# Controllable Attention for Structured Layered Video Decomposition

Jean-Baptiste Alayrac[1*]     João Carreira[1*]     Relja Arandjelović[1]     Andrew Zisserman[1,2]

{jalayrac,joaoluis}@google.com

[1]DeepMind     [2]VGG, Dept. of Engineering Science, University of Oxford

## Abstract

*The objective of this paper is to be able to separate a video into its natural layers, and to control which of the separated layers to attend to. For example, to be able to separate reflections, transparency or object motion.*

*We make the following three contributions: (i) we introduce a new structured neural network architecture that explicitly incorporates layers (as spatial masks) into its design. This improves separation performance over previous general purpose networks for this task; (ii) we demonstrate that we can augment the architecture to leverage external cues such as audio for controllability and to help disambiguation; and (iii) we experimentally demonstrate the effectiveness of our approach and training procedure with controlled experiments while also showing that the proposed model can be successfully applied to real-word applications such as reflection removal and action recognition in cluttered scenes.*

## 1. Introduction

"The more you look the more you see", is generally true for our complex, ambiguous visual world. Consider the everyday task of cleaning teeth in front of a mirror. People performing this task may first attend to the mirror surface to identify any dirty spots, clean them up, then switch attention to their mouth reflected in the mirror. Or they may hear steps behind them and switch attention to a new face now reflecting in the mirror. Not all visual possibilities can be investigated at once given a fixed computational budget and this creates the need for such controllable attention mechanisms.

Layers offer a simple but useful model for handling this complexity of the visual world [51]. They provide a compositional model of an image or video sequence, and cover a multitude of scenarios (reflections, shadows, occlusions, haze, blur, ...) according to the composition rule. For example, an additive composition models reflections, and occlu-
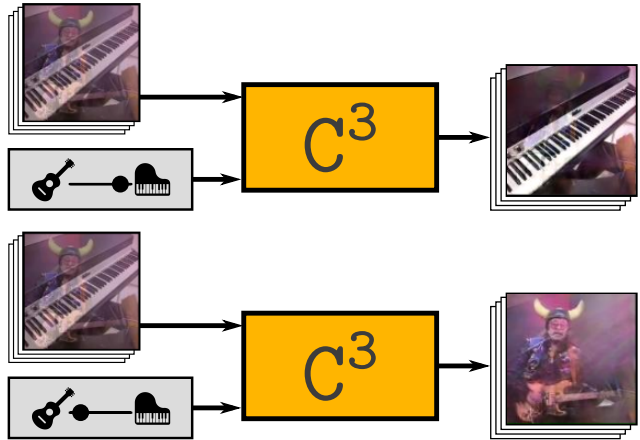


Figure 1: We propose a model, $C^3$, able to decompose a video into meaningful layers. This decomposition process is controllable through external cues such as *audio*, that can select the layer to output.

sion is modelled by superimposing opaque layers in a depth ordering. Given a a layered decomposition, attention can switch between the various layers as necessary for the task at hand.

Our objective in this paper is to separate videos into their constituent layers, and to *select* the layers to attend to as illustrated in Figure 1. A number of recent works have used deep learning to separate layers in images and videos [3, 12, 16, 18, 26, 58], with varying success, but the selection of the layers has either had to be hard coded into the architecture, or the layers are arbitrarily mapped to the outputs. For example, [3] considers the problem of separating blended videos into component videos, but because the mapping between input videos and outputs is arbitrary, training is forced to use a permutation invariant loss, and there is no control over the mapping at inference time. How can this symmetry between the composed input layers and output layers be broken?

The solution explored here is based on the simple fact that videos do not consist of visual streams alone, they also have an audio stream; and, significantly, the visual and au-

---

*Equal contribution.

dio streams are often correlated. The correlation can be strong (e.g. the synchronised sound and movement of beating on a drum), or quite weak (e.g. street noise that separates an outdoor from indoor scene), but this correlation can be employed to break the symmetry. This symmetry breaking is related to recent approaches to the cocktail party audio separation problem [2, 15] where visual cues are used to select speakers and improve the quality of the separation. Here we use audio cues to select the visual layers.

**Contributions:** The contributions of this paper are three-fold: **(i)** we propose a new structured neural network architecture that explicitly incorporates layers (as spatial masks) into its design; **(ii)** we demonstrate that we can augment the architecture to leverage external cues such as *audio* for controllability and to help disambiguation; and **(iii)** we experimentally demonstrate the effectiveness of our approach and training procedure with controlled experiments while also showing that the proposed model can be successfully applied to real-word applications such as *reflection removal* and *action recognition in cluttered scenes*.

We show that the new architecture leads to improved layer separation. This is demonstrated both qualitatively and quantitatively by comparing to recent general purpose models, such as the visual centrifuge [3]. For the quantitative evaluation we evaluate how the downstream task of human action recognition is affected by reflection removal. For this, we compare the performance of a standard action classification network on sequences with reflections, and with reflections removed using the layer architecture, and demonstrate a significant improvement in the latter case.

## 2. Related work

**Attention control.** Attention in neural network modelling has had a significant impact in natural language processing, such as machine translation, [5, 49] and vision [54], where it is implemented as a soft masking of features. In these settings attention is often not directly evaluated, but is just used as an aid to improve the end performance. In this paper we investigate models of attention in isolation, aiming for high consistency and controllability. By consistency we mean the ability to maintain the focus of attention on a particular target. By controllability we mean the ability to switch to a different target on command.

Visual attentional control is actively studied in psychology and neuroscience [14, 20, 28, 36, 48, 57] and, when malfunctioning, is a potentially important cause of conditions such as ADHD, autism or schizophrenia [33]. One of the problems studied in these fields is the relationship between attention control based on top-down processes that are voluntary and goal-directed, and bottom-up processes that are stimulus-driven (e.g. saliency) [27, 48]. Another in-

teresting aspect is the types of representations that are subject to attention, often categorized into location-based [42], object-based or feature-based [6]: examples of the latter include attending to anything that is red, or to anything that moves. Another relevant stream of research relates to the role of attention in multisensory integration [45, 47]. Note also that attention does not always require eye movement – this is called *covert* (as opposed to *overt*) attention. In this paper we consider covert attention as we will not be considering active vision approaches, and focus on feature-based visual attention control.

**Cross-modal attention control.** The idea of using one modality to control attention in the other has a long history, one notable application being informed audio source separation and denoising [7, 21, 39, 52]. Visual information has been used to aid audio denoising [21, 39], solve the cocktail party problem of isolating sound coming from different speakers [2, 15, 37, 52] or musical instruments [7, 19, 59]. Other sources of information used for audio source separation include text to separate speech [32] and score to separate musical instruments [25].

More relevant to this paper where audio is used for control, [4, 37, 40, 59] learn to attend to the object that is making the sound. However, unlike in this work, they do not directly output the disentangled video nor can they be used to remove reflections as objects are assumed to be perfectly opaque.

Other examples of control across modalities include temporally localizing a moment in a video using language [24], video summarization guided by titles [44] or query object labels [41], object localization from spoken words [23], image-text alignment [29], and interactive object segmentation via user clicks [9].

**Layered video representations.** Layered image and video representations have a long history in computer vision [50] and are an appealing framework for modelling 2.1D depth relationships [43, 56], motion segmentation [50], reflections [8, 12, 16, 17, 22, 26, 31, 35, 46, 55, 58], transparency [3, 18], or even haze [18]. There is also evidence that the brain uses multi-layered visual representations for modelling transparency and occlusion [53].

## 3. Approach

This section describes the two technical contributions of this work. First, in Section 3.1, a novel architecture for decomposing videos into layers. This architecture is built upon the visual centrifuge [3], a generic U-Net like encoder-decoder, but extends it with two structural changes tailored towards the layered video decomposition task. Second, in Section 3.2, the decomposition model is endowed with controllability – the ability of the network to use external cues
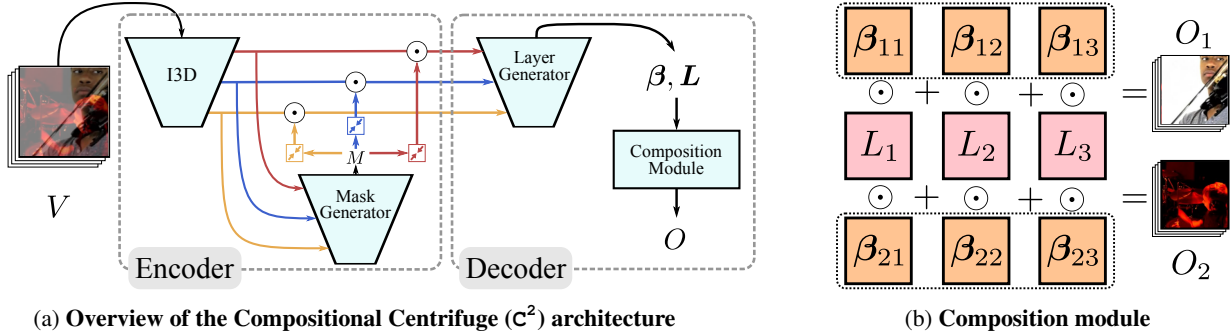
(a) **Overview of the Compositional Centrifuge ($\mathtt{C^2}$) architecture**

(b) **Composition module**

Figure 2: **Network architecture for layer decomposition** (3.1).

to control what it should focus on reconstructing. Here, we propose to use a natural video modality, namely *audio*, to select layers. Given this external cue, different mechanisms for controlling the outputs are investigated. Finally, in Section 3.3, we describe how this model can be trained for successful controllable video decomposition.

In the following, $V$ stands for an input video. Formally, $V \in \mathbb{R}^{T \times W \times H \times 3}$ where $T$ is the number of frames, $W$ and $H$ are the width and height of the frames, and there are 3 standard RGB channels. The network produces an $T \times W \times H \times (n \times 3)$ tensor, interpreted as $n$ output videos $O$, where each $O_i$ is of the same size as $V$.

### 3.1. Architecture for layer decomposition

We start from the visual centrifuge [3], a U-Net [38] encoder-decoder architecture, which separates an input video into $n$ output videos. The encoder consists of an I3D network [11] and the decoder is composed by stacking 3D up convolutions. However, the U-Net architecture used there is generic and not tailored to the layered video decomposition task (this is verified experimentally in Section 4.1). Therefore, we propose two structural modifications specifically designed to achieve layered decomposition, forming a new network architecture, *Compositional Centrifuge ($C^2$)*, shown in Figure 2a. Firstly, a bespoke gating mechanism is used in the encoder, which enables selection of scene segments across space/time, thereby making the decoder's task easier. Secondly, layer compositionality is imposed by constraining how the output videos are generated – the *layer generator* outputs multiple layers $L$ and their composing coefficients $\beta$ such that the output videos $O$ are produced as a linear combination of the layers. These modifications are described in detail next.

**Encoder.** We aim to recover layers in the presence of occlusions and transparent surfaces. In such cases there are windows of opportunity when objects are fully visible and their appearance can be modelled, and periods when the objects are temporarily invisible or indistinguishable and hence can only be tracked. We incorporate this intuition into a novel spatio-temporal encoder architecture. The core idea is that

the features produced by the I3D are gated with multiple ($m$) masks, also produced by the encoder itself. The gated features therefore already encode information about the underlying layers and this helps the decoder's task.

In order to avoid gating all features with all $m$ masks, which would be prohibitively expensive in terms of computation and memory usage, feature channels are divided into $m$ mutually-exclusive groups and each mask is applied only to the corresponding group.

More formally, the mask generator produces $M \in [0, 1]^{T \times W \times H \times m}$ which is interpreted as a set of $m$ spatio-temporal masks $M = (M^c)_{c=1}^m$. $M$ is constrained to sum to 1 along the channel dimension by using a softmax non-linearity. Denote $F_l$ the output feature taken at level $l$ in the I3D. We assume that $F_l \in \mathbb{R}^{T_l \times W_l \times H_l \times (m \times d_l)}$, *i.e.* the number of output channels of $F_l$ is a multiple of $m$. Given this, $F_l$ can be grouped into $m$ features $(F_l^c)_{c=1}^m$ where $F_l^c \in \mathbb{R}^{T_l \times W_l \times H_l \times d_l}$. The following transformation is applied to each $F_l^c$:

$$\tilde{F}_l^c = M_l^c \odot F_l^c, \tag{1}$$

where $M_l^c$ is obtained by downsampling $M^c$ to the shape $[T_l \times W_l \times H_l]$, $\odot$ refers to the Hadamard matrix product with a slight abuse of notation as the channel dimension is broadcast, *i.e.* the same mask is used across the channels. This process is illustrated in Figure 2a. Appendix B gives details on which feature levels are used in practice.

**Imposing compositionality.** In order to bias the decoder towards constructing layered decompositions, we split it into two parts – the *layer generator* produces $m$ layers $L$ and composing coefficients $\beta$ which are then combined by the *composition module* to form the final $n$ output videos $O$. The motivation is that individual layers should ideally represent independent scene units, such as moving objects, reflections or shadows, that can be composed in different ways into full scene videos. The proposed model architecture is designed to impose the inductive bias towards this type of compositionality.

More formally, the *layer generator* outputs a set of $m$

*layers* $\boldsymbol{L} = (L_j)_{j=1}^m$, where $L_j \in \mathbb{R}^{T \times H \times W \times 3}$, and a set of $n \times m$ *composing coefficients* $\boldsymbol{\beta} = (\boldsymbol{\beta}_{ij})_{(i,j) \in [\![1,n]\!] \times [\![1,m]\!]}$. These are then combined in the *composition module* (Figure 2b) to produce the final output videos $\boldsymbol{O}$:

$$O_i = \sum_j \boldsymbol{\beta}_{ij} \odot L_j. \tag{2}$$

## 3.2. Controllable symmetry breaking

The method presented in the previous section is inherently symmetric – the network is free to assign videos to output slots in any order. In this section, we present a strategy for controllable attention that is able to break the symmetry by making use of side-information, a *control signal*, provided as an additional input to the network. Audio is used as a natural control signal since it is readily available with the video. In our mirror example from the introduction, hearing speech indicates the attention should be focused on the person in the mirror, not the mirror surface itself. For the rest of this section, audio is used as the *control signal*, but the proposed approach remains agnostic to the control signal nature.

Next, we explain how to compute audio features, fuse them with the visual features, and finally, how to obtain the output video which corresponds to the input audio. The architecture, named *Controllable Compositional Centrifuge* ($C^3$), is shown in Figure 3.

**Audio network.** The audio first needs to be processed before feeding it as a control signal to the video decomposition model. We follow the strategy employed in [4] to process the audio. Namely, the log spectrogram of the raw audio signal is computed and treated as an image, and a VGG-like network is used to extract the audio features. The network is trained from scratch jointly with the video decomposition model.

**Audio-visual fusion.** To feed the audio signal to the video model, we concatenate audio features to the outputs of the encoder before they get passed to the decoder. Since visual and audio features have different shapes – their sampling rates differ and they are 3-D and 4-D tensors for audio and vision, respectively – they cannot be concatenated naively. We make the two features compatible by **(1)** average pooling the audio features over frequency dimension, **(2)** sampling audio features in time to match the number of temporal video feature samples, and **(3)** broadcasting the audio feature in the spatial dimensions. After these operations the audio tensor is concatenated with the visual tensor along the channel dimension. This fusion process is illustrated in Figure 3. We provide the full details of this architecture in Appendix B.

**Attention control.** We propose two strategies for obtaining the output video which corresponds to the input audio.
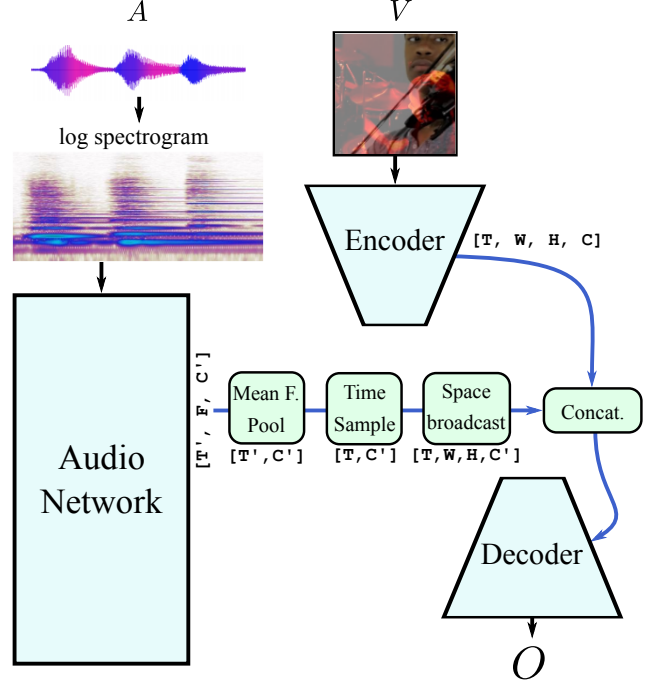


Figure 3: **The Controllable Compositional Centrifuge ($C^3$).** The Encoder-Decoder components are the same as in $C^2$ (Figure 2a). Audio features are extracted from the audio control signal and fused with the visual features before entering the decoder.

One is to use *deterministic control* where the desired video is forced to be output in a specific pre-defined output slot, without loss of generality $O_1$ is used. While simple, this strategy might be too rigid as it imposes too many constraints onto the network. For example, a network might naturally learn to output guitars in slot 1, drums in slot 2, *etc.*, while *deterministic control* is forcing it to change this ordering at will. This intuition motivates our second strategy – *internal prediction* – where the network is free to produce output videos in any order it sees fit, but it also provides a pointer to the output slot which contains the desired video. Internal prediction is trained jointly with the rest of the network, full details of the architecture are given in Appendix B. The training procedure and losses for the two control strategies are described in the next section.

## 3.3. Training procedure

**Training data.** Since it is hard to obtain supervised training data for the video decomposition problem, we adopt and extend the approach of [3] and synthetically generate the training data. This by construction provides direct access to one meaningful ground truth decomposition. Specifically, we start from two real videos $V_1, V_2 \in \mathbb{R}^{T \times W \times H \times 3}$. These videos are mixed together to generate a training video $V \in \mathbb{R}^{T \times W \times H \times 3}$:

$$V = \boldsymbol{\alpha} \odot V_1 + (\mathbf{1} - \boldsymbol{\alpha}) \odot V_2, \tag{3}$$

where $\boldsymbol{\alpha} \in [0,1]^{T \times W \times H}$ is a composing mask.

We explore two ways to generate the composing mask $\boldsymbol{\alpha}$. The first one is *transparent blending*, used by [3], where $\boldsymbol{\alpha} = \frac{1}{2}\mathbf{1}$. While attractive because of its simplicity, it does not capture the full complexity of the real world compositions we wish to address, such as occlusions. For this reason, we also explore a second strategy, referred to as *occlusion blending*, where $\boldsymbol{\alpha}$ is allowed to vary in space and takes values 0 or 1. In more detail, we follow the procedure of [13] where spatio-temporal SLIC superpixels [1] are extracted from $V_1$, and one is chosen at random. The compositing mask $\boldsymbol{\alpha}$ is set to 1 inside the superpixel and 0 elsewhere; this produces mixtures of completely transparent or completely opaque spatio-temporal regions. The impact of the $\boldsymbol{\alpha}$ sampling strategy on the final performance is explored in Section 4.1.

**Training loss: without control.** By construction, for an input training video $V$ we know that one valid decomposition is into $V_1$ and $V_2$. However, when training *without control*, there is no easy way to know beforehand the order in which output videos are produced by the network. We therefore optimize the network weights to minimize the following permutation invariant reconstruction loss [3]:

$$\mathcal{L}_{\text{pil}}\left(\{V_1, V_2\}, \boldsymbol{O}\right) = \min_{(i,j)|i \neq j} \ell(V_1, O_i) + \ell(V_2, O_j), \quad (4)$$

where $\ell$ is a video reconstruction loss, *e.g.* a pixel wise error loss (see Section 4 for our particular choice).

**Training loss: with control.** When training with audio as the control signal, the audio of one video ($V_1$ without loss of generality) is also provided. This potentially removes the need for the permutation invariant loss required in the no-control case, but the loss depends on the choice of control strategy. The two proposed strategies are illustrated in Figure 4 and described next.

*Deterministic control loss.* Here, the network is forced to output the desired video $V_1$ as $O_1$ so a natural loss is:

$$\mathcal{L}_{\text{det}}\left(\{V_1, V_2\}, \boldsymbol{O}\right) = \ell(V_1, O_1) + \ell(V_2, O_2). \quad (5)$$

Note that for this loss the number of output videos has to be restricted to $n = 2$. This limitation is another drawback of *deterministic control* as it allows less freedom to propose multiple output video options.

*Internal prediction loss.* In this strategy, the network freely decomposes the input video into outputs, and therefore the training loss is the same permutation invariant loss as for the no-control case (4). In addition, the network also points to the output which corresponds to the desired video, where the pointing mechanism is implemented as a module which outputs $n$ real values $\boldsymbol{s} = (s_i)_{i=1}^n$, one for each output video. These represent predicted dissimilarity between the
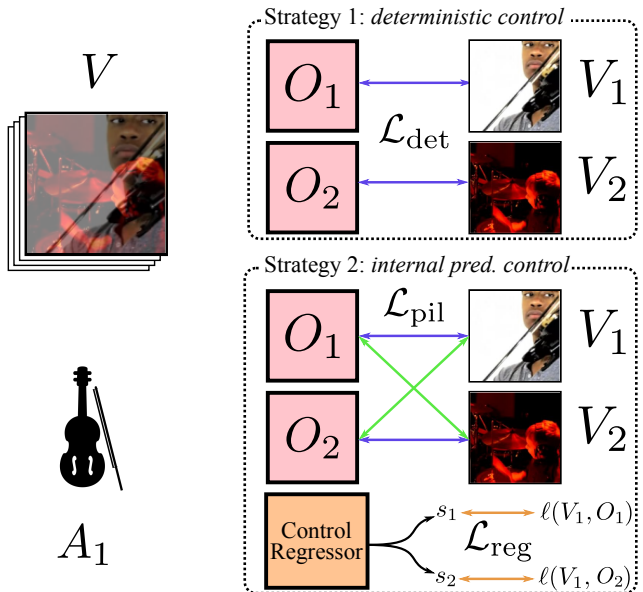


Figure 4: **Audio control strategies for video decomposition.** In this example, the inputs are the video $V$, a composition of $V_1$ showing a violin and $V_2$ showing drums, and an audio control signal, $A_1$, being the sound of the violin. With *deterministic control*, $V_1$ is forced to be put in output slot $O_1$ (and therefore $V_2$ in $O_2$). With *internal prediction control*, the network can freely order the output videos, so is trained with the permutation invariant loss, but it contains an additional control regressor module which is trained to point to the desired output.

desired video and output videos, and the attended output is chosen as $\arg\min_i s_i$. This module is trained with the following regression loss:

$$\mathcal{L}_{\text{reg}}\left(V_1, \boldsymbol{s}\right) = \sum_{i=1}^{n} |s_i - \ell(V_1, \mathbf{sg}(O_i))|, \quad (6)$$

where $\mathbf{sg}$ is the `stop_gradient` operator. Stopping the gradient flow is important as it ensures that the only effect of training the module is to learn to point to the desired video. Its training is not allowed to influence the output videos themselves, as if it did, it could sacrifice the reconstruction quality in order to set an easier regression problem for itself.

## 4. Experiments

This section evaluates the merits of the proposed Compositional Centrifuge ($\text{C}^2$) compared to previous work, performs ablation studies, investigates attention control via the audio control signal and the effectiveness of the two proposed attention control strategies of the Controllable Compositional Centrifuge ($\text{C}^3$), followed by qualitative decomposition examples on natural videos, and evaluation on the downstream task of action recognition.

| Model | Loss (Transp.) | Loss (Occl.) | Size |
|---|---|---|---|
| Identity | 0.364 | 0.362 | – |
| Centrifuge [3] | 0.149 | 0.253 | 22.6M |
| CentrifugePC [3] | 0.135 | 0.264 | 45.4M |
| $C^2$ w/o masking | 0.131 | 0.200 | 23.4M |
| $C^2$ | **0.120** | **0.190** | 27.1M |

Table 1: Model comparison in terms of average validation loss for synthetically generated videos with transp(arency) and occl(usions), as well as size in millions of parameters. All the results are obtained using models with $n = 4$ output layers. *CentrifugePC* is the predictor-corrector centrifuge [3], *Identity* is a baseline where the output videos are just copies of the input.

**Implementation details.** Following [3, 34], in all experiments we use the following video reconstruction loss, defined for videos $U$ and $V$ as:

$$\ell(U, V) = \frac{1}{2T} \left( \sum_t \|U_t - V_t\|_1 + \|\nabla(U_t) - \nabla(V_t)\|_1 \right),$$

where $\| \cdot \|_1$ is the L1 norm and $\nabla(\cdot)$ is the spatial gradient operator.

All models are trained and evaluated on the blended versions of the training and validation sets of the Kinetics-600 dataset [10]. Training is done using stochastic gradient descent with momentum for 124k iterations, using batch size 128. We employed a learning rate schedule, dividing by 10 the initial learning rate of 0.5 after 80k, 100k and 120k iterations. In all experiments we randomly sampled 64-frame clips at 128x128 resolution by taking random crops from videos whose smaller size being resized to 148 pixels.

## 4.1. Quantitative analysis

In this section, we evaluate the effectiveness of our approaches through quantitative comparisons on synthetically generated data using blended versions of the Kinetics-600 videos.

**Effectiveness of the $C^2$ architecture for video decomposition.** The baseline visual centrifuge achieves a slightly better performance (lower loss) than originally reported [3] by training on clips which are twice as long (64 vs 32 frames). As can be seen in Table 1, our proposed architecture outperforms both the *Centrifuge* baseline [3], as well as the twice as large predictor-corrector model of [3]. Furthermore, both of our architectural improvements – the masking and the composition module – improve the performance (recall that the baseline *Centrifuge* is equivalent to $C^2$ without the two improvements). The improvements are especially apparent for *occlusion blending* since our architecture is explicitly designed to account for more complicated real-world blending than the simple *transparency blending* used in [3].
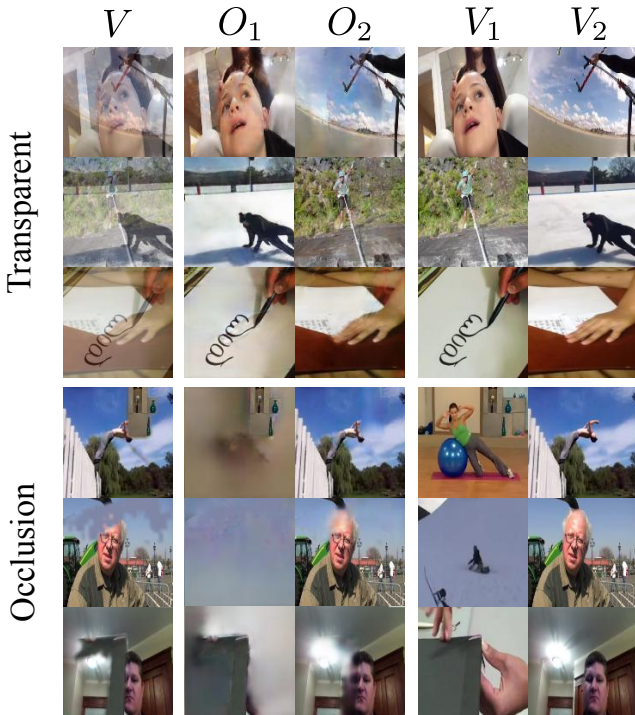


Figure 5: **Outputs of $C^2$ on blended Kinetics validation clips.** Each row shows one example via a representative frame, with columns showing the input blended clip $V$, two output videos $O_1$ and $O_2$, and the two ground truth clips $V_1$ and $V_2$. Top three rows show the network is able to successfully decompose videos with transparencies. Bottom three rows show synthetic occlusions – this is a much harder task where, apart from having to detect the occlusions, the network also has to inpaint the occluded parts of each video. $C^2$ performs satisfactory in such a challenging scenario.

| Model | Loss (Transp.) | Control Acc. |
|---|---|---|
| $C^2$ | 0.120 | 50% (chance) |
| $C^3$ w/ deterministic control | 0.191 | 79.1% |
| $C^3$ w/ internal prediction | 0.119 | 77.7% |

Table 2: Model comparison on average validation reconstruction loss and control accuracy. The controllable models, $C^3$, use audio as the control signal.

**Attention control.** The effectiveness of the two proposed attention control strategies using the audio control signal is evaluated next. Apart from comparing the reconstruction quality, we also contrast the methods in terms of their *control accuracy*, *i.e.* their ability to output the desired video into the correct output slot. For a given video $V$ (composed of videos $V_1$ and $V_2$) and audio control signal $A_1$, the output is deemed to be correctly controlled if the chosen output slot $O_c$ reconstructs the desired video $V_1$ well. Recall that the 'chosen output slot' is simply slot $O_c = O_1$ for the *deterministic control*, and predicted by the *control regressor* as $O_{\arg\min_i(s_i)}$ for the *internal pre-*
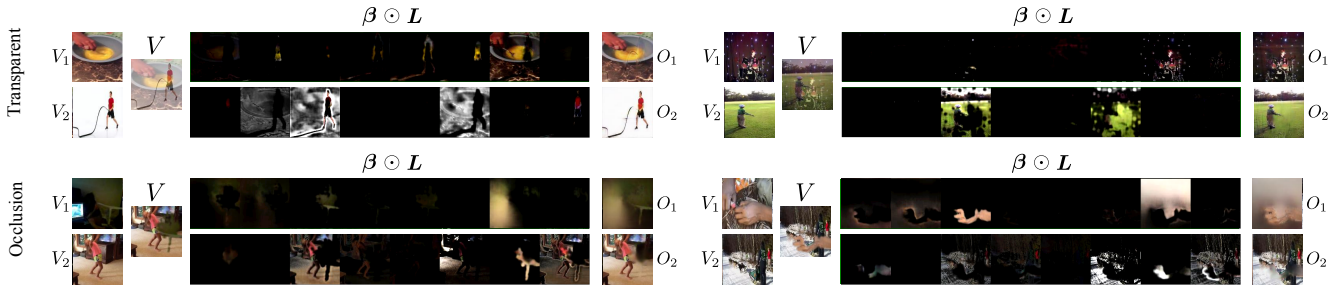
Figure 6: **Visualization of the internals of the compositional model.** Recall that the $C^2$ model produces the output videos via the *composition module* (Figure 2b) which multiplies the layers $\boldsymbol{L}$ with composing coefficients $\boldsymbol{\beta}$. Here we visualize the individual $\boldsymbol{\beta} \odot \boldsymbol{L}$ terms which when added together form the output videos. It can be observed that the layers and composing coefficient indeed decompose the input video $V$ into its constituent parts, for both the transparent and occlusion blending.

*diction control*. The chosen output video $O_c$ is deemed to reconstruct the desired video well if its reconstruction loss is the smallest out of all outputs (up to a threshold $t = 0.2 * (\max_i \ell(V_1, O_i) - \min_i \ell(V_1, O_i))$ to account for potentially nearly identical outputs when outputing more than 2 layers): $\ell(V_1, O_c) < \min_i \ell(V_1, O_i) + t$.

Table 2 evaluates control performance across different models with the *transparency blending*. It shows that the non-controllable $C^2$ network, as expected, achieves control accuracy equal to random chance, while the two controllable variants of $C^3$ indeed exhibit highly controllable behaviour. The two strategies are comparable on control accuracy, while *internal prediction control* clearly beats *deterministic control* in terms of reconstruction loss, confirming our intuition that *deterministic control* imposes overly tight constraints on the network.

## 4.2. Qualitative analysis

Here we perform qualitative analysis of the performance of our decomposition networks and investigate the internal layered representations.

Figure 5 shows the video decompositions obtained from our $C^2$ network for transparent and occlusion blending. The network is able to almost perfectly decompose the videos with transparencies, while it does a reasonable job of reconstructing videos in the much harder case where strong occlusions are present and it needs to inpaint parts of the videos it has never seen.

The internal representations produced by our *layer generator*, which are combined in the *composition module* to produce the output videos, are visualized in Figure 6. Our architecture indeed biases the model towards learning compositionality as the internal layers show a high degree of independence and specialize towards reconstructing one of the two constituent videos.

Finally, Figure 7 shows qualitative results for the best controllable network, $C^3$ with internal prediction, where audio is used as the control signal. The network is able to accurately predict which output slot corresponds to the desired
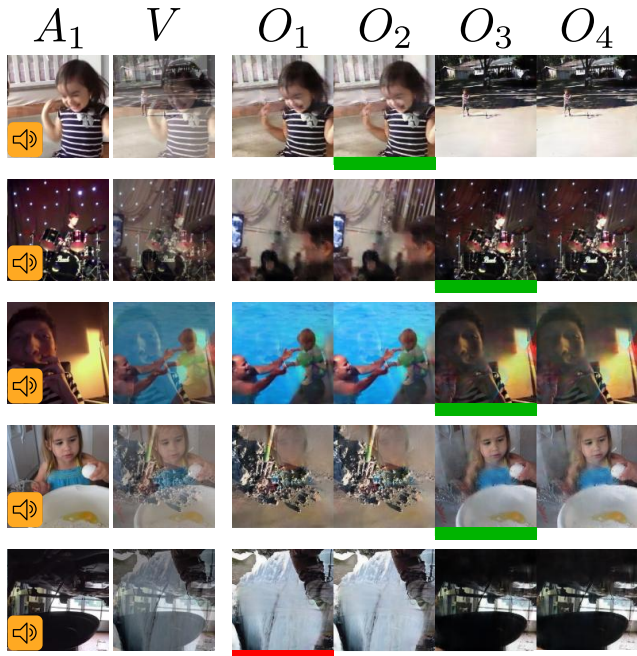


Figure 7: **Qualitative results of $C^3$ with internal prediction.** For visualization purposes, as it is hard to display sound, we show a frame of the video from which we use the *audio* as control on the left most column ($A_1$). $V$ (second column) represents the visual input to the model. The right 4 columns are the outputs of $C^3$. All examples exhibit good reconstruction error. The first four rows illustrate accurate control behaviour, where $C^3$ has correctly predicted the output that corresponds to the *control* signal (illustrated by a green marker under the frame). The last row illustrates an incorrect control (specified with a red marker under the wrongly chosen frame), where $C^3$ was fooled by a liquid sound that is plausible in the two scenarios.

video, making few mistakes which are often reasonable due to the inherent noisiness and ambiguity in the sound.

## 4.3. Downstream tasks

In the following, we investigate the usefulness of layered video decomposition as a preprocessing step for other

downstream tasks.

**Graphics.** Layered video decomposition can be used in various graphics applications, such as removal of reflections, specularities, shadows, *etc*. Figure 8 shows some examples of decompositions of real videos. Compared with previous work of [3], as expected from the quantitative results, the decompositions are better as the produced output videos are more pure.

**Action recognition.** A natural use case for video decomposition is action recognition in challenging scenarios with transparencies, reflections and occlusions. Since there are no action recognition datasets focused on such difficult settings, we again resort to using blended videos. A pre-trained I3D action recognition network [11] is used and its performance is measured when the input is pure unblended video, blended video, and decomposed videos, where the decomposition is performed using the best baseline model (predictor-corrector centrifuge, CentrifugePC [3]) or our Compositional Centrifuge ($C^2$). For the pure video performance, we report the standard top-1 accuracy.

For *transparency* blended videos, the desired outputs are both ground truth labels of the two constituent videos. Therefore, the models make two predictions and are scored 1, 0.5 and 0 depending on whether both predictions are correct, only one or none is, respectively. When I3D is applied directly on the blended video, the two predictions are naturally obtained as the two classes with the largest scores. For the decomposition models, each of the two output videos contributes their highest scoring prediction.

In the case of *occlusion* blended videos, the desired output is the ground truth label of $V_2$ because there is not enough signal to reconstruct $V_1$ as the blended video only contains a single superpixel from $V_1$. When I3D is applied directly on the blended video, the top prediction is used. The decomposition models tend to consistently reconstruct $V_2$ in one particular output slot, so we apply the I3D network onto the relevant output and report the top-1 accuracy.

Table 3 shows that decomposition significantly improves the action recognition performance, while our $C^2$ strongly outperforms the baseline CentrifugePC [3] for both blending strategies. There is still a gap between $C^2$ and the pure video performance, but this is understandable as blended videos are much more challenging.

## 5. Conclusion

General vision systems, that can serve a variety of purposes, will probably require controllable attention mechanisms. There are just too many possible visual narratives to investigate in natural scenes, for a system with finite computational power to pursue them all at once, always. In this paper we proposed a new compositional model for layered video representation and introduced techniques to make the resulting layers selectable via an external control signal – in
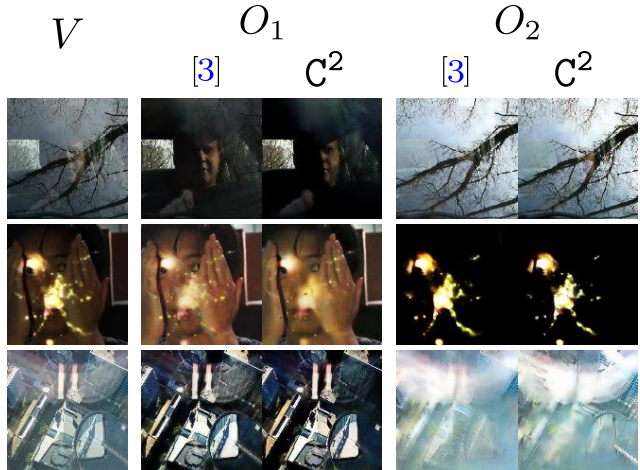


Figure 8: **Comparison of our $C^2$ model against [3] on real-world videos.** The input video is shown on the left, and the output videos of $C^2$ and [3] are interleaved in the remaining columns for easier comparison. While both models manage to decompose the videos reasonably well, $C^2$ achieves less leakage of one video into another. For example, $C^2$ versus [3] output $O_1$ (first row) removes the reflections of branches on the right side better, (second row) has fewer yellow circles of light, and (third row) makes the large circular reflection in the top half of the image much fainter.

| Mode | Acc. (Transp.) | Acc. (Occl.) |
|------|:---:|:---:|
| I3D – pure video | 59.5 | 59.5 |
| I3D | 22.1 | 21.3 |
| CentrifugePC [3] + I3D | 34.4 | 21.5 |
| $C^2$ + I3D | 40.1 | 24.7 |

Table 3: Action recognition accuracy on the Kinetics-600 validation set when the input to a pre-trained I3D classifier is a pure – non-blended – video (top row), a blended video directly passed through I3D, or a blended video that is first unblended using a layer decomposition model. The two columns show accuracies for two different blending processes: transparent and occluding.

this case sound. We showed that the proposed model can better endure automatically generated transparency and especially occlusions, compared to previous work, and that the layers are selected based on sound cues with accuracies of up to 80% on the blended Kinetics dataset. As future work we would like to train our model on more naturally-looking occlusions, possibly by generating the composing mask using supervised segmentations instead of unsupervised superpixels.

## References

[1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. SLIC superpixels compared to state-of-the-art superpixel methods. *TPAMI*, 2012.

[2] Triantafyllos Afouras, Joon Son Chung, and Andrew Zisserman. The Conversation: Deep Audio-Visual Speech En-

hancement. In *Interspeech*, 2018.

[3] Jean-Baptiste Alayrac, João Carreira, and Andrew Zisserman. The visual centrifuge: Model-free layered video representations. In *CVPR*, 2019.

[4] Relja Arandjelović and Andrew Zisserman. Objects that sound. In *ECCV*, 2018.

[5] Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.

[6] Daniel Baldauf and Robert Desimone. Neural mechanisms of object-based attention. *Science*, 2014.

[7] Zohar Barzelay and Yoav Y. Schechner. Harmony in motion. In *CVPR*, 2007.

[8] Efrat Be'ery and Arie Yeredor. Blind Separation of Superimposed Shifted Images Using Parameterized Joint Diagonalization. *Transactions on Image Processing*, 2008.

[9] Yuri Y. Boykov and Marie-Pierre Jolly. Interactive graph cuts for optimal boundary & region segmentation of objects in N-D images. In *ICCV*, 2001.

[10] Joao Carreira, Eric Noland, Andras Banki-Horvath, Chloe Hillier, and Andrew Zisserman. A short note about kinetics-600. In *arXiv preprint arXiv:1808.01340*, 2018.

[11] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017.

[12] Zhixiang Chi, Xiaolin Wu, Xiao Shu, and Jinjin Gu. Single Image Reflection Removal Using Deep Encoder-Decoder Network. In *arXiv preprint arXiv:1802.00094*, 2018.

[13] Carl Doersch and Andrew Zisserman. Sim2real transfer learning for 3D pose estimation: motion to the rescue. In *arXiv preprint arXiv:1907.02499*, 2019.

[14] Howard E. Egeth and Steven Yantis. Visual attention: Control, representation, and time course. *Annual review of psychology*, 1997.

[15] Ariel Ephrat, Inbar Mosseri, Oran Lang, Tali Dekel, Kevin Wilson, Avinatan Hassidim, William T. Freeman, and Michael Rubinstein. Looking to Listen at the Cocktail Party: A Speaker-Independent Audio-Visual Model for Speech Separation. In *SIGGRAPH*, 2018.

[16] Qingnan Fan, Jiaolong Yang, Gang Hua, Baoquan Chen, and David Wipf. A Generic Deep Architecture for Single Image Reflection Removal and Image Smoothing. In *ICCV*, 2017.

[17] Hany Farid and Edward H. Adelson. Separating reflections and lighting using independent components analysis. In *CVPR*, 1999.

[18] Yossi Gandelsman, Assaf Shocher, and Michal Irani. "Double-DIP": Unsupervised image decomposition via coupled deep-image-priors. In *CVPR*, 2019.

[19] Ruohan Gao, Rogeris Feris, and Kristen Grauman. Learning to separate object sounds by watching unlabeled video. In *CVPR*, 2018.

[20] Michael Gazzaniga and Richard B Ivry. *Cognitive Neuroscience: The Biology of the Mind: Fourth International Student Edition*. WW Norton, 2013.

[21] Laurent Girin, Jean-Luc Schwartz, and Gang Feng. Audio-visual enhancement of speech in noise. *The Journal of the Acoustical Society of America*, 2001.

[22] Xiaojie Guo, Xiaochun Cao, and Yi Ma. Robust Separation of Reflection from Multiple Images. In *CVPR*, 2014.

[23] David Harwath, Adrià Recasens, Dídac Surís, Galen Chuang, Antonio Torralba, and James Glass. Jointly discovering visual objects and spoken words from raw sensory input. In *ECCV*, 2018.

[24] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language. In *ICCV*, 2017.

[25] Romain Hennequin, Bertrand David, and Roland Badeau. Score informed audio source separation using a parametric model of non-negative spectrogram. In *ICASSP*, 2011.

[26] Daniel Heydecker, Georg Maierhofer, Angelica I. Aviles-Rivero, Qingnan Fan, Carola-Bibiane Schönlieb, and Sabine Süsstrunk. Mirror, mirror, on the wall, who's got the clearest image of them all? – A tailored approach to single image reflection removal. *arXiv preprint arXiv:1805.11589*, 2018.

[27] Laurent Itti, Christof Koch, and Ernst Niebur. A model of saliency-based visual attention for rapid scene analysis. *PAMI*, 1998.

[28] Laurent Itti, Geraint Rees, and John Tsotsos. *Neurobiology of attention*. Academic Press, 2005.

[29] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *CVPR*, 2015.

[30] Bruno Korbar, Du Tran, and Lorenzo Torresani. Cooperative learning of audio and video models from self-supervised synchronization. In *NeurIPS*, 2018.

[31] Kun Gai, Zhenwei Shi, and Changshui Zhang. Blind Separation of Superimposed Moving Images Using Image Statistics. *PAMI*, 2012.

[32] Luc Le Magoarou, Alexey Ozerov, and Ngoc Q.K. Duong. Text-informed audio source separation. Example-based approach using non-negative matrix partial co-factorization. *Journal of Signal Processing Systems*, 2015.

[33] Eric Mash and David Wolfe. *Abnormal child psychology*. Cengage Learning, 2012.

[34] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. In *ICLR*, 2016.

[35] Ajay Nandoriya, Mohamed Elgharib, Changil Kim, Mohamed Hefeeda, and Wojciech Matusik. Video Reflection Removal Through Spatio-Temporal Optimization. In *ICCV*, 2017.

[36] Aude Oliva, Antonio Torralba, Monica S Castelhano, and John M Henderson. Top-down control of visual attention in object detection. In *ICIP*, volume 1, pages I–253. IEEE, 2003.

[37] Andrew Owens and Alexei A. Efros. Audio-visual scene analysis with self-supervised multisensory features. In *ECCV*, 2018.

[38] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-Net: Convolutional Networks for Biomedical Image Segmentation. In *MICCAI*, 2015.

[39] Dana Segev, Yoav Y. Schechner, and Michael Elad. Example-based cross-modal denoising. In *CVPR*, 2012.

[40] Arda Senocak, Tae-Hyun Oh, Junsik Kim, Ming-Hsuan Yang, and In So Kweon. Learning to localize sound source in visual scenes. In *CVPR*, 2018.

[41] Aidean Sharghi, Boqing Gong, and Mubarak Shah. Query-focused extractive video summarization. In *ECCV*, 2016.

[42] Markus Siegel, Tobias H Donner, Robert Oostenveld, Pascal Fries, and Andreas K Engel. Neuronal synchronization along the dorsal visual pathway reflects the focus of spatial

attention. *Neuron*, 2008.

[43] Paul Smith, Tom Drummond, and Roberto Cipolla. Layered motion segmentation and depth ordering by tracking edges. *PAMI*, 26(4):479–494, 2004.

[44] Yale Song, Jordi Vallmitjana, Amanda Stent, and Alejandro Jaimes. TVSum: Summarizing web videos using titles. In *CVPR*, 2015.

[45] Charles Spence and Jon Driver. Audiovisual links in endogenous covert spatial attention. *Journal of Experimental Psychology: Human Perception and Performance*, 22(4):1005, 1996.

[46] Richard Szeliski, Shai Avidan, and P. Anandan. Layer extraction from multiple images containing reflections and transparency. In *CVPR*, 2000.

[47] Durk Talsma, Daniel Senkowski, Salvador Soto-Faraco, and Marty G Woldorff. The multifaceted interplay between attention and multisensory integration. *Trends in cognitive sciences*, 2010.

[48] John Tsotsos, Scan M. Culhane, Winky Yan Kei Wai, Yuzhong Lai, Neal Davis, and Fernando Nuflo. Modeling visual attention via selective tuning. *Artificial Intelligence*, 1995.

[49] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *NIPS*, 2017.

[50] John YA Wang and Edward H Adelson. Representing moving images with layers. *IEEE Transactions on Image Processing*, 3(5):625–638, 1994.

[51] John Y. A. Wang and Edward H. Adelson. Representing moving images with layers. *Transactions on image processing*, 1994.

[52] Wenwu Wang, Darren Cosker, Yulia Hicks, Saeid Sanei, and Jonathon. Chambers. Video assisted speech source separation. In *ICASSP*, 2005.

[53] Michael A Webster. Color vision: Appearance is a many-layered thing. *Current Biology*, 2009.

[54] Kelvin Xu, Jimmy Ba, Ryan Kiros, Kyunghyun Cho, Aaron Courville, Ruslan Salakhudinov, Rich Zemel, and Yoshua Bengio. Show, attend and tell: Neural image caption generation with visual attention. In *ICML*, pages 2048–2057, 2015.

[55] Tianfan Xue, Michael Rubinstein, Ce Liu, and William T. Freeman. A Computational Approach for Obstruction-Free Photography. In *SIGGRAPH*, 2015.

[56] Yi Yang, Sam Hallman, Deva Ramanan, and Charless Fowlkes. Layered object detection for multi-class segmentation. In *CVPR*, 2010.

[57] Steven Yantis. Control of visual attention. *Attention*, 1998.

[58] Xuaner Zhang, Ren Ng, and Qifeng Chen. Single image reflection separation with perceptual losses. *CVPR*, pages 4786–4794, 2018.

[59] Hang Zhao, Can Gan, Andrew Rouditchenko, Carl Vondrick, Josh McDermott, and Antonio Torralba. The sound of pixels. In *ECCV*, 2018.

## Overview

In this appendix, we cover three additional aspects: (i) in Section A we include an additional qualitative comparison on real videos, for which there was not space in the original manuscript; (ii) Section B provides details for the network architecture; and finally (iii) in Section C we study how the network uses the audio for control, by perturbing the audio.

## A. Additional comparison of $C^2$ with previous work

We compare to previous work [3, 55] on the task of reflection removal in Figure 9. One of the baselines [55] uses geometrical modelling and optimization but under strict assumptions (e.g. rigid motion). The second baseline [3] is trained on the same data as our model. The proposed model generates a sharp video with little reflection left.

## B. Architecture details

Figure 11 illustrates the architecture employed for $C^2$ while Figure 12 provides full details about the architecture employed for $C^3$ with *internal prediction* control strategy.

## C. Additional quantitative study for $C^3$

What aspect of the audio control signal is used for the controlled decomposition? One hypothesis is that the network latches onto low-level synchronization cues, so that the desired output video is identified as the one that is in sync with the audio. An alternative is that the desired video is the one whose semantic content matches the audio.

To answer this question, we use the best trained $C^3$ network with internal prediction control and evaluate its performance with respect to varying degrees of audio offset. The experiment is performed on the validation set of Kinetics-600. Reconstruction loss remains completely unaffected by shifting audio, while control accuracy deteriorates slightly as the offset is increased, as shown in Figure 10. The results suggest that the network predominantly uses the semantic information contained in the audio signal as control accuracy only decreases by 1.4 percentage points with the largest offsets where the audio does not overlap with the visual stream. However, some synchronization information is probably used as audio offset does have an adverse effect on control accuracy, and there is a sharp drop at relatively small offsets of 0.5-1s. There is scope for exploiting the synchronization signal further as it might provide a boost in control accuracy. A potential approach includes using a training curriculum analogous to [30].
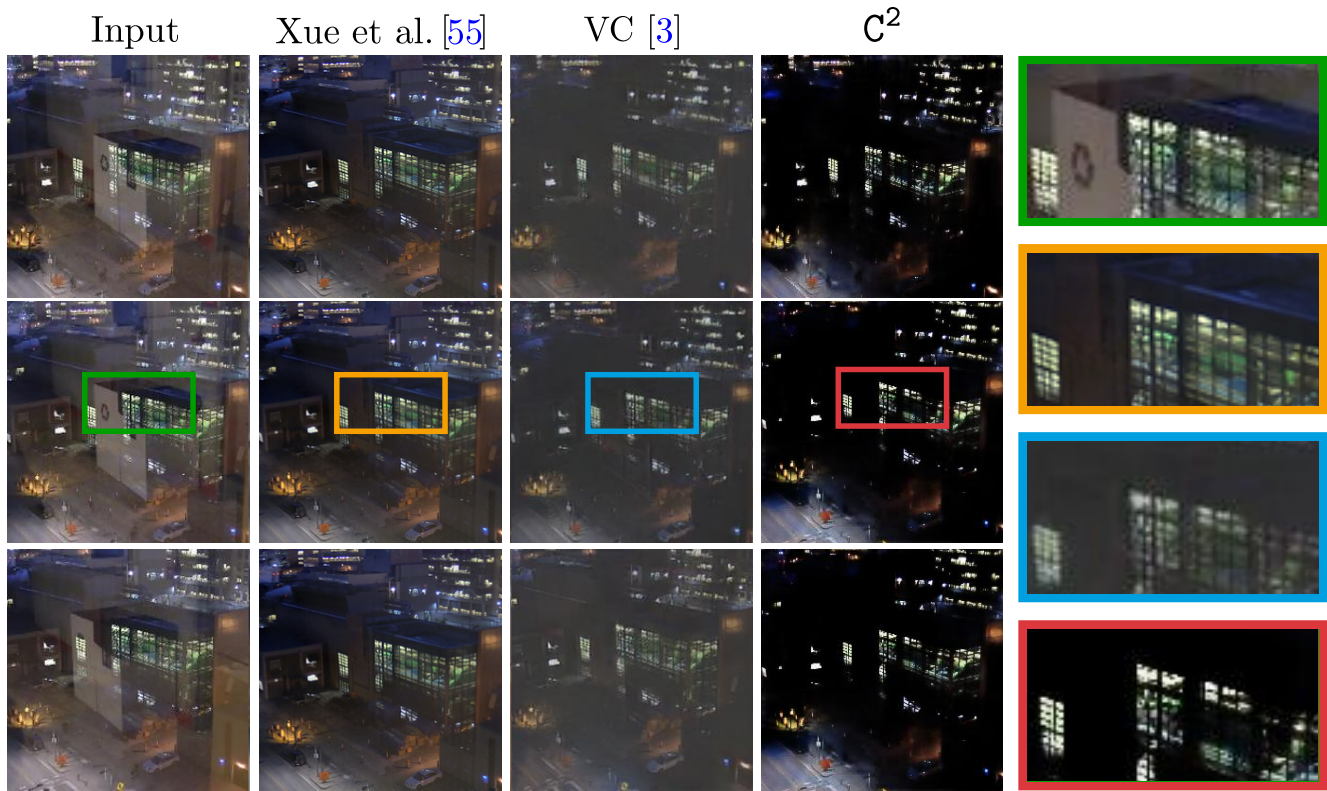
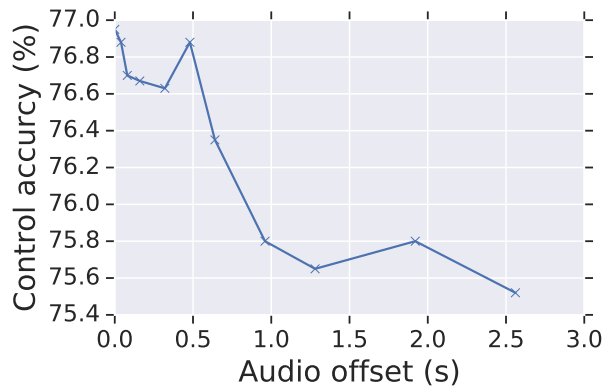Figure 9: Qualitative comparison of $C^2$ with other works.



Figure 10: Effect of shifting the control audio signal on the control accuracy. Note that the network was trained and tested on 2.56 second clips, so a 2.56 second offset corresponds to no overlap between the audio and visual streams.
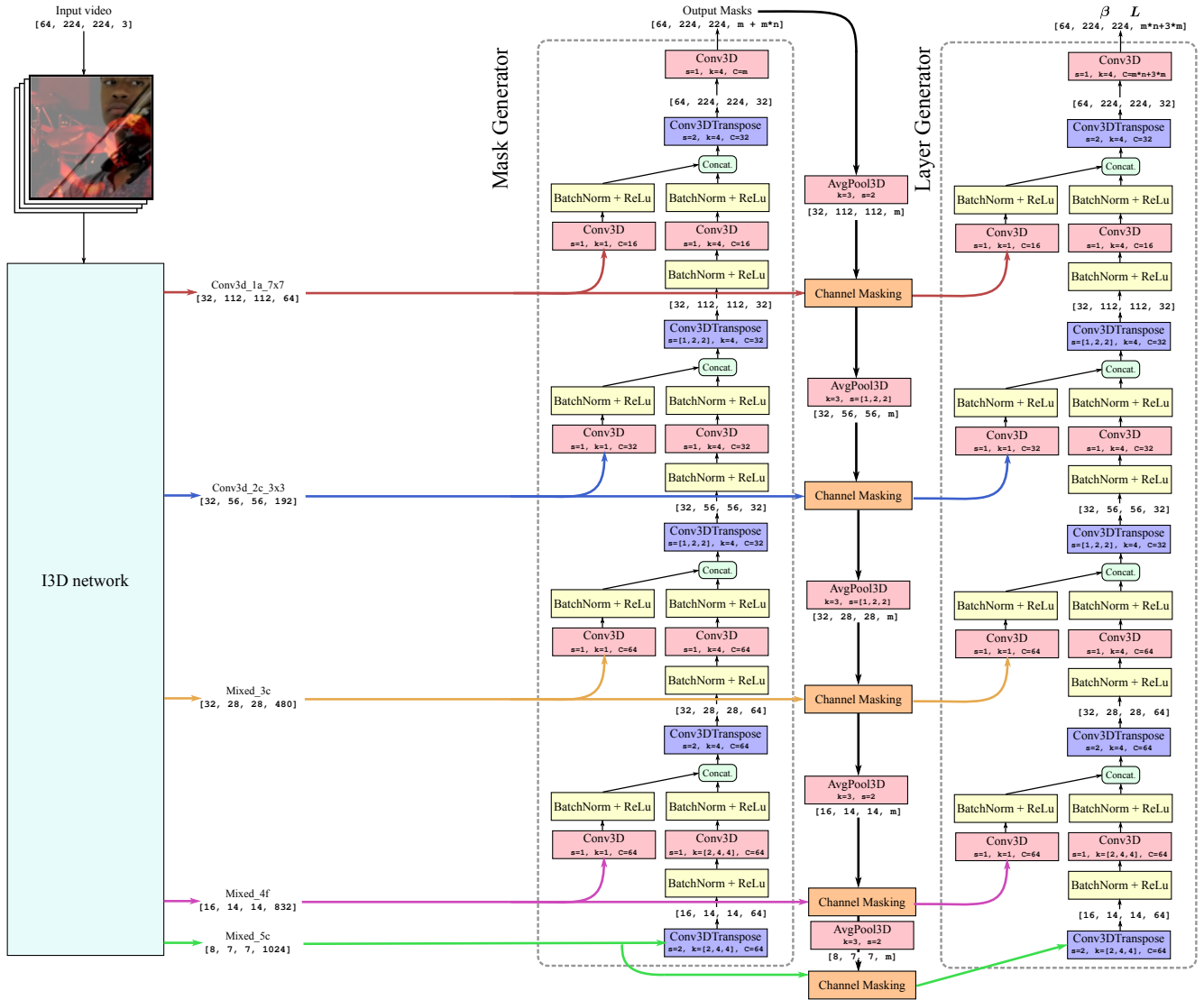
Figure 11: Details of the architecture used for $C^2$. The 'Channel Masking' block corresponds to the masking procedure described in equation (1) of the main paper.

Figure 12: Details of the architecture used for $C^3$ with *internal prediction* control. The 'Channel Masking' block corresponds to the masking procedure described in equation (1) of the main paper. The 'Audio/Visual Fusion' block matches the **Audio-visual fusion** procedure described in Section 3.2 of the main paper and illustrated in Figure 3.