 | 0.8410 | 0.7537 |
| TransFuse-S [16] | 0.6738 | 0.8539 | 0.7639 |
| CATS | **0.7136** | **0.8618** | **0.7877** |

### 3.4. Decathlon-5 results

We compare the nnUnet and TransFuse for the prostate segmentation in Tab. 3. Our proposed method has the highest Dice scores on all labels. Moreover, we improved the peripheral zone (PZ) nearly $4\%$ compared to the performance of the TransFuse model.

### 4. DISCUSSION AND CONCLUSIONS

In this paper, we propose a convolutional neural network with a transformer as an independent encoder. The transformer can complement the CNN by modeling long-range dependencies and capturing low-level details. We evaluate our proposed method on three public datasets: (1) BTCV, (2) MoDA and (3) Decathlon. Compared to the state-of-the-art models which also attempt to incorporate transformers into the segmentation networks in various ways, our proposed model has superior performance on each task. We believe this is due to our efficient integration of the transformer and CNN encoders, as well as our use of a 3D architecture compared to 2D models. In future work, we will apply our method to larger public datasets. Additionally, others transformer layers may be helpful to further improve the performance.

### 5. COMPLIANCE WITH ETHICAL STANDARDS

This research study was conducted retrospectively using human subject data made available in open access by BTCV,

MoDA and Decathlon (see detailed information in Sec. 3.1). Ethical approval was not required as confirmed by the license attached with the open access data.

## 6. ACKNOWLEDGMENTS

## 7. REFERENCES

[1] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger, "3d u-net: learning dense volumetric segmentation from sparse annotation," in *International conference on medical image computing and computer-assisted intervention*. Springer, 2016, pp. 424–432.

[2] Jo Schlemper, Ozan Oktay, Michiel Schaap, Mattias Heinrich, Bernhard Kainz, Ben Glocker, and Daniel Rueckert, "Attention gated networks: Learning to leverage salient regions in medical images," *Medical image analysis*, vol. 53, pp. 197–207, 2019.

[3] Fabian Isensee, Paul F Jäger, Simon AA Kohl, Jens Petersen, and Klaus H Maier-Hein, "Automated design of deep learning methods for biomedical image segmentation," *arXiv preprint arXiv:1904.08128*, 2019.

[4] Jonathan Shapey, Guotai Wang, Reuben Dorent, Alexis Dimitriadis, Wenqi Li, Ian Paddick, Neil Kitchen, Sotirios Bisdas, Shakeel R Saeed, Sebastien Ourselin, et al., "An artificial intelligence framework for automatic segmentation and volumetry of vestibular schwannomas from contrast-enhanced t1-weighted and high-resolution t2-weighted mri," *Journal of neurosurgery*, vol. 134, no. 1, pp. 171–179, 2019.

[5] Dewei Hu, Can Cui, Hao Li, Kathleen E Larson, Yuankai K Tao, and Ipek Oguz, "Life: A generalizable autodidactic pipeline for 3d oct-a vessel segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 514–524.

[6] Hao Li, Huahong Zhang, Hans Johnson, Jeffrey D Long, Jane S Paulsen, and Ipek Oguz, "Mri subcortical segmentation in neurodegeneration with cascaded 3d cnns," in *Medical Imaging 2021: Image Processing*. International Society for Optics and Photonics, 2021, vol. 11596, p. 115960W.

[7] Huahong Zhang, Hao Li, and Ipek Oguz, "Segmentation of new ms lesions with tiramisu and 2.5 d stacked slices," *MSSEG-2 challenge proceedings: Multiple sclerosis new lesions segmentation challenge using a data management and processing infrastructure*, p. 61, 2021.

[8] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al., "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[9] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo, "Swin transformer: Hierarchical vision transformer using shifted windows," *arXiv preprint arXiv:2103.14030*, 2021.

[10] Yutong Xie, Jianpeng Zhang, Chunhua Shen, and Yong Xia, "Cotr: Efficiently bridging cnn and transformer for 3d medical image segmentation," in *International conference on medical image computing and computer-assisted intervention*. Springer, 2021, pp. 171–180.

[11] Jeya Maria Jose Valanarasu, Poojan Oza, Ilker Hacihaliloglu, and Vishal M Patel, "Medical transformer: Gated axial-attention for medical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 2021, pp. 36–46.

[12] Fahad Shamshad, Salman Khan, Syed Waqas Zamir, Muhammad Haris Khan, Munawar Hayat, Fahad Shahbaz Khan, and Huazhu Fu, "Transformers in medical imaging: A survey," *arXiv preprint arXiv:2201.09873*, 2022.

[13] Jieneng Chen, Yongyi Lu, Qihang Yu, Xiangde Luo, Ehsan Adeli, Yan Wang, Le Lu, Alan L Yuille, and Yuyin Zhou, "Transunet: Transformers make strong encoders for medical image segmentation," *arXiv preprint arXiv:2102.04306*, 2021.

[14] Ali Hatamizadeh, Dong Yang, Holger Roth, and Daguang Xu, "Unetr: Transformers for 3d medical image segmentation," *arXiv preprint arXiv:2103.10504*, 2021.

[15] Hu Cao, Yueyue Wang, Joy Chen, Dongsheng Jiang, Xiaopeng Zhang, Qi Tian, and Manning Wang, "Swin-unet: Unet-like pure transformer for medical image segmentation," *arXiv preprint arXiv:2105.05537*, 2021.

[16] Yundong Zhang, Huiye Liu, and Qiang Hu, "Transfuse: Fusing transformers and cnns for medical image segmentation," *arXiv preprint arXiv:2102.08005*, 2021.

[17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.