

DISCOHEAD: AUDIO-AND-VIDEO-DRIVEN TALKING HEAD GENERATION BY DISENTANGLED CONTROL OF HEAD POSE AND FACIAL EXPRESSIONS

Geumbyeol Hwang*, Sunwon Hong*, Seunghyun Lee, Sungwoo Park, Gyeongsu Chae†

DeepBrain AI Inc., Seoul, Korea

ABSTRACT

For realistic talking head generation, creating natural head motion while maintaining accurate lip synchronization is essential. To fulfill this challenging task, we propose DisCoHead, a novel method to disentangle and control head pose and facial expressions without supervision. DisCoHead uses a single geometric transformation as a bottleneck to isolate and extract head motion from a head-driving video. Either an affine or a thin-plate spline transformation can be used and both work well as geometric bottlenecks. We enhance the efficiency of DisCoHead by integrating a dense motion estimator and the encoder of a generator which are originally separate modules. Taking a step further, we also propose a neural mix approach where dense motion is estimated and applied implicitly by the encoder. After applying the disentangled head motion to a source identity, DisCoHead controls the mouth region according to speech audio, and it blinks eyes and moves eyebrows following a separate driving video of the eye region, via the weight modulation of convolutional neural networks. The experiments using multiple datasets show that DisCoHead successfully generates realistic audio-and-video-driven talking heads and outperforms state-of-the-art methods. Project page: <https://deepbrainai-research.github.io/discohead/>

Index Terms— Talking head generation, audio-and-video-driven, disentanglement, head pose, facial expressions

1. INTRODUCTION

Talking head generation is a method to synthesize facial video according to speech audio. It has broad applications like virtual videotelephony, automated video production, and character animation to name a few. Although deep learning technology has accelerated the advancement of talking head generation, creating natural head motion while maintaining accurate lip synchronization is still a challenge.

A way of talking face synthesis is inpainting the mouth parts in existing videos using audio information [1, 2, 3, 4].

While head motion itself is natural in this approach, it is hard to manipulate the motion in existing videos, particularly to match the speech content.

Warping a facial frame based on a sequence of audio features is another way [5, 6]. Speech and head motion are better aligned in this case, but the generated motions and nonspeech facial expressions lack natural dynamics. Following studies improve the audio-driven head motion and facial expressions by introducing sequence discriminators [7, 8] and a noise generator [7]. [9] adds diversity to the motion and expressions using a variational autoencoder framework. Other researches create head motion from audio via explicit structural representation of the face [10, 11, 12, 13].

Few audio-driven talking face models control head pose across identities independent of lip synchronization. [14] modularizes audio-visual representation into identity, head pose, and speech content using elaborate image augmentation and contrastive learning without structural intermediate information. Whereas [15, 16] use an external detector to extract shape-independent information of head pose and eye blinks, then fuse them with audio signals to reenact faces.

Meanwhile, in the field of image animation, a series of studies successfully decouple appearance and motion using unsupervised geometric-transformation-based approaches [17, 18, 19]. In these approaches, the parameters of multiple affine or thin-plate spline (TPS) transformations between a source image and driving video frames are extracted, and dense motion (i.e., pixel-wise optical flow) is estimated based on the transformation parameters. Then, the source image is warped and inpainted to generate output frames.

Inspired by the geometric-transformation-based image animation, we propose DisCoHead, a novel method enabling disentangled control of head pose and facial expressions for audio-and-video-driven talking head generation.

Our novelties are threefold. First, we use a single geometric transformation instead of multiple ones. A geometric transformation computed from a source and a driving frames serves as a bottleneck to isolate head pose from facial expressions. We find both affine and TPS transformations work well as geometric bottlenecks. Second, we integrate a dense motion estimator and the encoder of a generator into one. Unlike the existing image animation models, we do not need to combine multiple geometric transformations to estimate

*Equal contribution. †Correspondence to: gc@deepbrain.io

This work was supported by the Institute of Information & Communications Technology Planning & Evaluation (IITP) grant funded by the Ministry of Science and ICT (MSIT) of South Korea (No. 2021-0-00888).

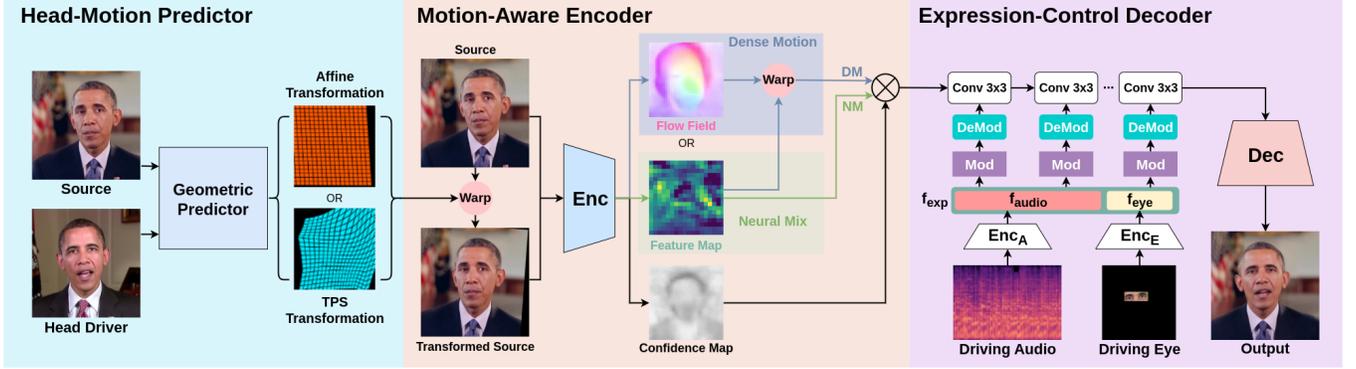


Fig. 1. Architecture of our model. Head-Motion Predictor estimates an affine or TPS transformation by separately forwarding source and driving frames. With the original and transformed sources, Motion-Aware Encoder extracts head-motion-applied source features. Expression-Control Decoder then manipulates facial expressions according to driving audio and eye features.

dense motion. With proper tweaks to input-output formulation, the encoder can effectively extract source features and dense motion simultaneously. Taking a step further to efficient architecture design, we also suggest a neural mix of the features from a source and a geometrically transformed source instead of explicit warping based on dense motion. Third, after warping or mixing the source features to apply head motion, we control facial expressions using speech audio and video frames of the eye region. We adopt the weight modulation of convolutional neural networks [20], where the audio and eye signals change the values of convolution filters to modify the facial expressions of the warped features.

Our proposed method allows us to separately control head pose, lip synchronization, and nonspeech facial expressions in an identity-agnostic manner. Our extensive experiments using multiple datasets show that DisCoHead successfully generates realistic audio-and-video-driven talking heads and outperforms state-of-the-art methods.

2. METHODS

The entire pipeline of DisCoHead is depicted in Fig. 1. We first introduce the head-motion predictor (Sec. 2.1). Then, we explain the motion-aware encoder (Sec. 2.2) and the expression-control decoder (Sec. 2.3) of the generator.

2.1. Head-Motion Predictor

The geometric transformation predictor representing head motion has two variants: affine predictor and TPS predictor.

We employ the PCA-based affine predictor [18]. First, a single-channel heatmap H is predicted from the input image X . The expectation of pixel locations Z with probabilities H can be interpreted as the translation component t of an affine transformation:

$$t = \sum_{z \in Z} H(z)z. \quad (1)$$

The rotation and scaling components are computed with the principal axes of H , which can be extracted from a covariance matrix by the singular value decomposition (SVD):

$$U\Sigma V = \sum_{z \in Z} H(z)(z-t)(z-t)^T. \quad (2)$$

Here, U and V are unitary matrices, and Σ is the diagonal matrix of singular values. Then we can configure the affine transformation as $A_{X \leftarrow R} = [U\Sigma^{1/2}, t]$, where R is an abstract reference frame. The affine transformation for the motion from the head driver D to the source S is computed as:

$$A_{S \leftarrow D} = A_{S \leftarrow R} A_{D \leftarrow R}^{-1}. \quad (3)$$

For the TPS predictor [19], N keypoints are predicted from the input frame. Given the N pairs of keypoints from the source and the driving frames, we can get the corresponding TPS transformation as follows:

$$T(p) = A \begin{bmatrix} p \\ 1 \end{bmatrix} + \sum_{i=1}^N \omega_i \phi(\|P_i^D - p\|_2), \quad (4)$$

where $p = (x, y)^T$ is pixel coordinates, P_i^D is the i th keypoint of the driving frame D , $A \in \mathbb{R}^{2 \times 3}$ and $\omega_i \in \mathbb{R}^{2 \times 1}$ are the TPS coefficients obtained by solving the minimum distortion equations in [21], and a radial basis function $\phi(r)$ represents the influence of each keypoint on the pixel at p :

$$\phi(r) = r^2 \log r^2. \quad (5)$$

2.2. Motion-Aware Encoder

We integrate a dense motion estimator and the encoder of a generator into a motion-aware encoder that combines the source appearance with the driver’s head pose and outputs a head-pose-aligned source feature to be used to generate the

output frame. Along with the dense motion estimation approach (dense motion variant), we also propose a more compact approach that implicitly applies dense motion by neural networks (neural mix variant).

First, we warp the source frame S with the predicted geometric transformation to make a transformed source S_T . The motion-aware encoder receives both S and S_T and produces a source feature F , a confidence map C , and optionally a motion mask M for the dense motion variant:

$$F, C, (M) = Enc(S, S_T). \quad (6)$$

For the dense motion variant, the motion mask M is used to produce a pixel-wise optical flow O_P with the weighted sum of the identity flow O_I and the coarse flow O_T computed from the predicted geometric transformation. Then we warp F with O_P to make the aligned source feature F_A :

$$O_P = (1 - M) \circ O_I + M \circ O_T, \quad (7)$$

$$F_A = Warp(F, O_P). \quad (8)$$

For the neural mix variant, we simply omit the warping based on the dense flow O_P supposing that neural networks are able to implicitly estimate and apply dense motion by mixing the features from the source and the transformed source:

$$F_A = F. \quad (9)$$

The final encoder output E is computed by applying the confidence map C to the aligned feature F_A to inform the decoder where and how much the local details should be painted:

$$E = C \circ F_A. \quad (10)$$

In Eq. (7) and Eq. (10), \circ denotes the Hadamard product.

2.3. Expression-Control Decoder

An audio spectrogram corresponding to a target frame is encoded into an audio feature f_{audio} , and a masked eye driving frame is encoded into an eye feature f_{eye} . We concatenate the two features to obtain a facial expression feature f_{exp} . Then we manipulate the expressions of the encoder output E by modulating the convolution kernel weights [20] with f_{exp} :

$$\omega'_{ijk} = \frac{s_i \cdot \omega_{ijk}}{\sqrt{\sum_{i,k} (s_i \cdot \omega_{ijk})^2 + \epsilon}}, \quad (11)$$

where ω and ω' are the original and modulated weights, s is the scale value predicted from f_{exp} , and i , j , and k indicate the input channel, output channel, and spatial position of the convolution kernel, respectively. ϵ is a small constant to prevent numerical issues.

After the modulated convolution blocks modify the facial expressions, DisCoHead generates the output frame with bilinear upsampling convolution blocks.

3. EXPERIMENTS

3.1. Datasets

We use the Obama dataset [1], the GRID dataset [22], and the Korean election broadcast addresses dataset (KoEBA). The Obama dataset contains the weekly presidential addresses of Barrack Obama. GRID is a set of video clips of 34 speakers pronouncing short utterances. KoEBA is a high-quality multi-speaker audio-video dataset composed of official broadcast addresses of Korean politicians.¹ The facial regions of the video frames are cropped and resized to 256×256 . The driving audio is 400 ms long and centered on each video frame. We split each dataset into training (80%) and test (20%) sets.

3.2. Implementation Details

We adopt the hourglass network [23] for the affine predictor and ResNet-18 [24] for the TPS predictor. The motion-aware encoder has three downsample convolution blocks, followed by three convolution layers with a kernel size of 1. The audio encoder consists of four 1D convolution layers, an LSTM layer, and two fully-connected layers. The eye encoder has five downsample blocks, followed by a global average pooling layer and two fully-connected layers. The expression-control decoder has six modulated residual convolution blocks, followed by three upsample convolution blocks and an output convolution layer with a kernel size of 7. We combine L1 loss and the perceptual loss [25] based on the pre-trained VGG-19 [26] and use Adam optimizer with a learning rate of $1e-4$ and a batch size of 16:

$$L_{total} = L_1 + \lambda \cdot L_{perceptual}. \quad (12)$$

3.3. Quantitative Results

Table 1 shows the reconstruction performance of DisCoHead on the three datasets. We use peak signal-to-noise ratio (PSNR), structural similarity index measure (SSIM), Fréchet inception distance (FID), and learned perceptual image patch similarity (LPIPS) to assess generation quality. Average key-point distance (AKD) measures the accuracy of pose and expressions using facial keypoints, and average Euclidean distance (AED) evaluates identity preservation based on a face recognition model. Comparisons with state-of-the-art baselines show that our model achieves the best results for all metrics and datasets. Both affine and TPS transformations work well as geometric bottlenecks to isolate and steer head motion. Also, the performance of the neural mix variant is almost on a par with the dense motion variant.

3.4. Qualitative Results

Fig. 2 demonstrates DisCoHead’s capability to disentangle head motion and facial expressions. The generated head pose,

¹<https://github.com/deepbrainai-research/koeba>

Dataset	Obama						GRID						KoEBA					
Method	PSNR \uparrow	SSIM \uparrow	FID \downarrow	LPIPS \downarrow	AKD \downarrow	AED \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	LPIPS \downarrow	AKD \downarrow	AED \downarrow	PSNR \uparrow	SSIM \uparrow	FID \downarrow	LPIPS \downarrow	AKD \downarrow	AED \downarrow
APB2FaceV2 [16]	17.84	0.494	60.38	0.596	3.746	0.266	28.12	0.662	62.24	0.261	2.952	0.116	22.04	0.706	58.59	0.256	2.972	0.244
Wav2Lip [2]	22.78	0.806	50.76	0.145	1.711	0.075	29.02	0.903	70.77	0.152	1.155	0.024	25.61	0.876	32.85	0.086	1.834	0.080
PC-AVS [14]	22.14	0.491	8.493	0.227	2.761	0.195	25.31	0.581	21.23	0.213	2.357	0.102	23.36	0.556	37.45	0.101	2.174	0.128
DisCoHead-Affine-DM	28.39	0.904	0.618	0.051	0.915	0.031	34.04	0.926	0.340	0.056	1.047	0.015	28.87	0.906	0.857	0.054	1.061	0.023
DisCoHead-Affine-NM	27.07	0.883	0.736	0.052	1.002	0.037	34.14	0.928	0.406	0.076	1.007	0.012	28.18	0.900	4.602	0.066	1.434	0.043
DisCoHead-TPS-DM	26.33	0.878	1.111	0.075	0.900	0.041	33.80	0.923	0.781	0.048	1.060	0.015	28.66	0.899	0.746	0.057	1.297	0.033
DisCoHead-TPS-NM	26.35	0.879	0.866	0.074	1.397	0.062	33.81	0.924	0.849	0.046	0.983	0.012	29.25	0.903	0.758	0.059	0.987	0.021

Table 1. Quantitative results on the three datasets (DM: dense motion, NM: neural mix).

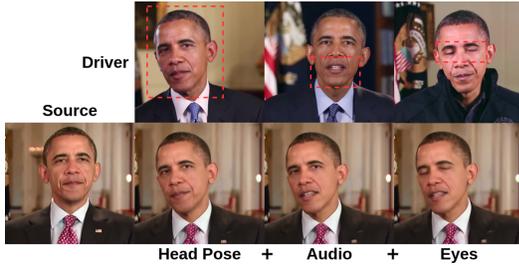


Fig. 2. Disentangled control of head pose, lip movements, and eye expressions on the Obama dataset.

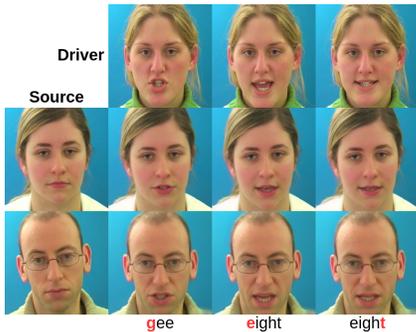


Fig. 3. Qualitative results on the GRID dataset.



Fig. 4. Qualitative results on the KoEBA dataset.

articulatory expression, and eye blink conform to the successively applied head pose, speech audio (corresponding to the shown driving frame), and eye blink of the excerpts from three different video clips without interfering with each other.

Fig. 3 and Fig. 4 present our model’s ability to maintain the source identity on the GRID and KoEBA datasets. In Fig. 3, the generated faces well maintain their own shapes, and the lip expressions pronouncing the same characters (marked in red) are correct but differ for each identity. In addition to precise motion control and identity preservation, subtle emotional expressions (i.e., slight frown) on the eye region of the driver are also transferred to the generated frames in Fig. 4.

4. DISCUSSION

DisCoHead’s geometric bottleneck successfully disentangles identity, head pose, and facial expressions, but modeling extreme poses can be difficult. Its input audio and images of the eye region provide rich information to modulate expressions.

On the other hand, APB2FaceV2 [16] acquires head pose and eye blink information from an external facial landmark detector. Therefore, nonspeech facial expressions other than eye blinks cannot be delivered. Wav2Lip [2] is unable to separate identity, head pose, and nonspeech expressions because it restores the lower half of a face according to the upper half and speech audio. Moreover, its output resolution is 96×96 , much smaller than 256×256 of DisCoHead. PC-AVS [14] can handle extreme poses but only decomposes identity, head pose, and speech content. As a consequence, it lacks the ability to control nonspeech facial expressions.

5. CONCLUSION

We design a novel method to disentangle and separately control head pose and facial expressions for audio-and-video-driven talking head generation. DisCoHead uses a single geometric transformation as a medium to represent head motion and the weight modulation of convolutional layers to manipulate speech and nonspeech facial expressions. We enhance the efficiency of DisCoHead by fusing dense motion estimation and video frame generation which are originally formulated as a sequential process. In our experiments using multiple datasets, DisCoHead excels state-of-the-art methods.

6. REFERENCES

- [1] Supasorn Suwajanakorn, Steven M. Seitz, and Ira Kemelmacher-Shlizerman, “Synthesizing obama: Learning lip sync from audio,” *ACM Trans. Graph.*, vol. 36, no. 4, jul 2017.
- [2] K R Prajwal, Rudrabha Mukhopadhyay, Vinay P. Namboodiri, and C.V. Jawahar, “A lip sync expert is all you need for speech to lip generation in the wild,” in *Proceedings of the 28th ACM International Conference on Multimedia*, 2020, p. 484–492.
- [3] Justus Thies, Mohamed Elgharib, Ayush Tewari, Christian Theobalt, and Matthias Nießner, “Neural voice puppetry: Audio-driven facial reenactment,” in *Computer Vision – ECCV 2020*, 2020, pp. 716–731.
- [4] Xin Wen, Miao Wang, Christian Richardt, Ze-Yin Chen, and Shi-Min Hu, “Photorealistic audio-driven video portraits,” *IEEE Transactions on Visualization and Computer Graphics*, vol. 26, no. 12, pp. 3457–3466, 2020.
- [5] Joon Son Chung, Amir Jamaludin, and Andrew Zisserman, “You said that?,” in *Proceedings of the British Machine Vision Conference (BMVC)*, September 2017, pp. 109.1–109.12.
- [6] Yang Song, Jingwen Zhu, Dawei Li, Andy Wang, and Hairong Qi, “Talking face generation by conditional recurrent adversarial network,” in *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, 7 2019, pp. 919–925.
- [7] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic, “End-to-end speech-driven realistic facial animation with temporal gans,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*, June 2019.
- [8] Konstantinos Vougioukas, Stavros Petridis, and Maja Pantic, “Realistic speech-driven facial animation with gans,” *Int. J. Comput. Vision*, vol. 128, no. 5, pp. 1398–1413, may 2020.
- [9] Ravindra Yadav, Ashish Sardana, Vinay P. Namboodiri, and Rajesh M. Hegde, “Stochastic talking face generation using latent distribution matching,” in *Proc. Interspeech 2020*, 2020, pp. 1311–1315.
- [10] Lele Chen, Guofeng Cui, Celong Liu, Zhong Li, Ziyi Kou, Yi Xu, and Chenliang Xu, “Talking-head generation with rhythmic head motion,” in *Computer Vision – ECCV 2020*, 2020, pp. 35–51.
- [11] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li, “Makelttalk: Speaker-aware talking-head animation,” *ACM Trans. Graph.*, vol. 39, no. 6, nov 2020.
- [12] Zhimeng Zhang, Lincheng Li, Yu Ding, and Changjie Fan, “Flow-guided one-shot talking face generation with a high-resolution audio-visual dataset,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 3661–3670.
- [13] Suzhen Wang, Lincheng Li, Yu Ding, Changjie Fan, and Xin Yu, “Audio2head: Audio-driven one-shot talking-head generation with natural head motion,” in *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, 8 2021, pp. 1098–1105.
- [14] Hang Zhou, Yasheng Sun, Wayne Wu, Chen Change Loy, Xiaogang Wang, and Ziwei Liu, “Pose-controllable talking face generation by implicitly modularized audio-visual representation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 4176–4186.
- [15] Jiangning Zhang, Liang Liu, Zhucun Xue, and Yong Liu, “Apb2face: Audio-guided face reenactment with auxiliary pose and blink signals,” in *ICASSP 2020 - 2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2020, pp. 4402–4406.
- [16] Jiangning Zhang, Xianfang Zeng, Chao Xu, and Yong Liu, “Real-time audio-guided multi-face reenactment,” *IEEE Signal Processing Letters*, vol. 29, pp. 1–5, 2022.
- [17] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe, “First order motion model for image animation,” in *Advances in Neural Information Processing Systems*, 2019, vol. 32.
- [18] Aliaksandr Siarohin, Oliver J. Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov, “Motion representations for articulated animation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2021, pp. 13653–13662.
- [19] Jian Zhao and Hui Zhang, “Thin-plate spline motion model for image animation,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2022, pp. 3657–3666.
- [20] Tero Karras, Samuli Laine, Miika Aittala, Janne Hellsten, Jaakko Lehtinen, and Timo Aila, “Analyzing and improving the image quality of stylegan,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020.
- [21] F.L. Bookstein, “Principal warps: thin-plate splines and the decomposition of deformations,” *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 11, no. 6, pp. 567–585, 1989.
- [22] Martin Cooke, Jon Barker, Stuart Cunningham, and Xu Shao, “An audio-visual corpus for speech perception and automatic speech recognition,” *The Journal of the Acoustical Society of America*, vol. 120, pp. 2421–2424, November 2006.
- [23] Alejandro Newell, Kaiyu Yang, and Jia Deng, “Stacked hour-glass networks for human pose estimation,” in *Computer Vision – ECCV 2016*, 2016, pp. 483–499.
- [24] 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 (CVPR)*, June 2016.
- [25] Justin Johnson, Alexandre Alahi, and Li Fei-Fei, “Perceptual losses for real-time style transfer and super-resolution,” in *Computer Vision – ECCV 2016*, 2016, pp. 694–711.
- [26] Karen Simonyan and Andrew Zisserman, “Very deep convolutional networks for large-scale image recognition,” in *International Conference on Learning Representations*, 2015.