

# Self-Supervised Learning of Perceptually Optimized Block Motion Estimates for Video Compression

Somdyuti Paul, Andrey Norkin, and Alan C. Bovik

**Abstract**—Block based motion estimation is integral to inter prediction processes performed in hybrid video codecs. Prevalent block matching based methods that are used to compute block motion vectors (MVs) rely on computationally intensive search procedures. They also suffer from the aperture problem, which tends to worsen as the block size is reduced. Moreover, the block matching criteria used in typical codecs do not account for the resulting levels of perceptual quality of the motion compensated pictures that are created upon decoding. Towards achieving the elusive goal of perceptually optimized motion estimation, we propose a search-free block motion estimation framework using a multi-stage convolutional neural network, which is able to conduct motion estimation on multiple block sizes simultaneously, using a triplet of frames as input. This composite block translation network (CBT-Net) is trained in a self-supervised manner on a large database that we created from publicly available uncompressed video content. We deploy the multi-scale structural similarity (MS-SSIM) loss function to optimize the perceptual quality of the motion compensated predicted frames. Our experimental results highlight the computational efficiency of our proposed model relative to conventional block matching based motion estimation algorithms, for comparable prediction errors. Further, when used to perform inter prediction in AV1, the MV predictions of the perceptually optimized model result in average Bjøntegaard-delta rate (BD-rate) improvements of -1.73% and -1.31% with respect to the MS-SSIM and Video Multi-Method Assessment Fusion (VMAF) quality metrics, respectively, as compared to the block matching based motion estimation system employed in the SVT-AV1 encoder.

**Index Terms**—Block motion estimation, inter prediction, video compression, convolutional neural networks, motion vectors, AV1.

## I. INTRODUCTION

**B**LOCK based motion estimation and compensation techniques are crucial to systems that perform such diverse tasks as video coding, frame rate up-conversion, 3D scene reconstruction etc. Motion estimation plays a key role in block based hybrid video codecs by facilitating the process of inter prediction, which reduces temporal redundancies inherent in natural videos. In fact, a block based hybrid encoder derives much of its compression capability from the inter prediction process. The use of a wide range of variable block sizes in modern codecs has further increased the complexity of the motion estimation process as compared to legacy codecs that

used block sizes from a more limited range. Such codecs thereby avoid the extreme expense of exhaustive search when estimating block motion vectors (MVs), by instead using fast block matching algorithms that greatly reduce computational complexity at the expense of higher prediction errors. However, optimization of the motion estimation process still largely relies on traditional pixelwise error (match) criteria, such as the sum of squared errors (SSE) or the sum of absolute differences (SAD). Simple pointwise measures of signal fidelity like these are unreliable predictors of perceptual picture quality [1], which is the visual attribute signal processing systems ultimately should optimize. Thus, it is quite plausible that the efficiency of the motion estimation modules used in current video encoders could be improved by incorporating measures of perceptual quality into their design.

While significant research efforts have been directed towards improving search based fast block matching algorithms [2], [3], the premise of using data driven deep learning techniques to learn block MVs has not received much research attention, despite successes attained on the related problem of dense optical flow estimation. Here, we address the question of perceptual optimization of block motion estimation, by introducing a data-driven, self-supervised framework for learning block translations for video encoding. We show that our method can be directly used in the inter prediction step of an AV1 encoder [4], [5], the latest royalty free video coding format. A multi-stage convolutional neural network (CNN) model is used to predict MVs for all block sizes required for AV1 inter-prediction. This composite block translation network, which we refer to as CBT-Net, is designed to predict the MVs of bidirectionally predicted frames (B-frames) using a frame triplet as its input. Each input triplet comprises a B-frame that is to be predicted through motion compensation, along with past and future reference frames. At each of its output stages, the model predicts MVs of blocks of a certain size, starting from the largest blocks at the first stage. Motion compensated frame predictions are generated using a spatial transformer module [6], allowing the model to be trained in a self-supervised manner, using the multi-scale structural similarity (MS-SSIM) measure [7] which is a widely used, accurate and differentiable measure of perceptual visual quality. In this way, the trained CBT-Net is optimized for MV estimation, and used to replace the search based motion estimation of the SVT-AV1 encoder [8]–[10]. We show that making this substitution results in significant perceptual rate-distortion (RD) gains when encoding AV1 bitstreams. The

Somdyuti Paul and Alan C. Bovik are with the Department of Electrical and Computer Engineering, University of Texas at Austin, Austin, TX, 78712 USA (email: somdyuti@utexas.edu, bovik@ece.utexas.edu).

Andrey Norkin is with Netflix Inc. Los Gatos, CA, 95032 USA (email: anorkin@netflix.com).

This work is supported by Netflix Inc.

©2022 IEEE. Personal use of this material is permitted. Permission from including reprinting/republishing this material for advertising or promotional purposes, creating new collective works, for resale or redistribution to servers or lists, or reuse of any copyrighted component of this work in other works.”

IEEE must be obtained for all other uses, in any current or future media, purposes, creating new collective works, for resale or redistribution to servers

overall proposed MV estimation framework is schematically outlined in Fig. 1.

We begin by noting that it is difficult to apply objective perceptual quality models like MS-SSIM to optimize block motion estimation algorithms used in video compression, since objective measurements of perceptual quality are less meaningful on the relatively small prediction block sizes that remain prevalent in current video codecs, alongside larger prediction block sizes. However, using CNN based model design, we are able to maximize the perceptual quality on entire predicted frames, instead of minimizing pixelwise prediction errors on individual blocks. This makes it possible to exploit differentiable perceptual quality models, like SSIM or MS-SSIM in the task of motion estimation, leading to improved RD performance not only against MS-SSIM, but also with respect to the equally widely-used perceptual Video Multi-Method Assessment Fusion (VMAF) [11] algorithm, and even against the traditional peak signal-to-noise ratio (PSNR). The primary contributions of our work can be summarized as follows:

- 1) We developed a composite model that learns to predict MVs over multiple block sizes in a hierarchical fashion, eliminating search processes required in previous block matching based ME approaches.
- 2) We exploited the spatial transformer module to achieve block translations in a continuous and learnable manner.
- 3) As an example, we substituted the MVs computed using hierarchical motion estimation and integer full search in the SVT-AV1 encoder with the trained CBT-Net predictions, and showed that it leads to a significant improvement in the RD performance.

We selected the SVT-AV1 encoder as the baseline since it is an open source AV1 encoder which performs well in terms of both compression efficiency and computational complexity, and since its motion estimation process uses the original video frames as reference frames rather than decoded ones (open loop motion estimation) for deriving integer pixel precision MVs, and thus fits our proposed motion estimation framework. Conversely, the standardized reference encoders generally deploy closed loop motion estimation architectures, where the decoded and reconstructed frames are used as reference frames. Nevertheless, many emerging practical encoders adopt the open loop architecture due to the advantage it imparts in terms of efficient parallelization and real-time processing performance, while the sequential nature of a typical closed loop architecture limits the parallelization that can be achieved. In SVT-AV1, only the fractional-pel refinement stage that follows the open loop integer precision MV estimation takes place in closed loop, as explained in Section IV-D. Further, the open loop motion estimation architecture substantially simplifies the process of generating the reference frames necessary to train data driven motion estimation models, since the source videos need not be encoded and decoded to generate reference frames as in the case of a closed loop architecture. Thus, our motion estimation framework can also be directly applied to other practical and scalable encoders such as SVT-HEVC [12] and SVT-VP9 [13], which use the open loop architecture, while

its application to closed loop encoder architectures such as HM [14] and VTM [15] necessitates a different training data generation policy from the one proposed here.

The rest of the paper is organized as follows. We review existing related work in Section II. Section III describes the frame triplet database which is essential to our self-supervised MV estimation approach. We introduce the proposed framework for learned motion estimation culminating in the CBT-Net model in Section IV. Section V presents the experimental results, based on which we draw conclusions in Section VI.

## II. RELATED WORK

Deep learning based methodologies are being increasingly explored as potentially efficacious alternatives to the traditional block based hybrid video coding framework. End-to-end deep video compression schemes such as [16], [17] have been shown to achieve a compression efficiency comparable to or better than mainstream codecs such as HEVC and VP9, adding to the considerable progress made in this direction. Many end-to-end video compression schemes such as [16] have employed CNN models to conduct optical flow estimation that is used to generate temporal predictions, subsequently compressing the motion information and residual signals using learned autoencoders. In [17] the motion estimation process was generalized to allow for efficient representation of complex motion types that cannot be captured by simple pixel translations, within an end-to-end compression framework. A bidirectional inter frame interpolation technique for video compression was developed in [18], where a pretrained optical flow estimation model was combined with a learned encoder-decoder pair to perform temporal prediction, and the residuals were computed in the latent space instead of pixel space to enable the use of existing deep image compression models to jointly encode the latent representations of key frames as well as residuals. Explicit motion estimation was avoided in the video compression framework introduced in [19] using displaced frame differences to capture motion regularities, along with a learned encoder-decoder network to achieve spatiotemporal reconstruction.

Learned video coding tools constitute another salient way in which deep learning based approaches have contributed to the advancement of video compression technology. Several deep learning based coding tools have been developed to optimize specific functions of hybrid video codecs, such as fractional interpolation [20], in loop filtering [21], [22], block partition prediction [23], [24], angular intra prediction [25], among others, to improve compression efficiency and/or reduce encoding complexity. Nevertheless, data driven learning based solutions to improve block based inter prediction have been less explored, despite this process being a key step in the traditional video coding pipeline. A combination of a fully connected network and a CNN was used to fuse spatial and temporal information for inter prediction to achieve coding gains in [26]. Learning based solutions were also developed to reduce blocking artifacts that often arise from block motion compensated prediction at low bitrates, using a CNN model to non-linearly combine bidirectional predictions in [27] and to refine predicted blocks in [28].

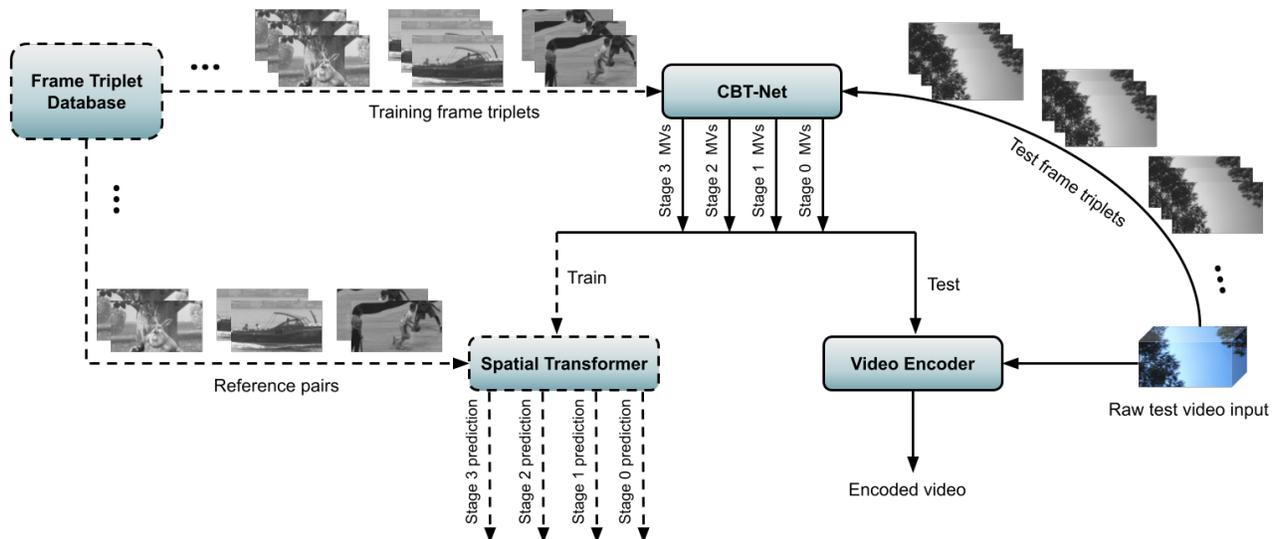


Fig. 1. Flow diagram of the proposed MV estimation framework.

The veritable improvements that image and video compression systems can derive by exploiting principles of human perception of visual information have been extensively emphasized by researchers [29]. However, rapid advancement in this direction is often stymied by challenges arising from our limited understanding of the complex phenomena that govern human perception, the computational complexity of models that seek to capture such perceptual effects, and by the intricacies of hybrid codec design that make it difficult to seamlessly incorporate and integrate novel components into the coding pipeline [30]. In [31], a deep learning based just-noticeable-quantization-distortion model was developed to eliminate perceptual redundancy in video coding, resulting in substantial coding gains. Properties of the human visual system, such as contrast sensitivity and visual masking, have also been used to enhance the quality of encoded videos by adaptively varying the quantization parameter [32], [33]. Saliency and visual attention models constitute yet another way of conducting bit allocation to achieve perceptual coding gains [34], [35]. Other methods of perceptual coding have employed reliable objective visual quality models like the structural similarity (SSIM) index [36] and VMAF [11]; a SSIM-based divisive normalization scheme was proposed in [37] that enhanced the coding efficiency and visual quality of hybrid codecs, while [38] substituted a proxy for the objective VMAF metric using a CNN regression model to derive a perceptual loss for training an end-to-end deep compression model.

Several block matching approaches have been developed to improve tradeoffs between the computational complexity of the motion search against the prediction error. A majority of these rely on using efficient search patterns such as diamond search (DS) [2] and adaptive rood pattern search (ARPS) [3], while few approaches employ learned models. Instead of searching for best matching blocks in pixel space, the authors of [39] proposed block matching in feature space, where convolutional features were extracted from block pixels,

and a simple average operator was used to perform representative matching in feature space to reduce the computational complexity of block motion estimation. In [40], a CNN model was trained to predict the similarity of two input blocks, where the trained model was subsequently used to perform block matching for frame interpolation. However, none of these MV estimation methods completely eliminate the search based block matching process. Furthermore, learned approaches like those in [39], [40] only predict the MVs for a single block size. Consequently, such MV estimation methods are not directly applicable for inter-prediction in hybrid video codecs, which rely on variable block sizes to obtain precise predictions while maintaining bitrate constraints. Additionally, the design of such models does not account for the perceptual quality of the resulting frame predictions.

The remarkable success of deep CNNs such as [41]–[45] on the related problem of dense optical flow estimation motivates our composite model design. Unlike existing models and algorithms that estimate block MVs, our proposed CBT-Net model entirely eliminates the search process required for block matching, and can collectively predict MVs over multiple block sizes required for the inter prediction process. Moreover, our approach also takes cognizance of the overall perceptual quality of the resulting motion compensated frame predictions by optimizing the MS-SSIM values of the predicted frames, rather than relying on perceptually agnostic pixelwise prediction error optimization that has been exclusively used as matching or optimization criteria in previous works.

### III. CONSTRUCTION OF FRAME TRIPLET DATABASE

A substantially large amount of training data is necessary to learn block MVs using a deep neural network. However, generating a large number of exact block MVs corresponding to each prediction block size used by the encoder is an onerous task, due to the extremely high computational complexity of the associated block matching process. Moreover, motion

estimation is inherently an ill-posed problem, due to the projection of the three dimensional true motion on a two dimensional plane, and the exact MVs between two frames cannot be theoretically determined without involving additional constraints. This signifies that even exact block matching algorithms can fail to estimate true motion in the presence of effects such as changes in illumination and occlusion. They are also impaired by the *aperture problem* which refers to the ambiguity in inferring the true direction of motion while viewing only a part of a moving object through a small aperture. This problem intensifies at smaller block sizes, as shown in Section V-B. Taking into account the shortcomings of block matching based motion estimation, here we instead propose a self-supervised training approach, which does not require exact block MVs as ground truth data for training the CBT-Net model. Instead, triplets of original frames are used to implicitly learn block translations, as elaborated further in Section IV.

#### A. Source Content

In order to create a database of frame triplets to perform self-supervised training of the CBT-Net model, we collected a diverse set of 109 uncompressed video sequences having native resolutions of  $1920 \times 1080$  pixels (1080p) or higher from publicly available sources [46]–[49]. We selected the sources to encompass different motion types, including camera motion, sports content, and computer graphics image (CGI) content, and different framerates ranging from 24 to 120 frames per second. A list of the contents selected for this database is provided in [50], specifying the sources from which the contents were collected, along with relevant information describing the lengths, original resolutions and framerates of each content. Some of the contents comprised multiple video shots. All sources were converted to 8 bit representation, and the luma channels were extracted. The sources which were at a higher resolution were downsampled to  $1920 \times 1080$  using Lanczos resampling, with kernel width  $\alpha = 3$ . If a source content comprised multiple scenes, it was split into two or more segments such that no segment contained any scene changes.

The sequences used to construct the frame-triplet database include many videos containing camera (ego) motion, as well as combinations of camera and object motions. Several videos were also captured with static cameras (object motion). For example, the training set was constructed using 87 video shots that have both object motion and ego motion, 10 video shots that exclusively contained camera motion over static scenes, and 26 video shots that were captured with static cameras.

#### B. Frame Triplet Extraction

Contemporary hybrid codecs typically use hierarchical prediction structures to conduct bidirectional inter prediction. A group of pictures (GOP) is subdivided into mini-GOPs, where the mini-GOP size is the number of frames between two consecutive unidirectionally predicted frames within a GOP. For example, a five-layer hierarchical temporal structure with a mini-GOP size of 16, as used in the SVT-AV1 encoder, is

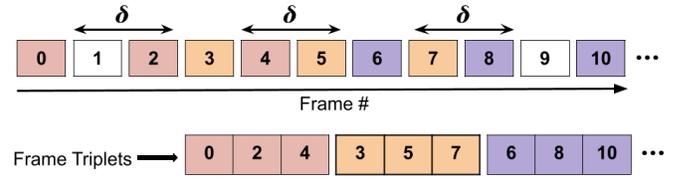


Fig. 2. Frame triplet extraction scheme used when creating the database (at layer 3).

TABLE I  
SUMMARY OF FRAME TRIPLET DATABASE FOR MV ESTIMATION.

Partition	# of contents	# of triplets samples for each layer			
		$\mathcal{S}_1$	$\mathcal{S}_2$	$\mathcal{S}_3$	$\mathcal{S}_4$
Training	97	13,943	16,959	17,169	13,166
Validation	12	1308	1361	1388	1058

illustrated in Fig. 3, where the top four layers correspond to B-frames; each B-frame uses a pair of reference frames, one aligned along each temporal direction. We used this prediction structure to guide our frame triplet database creation. Specifically, we extracted frame triplets to correspond to B-frames present at the top four temporal layers, such that the distances between the B-frames and their references at each layer are the same as those shown in Fig. 3.

Let  $(R_P, Q, R_F)_k$  constitute a frame triplet at layer  $k$  for  $k = 1 \dots 4$ , where  $Q$  denotes the intermediate frame that is to be predicted using block MVs computed with respect to a past reference frame  $R_P$  and a future reference frame  $R_F$ . Let  $d(F_i, F_j)$  denote the number of frames between any two frames  $F_i$  and  $F_j$  of a video sequence, i.e. the distance between the two frames. A set of frame triplets extracted at each temporal layer can then be represented as:

$$\mathcal{S}_k = \{(R_P, Q, R_F)_k : d(R_P, Q) = d(R_F, Q) = 2^{4-k}\}, \quad (1)$$

for  $k = 1 \dots 4$ . Since the prediction structure shown in Fig. 3 is symmetric,  $d(R_P, Q) = d(R_F, Q) \forall k$ . Our database  $\mathcal{S}$  is thus composed of four sets of frame triplets described by (1), i.e.  $\mathcal{S} = \{\mathcal{S}_1, \dots, \mathcal{S}_4\}$ .

We extracted frame triplets from each source content according to (1) to constitute each set  $\mathcal{S}_1, \dots, \mathcal{S}_4$ . To reduce the temporal redundancies between adjacent triplets extracted from the same content, we imposed a gap of  $\delta$  frames between adjacent triplets, as illustrated in Fig. 2 using layer 3 as an example. We used  $\delta = 2$  frames for layers 1, 2 and 3, and  $\delta = 3$  frames for layer 4. The sets  $\mathcal{S}_1, \dots, \mathcal{S}_4$  are individually partitioned into training and validation sets, with distinct contents in each partition, as summarized in Table I. The contents were divided between the training and validation partitions, such that each partition contains a range of different motion intensities and content types, such as CGI and sports contents. The contents included in each partition are also separately listed in [50].

As the distance between a B-frame to be predicted and its reference frames progressively increases from the top towards the bottom of the hierarchical prediction structure, it is reasonable to expect that the amounts of block displacements

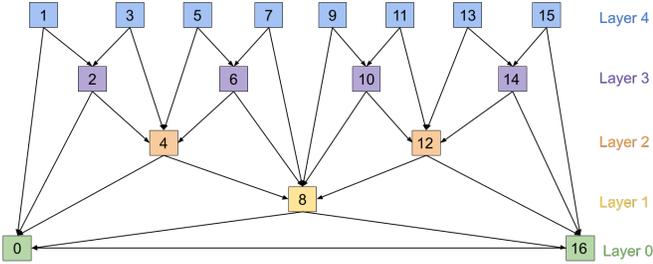


Fig. 3. Five-layer hierarchical temporal prediction structure used in SVT-AV1.

might differ among the temporal layers, with lower temporal layers being associated with larger block displacements. Thus, the sets  $\mathcal{S}_1, \dots, \mathcal{S}_4$  whose elements conform to the reference-to-predicted frame distances specified by the target prediction structure of Fig. 3, provide a systematic way of adequately representing the distributions of block displacements at each temporal layer. By independently training a separate instance of the CBT-Net model with the frame triplets from each of the four sets, the disparate amounts of block displacements corresponding to the different temporal layers can be effectively captured. Although the frame triplets comprising the database can also be extracted by following other policies, such as randomly choosing the distance between a B-frame and its references (within the limit of the mini-GOP size), to train a single instance of the model for all temporal layers, dividing the database among the four temporal layers, allows for more efficient training, since the separate instances of the model can be trained simultaneously with their corresponding datasets.

#### IV. PROPOSED METHOD

While our method is general, and could be applied to any number of prediction layers and different possible block sizes, we continue to exemplify our method by applying it in the SVT-AV1 framework. The open loop motion estimation process of the SVT-AV1 encoder uses four different prediction block sizes ( $64 \times 64$ ,  $32 \times 32$ ,  $16 \times 16$  and  $8 \times 8$ ) to generate inter predictions. Accordingly, we designed the CBT-Net model to concurrently predict the MVs of all non-overlapping  $64 \times 64$ ,  $32 \times 32$ ,  $16 \times 16$  and  $8 \times 8$  blocks of the input B-frame with respect to both past and future reference frames using a multi-stage process. The frames comprising the input triplet are padded with zeroes such that each frame contains an integral number of blocks of each size. Thus, each input triplet  $(R_P, Q, R_F)$  is of dimension  $W \times H \times 3$ , where  $W$  and  $H$  are the padded spatial width and height, respectively, that are divisible by 64.

Traditional MV estimation algorithms perform block matching between a pair of frames. However, using a triplet of frames as input to the CBT-Net model endows it with a notable computational advantage over using frame pairs for motion estimation. Using the frame pairs  $(R_P, Q)$  and  $(R_F, Q)$  as two separate inputs to the model to compute the forward and backward MVs, would require that all of the features of a B-frame  $Q$  be computed for each pair. Using the triplet

$(R_P, Q, R_F)$  as the model input instead expedites processing, since the features are extracted from  $Q$  only once, and are used to compute both the forward and backward MVs in the same forward pass. Indeed, we have observed  $\sim 30\%$  reduction in inference time using this strategy, with a marginal decline in prediction performance, as compared to using frame pairs as the model input.

#### A. CBT-Net Architecture

The architecture of the proposed CBT-Net is shown in Fig. 4, including the spatial transformer module that is used to translate blocks of the reference frames, using the predicted MVs to generate B-frame predictions. As Fig. 4 depicts, this model has four prediction stages, where the MVs of non-overlapping square blocks, of one of the four different sizes is predicted at each stage. The MVs of the largest blocks of size  $64 \times 64$  that are used for inter prediction by the SVT-AV1 encoder are predicted at the first stage of the CBT-Net model, followed by the MV predictions on progressively smaller blocks at the succeeding stages. Our motivation for this multi-stage prediction arises from the fact that larger blocks typically contain more distinctive features, making it easier for a model to learn good matches against blocks of the same size present in other frames. Conversely, on smaller blocks which contain fewer features, the *aperture problem* is exacerbated, rendering the block matching problem more ambiguous. With this in mind, the MVs of larger blocks can be used to guide the MV estimation of co-located smaller blocks in a coarse-to-fine manner, thereby helping to alleviate the aperture problem. Thus, each stage of the motion estimation process, from the second stage onward, relies on both the MVs estimated at the preceding stages, as well as on the feature maps from the preceding stages to make its predictions.

The CBT-Net model has three types of layers with distinct functionalities, as indicated in Fig. 4. The *feature extraction layers* are the convolutional layers that extract convolutional features from the input frame triplet. There are nine feature extraction layers in the CBT-Net model, each followed by the customary batch normalization and non-linearity operation, where the latter is implemented using rectified linear units (ReLU). A  $5 \times 5 \times 4$  convolutional layer is used as the last layer of each stage to act as the *prediction layer* that produces the MV predictions for that stage using the feature map it receives as input. By design, the first two channels of the output of each prediction layer correspond to the MVs of the past reference  $R_P$ , while the last two channels correspond to the MVs of the future reference  $R_F$ . The prediction layer of the first stage uses the feature map generated by the last feature extraction layer of the model, to predict translations of the  $64 \times 64$  blocks of the two reference frames  $R_P$  and  $R_F$  with respect to the frame  $Q$ . In the second stage, the feature map used as the input to the prediction layer of the first stage, along with the MV predictions from the first stage are individually passed to *upsampling layers* that upsample them by a factor of two, using transposed convolutions [41]. These two upsampled volumes and the convolutional feature map of the eighth feature extraction layer are concatenated and used as

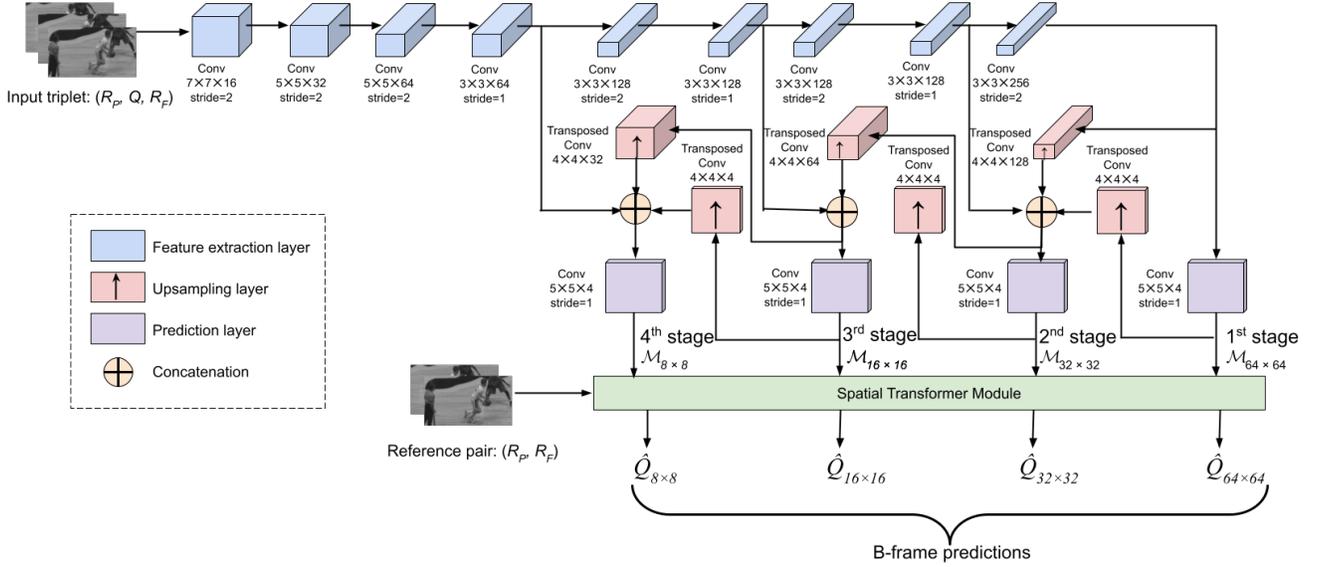


Fig. 4. Architecture of the multi-stage CBT-Net model, along with the spatial transformer module used to generate frame predictions.

the input to the prediction layer of the second stage, to generate the MV predictions on  $32 \times 32$  blocks. The MV outputs of the remaining stages are predicted by similarly concatenating the feature map of that stage with the upsampled versions of the feature map and the MV predictions from the previous stage, followed by applying the concatenated volume as input to a prediction layer, as shown in Fig. 4. Thus, by effectively re-using the motion information from the previous stages, the CBT-Net is able to predict the MVs of multiple block sizes, while incurring very little additional computational overhead as compared to traditional block motion estimation algorithms except the hierarchical ones, where a new search based block matching process must be performed for each different block size.

The computational complexity of conventional block matching based motion estimation algorithms increases as the search range used for matching is increased. However, limiting the search range to small values for the sake of computational tractability precludes accurate motion estimation, especially in the presence of fast moving objects in the scene. By contrast, the effective search range of our proposed approach is limited by the receptive field size of the last feature extraction layer of the model. Thus, the effective search range of the model can be increased by using more feature extraction layers, or by increasing the stride of the convolutional filters at each layer, allowing larger motions to be effectively captured with lower additional computational overhead, as compared to traditional methods. The size of the effective receptive field of the last layer of the CBT-Net model depicted in Fig. 4 is 255, as shown in Appendix A. Consequently, the effective search range of the model is  $\pm 127$  (even though no traditional search is conducted), and the outputs of the prediction layers at each stage were also clipped to lie in this range to give the composite output of the CBT-Net model, denoted by a set of four matrices  $\{\mathcal{M}_{64 \times 64}, \mathcal{M}_{32 \times 32}, \mathcal{M}_{16 \times 16}, \mathcal{M}_{8 \times 8}\}$ . If the spatial resolution of the input triplet is  $W \times H$ , then  $\mathcal{M}_{S \times S}$ ,

the output MV matrix for  $S \times S$  blocks is of dimension  $W/S \times H/S \times 4$ .

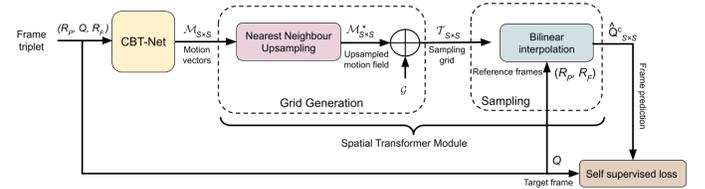


Fig. 5. Components of the spatial transformer module.

### B. Spatial Transformer Module

The CBT-Net model gives real-valued MV predictions denoted by  $\mathcal{M}_{S \times S}$ . However, the block translations to generate the frame predictions as performed in conventional block motion estimation algorithms would require integer valued MVs since the pixels of the reference frames are only defined at integer locations. Simply rounding the real valued MVs to the nearest integer for the purpose of block translations does not solve the problem since the gradients of the round function are zero almost everywhere which is not amenable to learning using back propagation. We propose a solution that overcomes this difficulty using a spatial transformer module. We designed the spatial transformer module to obtain the prediction signals through block translations in a manner that is conducive to back-propagation. A spatial transformer is a differentiable module that can be inserted within a CNN model to perform generic geometric transformations of an image or feature map [6]. In the context of motion estimation, spatial transformers have been employed to perform feature warping [43] and frame warping [44], [45] to learn dense optical flows. It was also used in [51] to learn affine transformations to conduct whole-frame predictions. Unlike these previous

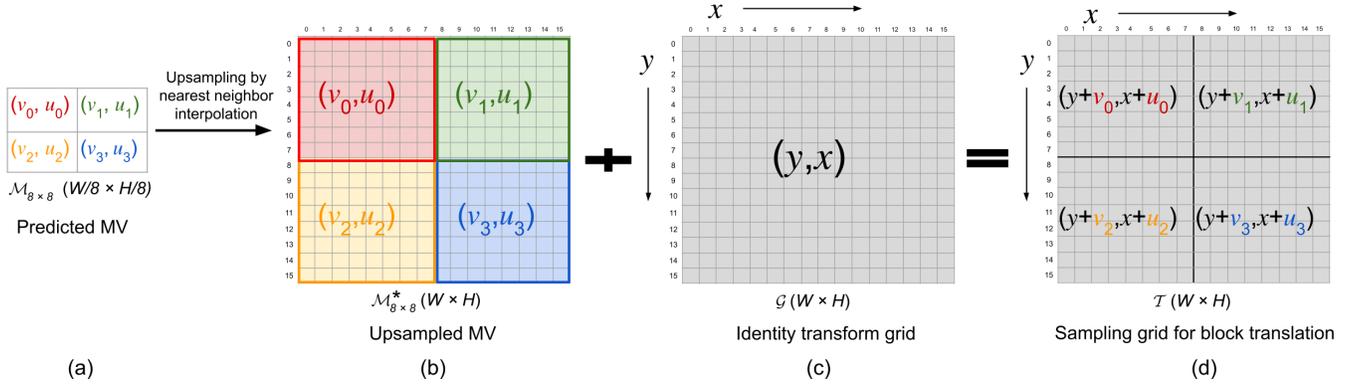


Fig. 6. Illustration of the grid generation process for blocks of size  $8 \times 8$  and  $W = H = 16$ , where a small spatial dimension is chosen for ease of illustration, and  $u_k$  and  $v_k$  are the vertical and horizontal components of the estimated MV of the  $k^{\text{th}}$  block with respect to a single reference frame.

methods, in our model, the spatial transformer module estimates differentiable local block translations, by an efficient grid generation process as described later in this section. Specifically, it serves to translate non-overlapping blocks of the reference frames  $R_P$  and  $R_F$  according to the predicted set of MVs  $\{\mathcal{M}_{64 \times 64}, \mathcal{M}_{32 \times 32}, \mathcal{M}_{16 \times 16}, \mathcal{M}_{8 \times 8}\}$  in a differentiable manner. This module thus generates the predictions of frame  $Q$  corresponding to each block size and reference frame in order to compute the model’s self-supervised loss. The components of the spatial transformer module as well as its role in the proposed framework are illustrated in Fig. 5. The inputs to the spatial transformer module are the set of MVs generated by the CBT-Net and the two reference frames. The two steps that constitute the spatial transformer module are described next:

- *Grid generation*: the purpose of this step is to find the locations of an input image or feature map that are to be sampled to generate the spatially transformed output image or feature map [6]. While the predicted elementwise motion parameters can be directly applied to perform warping through coordinate transformation in [43]–[45], [51], the CBT-Net generates the motion parameters at the block level. The grid generation process employed to accomplish block translations can be elucidated with the help of Fig. 6 for a single block size and reference frame. The predicted motion fields as shown in Fig. 6(a) are first upsampled to the spatial resolution of the input frames using nearest neighbor interpolation. Fig. 6(b) represents the upsampled motion field, where all pixels within the individual  $S \times S$  blocks get the same value of the pixelwise MVs as indicated by the different color coded regions; the block colors correspond to the ones used to represent the corresponding MVs in Fig. 6(a). This step effectively ensures that all pixels within the same block get the same MV  $(u_i, v_i)$ , thereby achieving the desired block translations. Let the upsampled MV matrix for  $S \times S$  blocks of size  $W \times H \times 4$  be  $\mathcal{M}_{S \times S}^*$ . If  $\mathcal{G}(x, y) = (x, y)$  is the sampling grid for the identity transformation as shown in Fig. 6(c), then the sampling

grid for translating the  $S \times S$  blocks is derived as

$$\begin{aligned} \mathcal{T}_{S \times S}(x, y, c) &= \mathcal{G}(x, y) + \mathcal{M}_{S \times S}^*(x, y, 2c : 2c + 1) \\ &= \left( x + \mathcal{M}_{S \times S}^*(x, y, 2c), \right. \\ &\quad \left. y + \mathcal{M}_{S \times S}^*(x, y, 2c + 1) \right) \end{aligned} \quad (2)$$

where  $c = 0$  for reference frame  $R_P$  and  $c = 1$  for reference frame  $R_F$ . Fig. 6(d) represents this sampling grid, and when the pixels of the reference frames are sampled according to this sampling grid  $\mathcal{T}_{S \times S}$ , all pixels within the non-overlapping  $S \times S$  blocks are translated by the amount determined by the block MVs  $(u_i, v_i)$ .

- *Sampling*: since the MV matrix  $\mathcal{M}_{S \times S}^*$  is real valued, the locations given by the corresponding sampling grid  $\mathcal{T}_{S \times S}$  are as well. We use bilinear interpolation to sample the pixels of the reference frame according to the locations given by  $\mathcal{T}_{S \times S}$  to obtain the corresponding prediction  $\hat{Q}_{S \times S}^c$  of frame  $Q$  obtained by translating  $S \times S$  blocks of the reference frame indexed by  $c$ .

Thus, for each reference frame, the output of the spatial transformer module is the set of predicted frames  $\mathcal{Q}^c = \{\hat{Q}_{64 \times 64}^c, \hat{Q}_{32 \times 32}^c, \hat{Q}_{16 \times 16}^c, \hat{Q}_{8 \times 8}^c\}$ , which are used to compute the self-supervised loss of the CBT-Net model, as described next.

### C. Loss Function

Self-supervised learning refers to machine learning schemes where the inputs signals of the models are reused to derive pseudo-labels that act as supervision signals for training the model. Such learning schemes have been proven to be useful to solve such problems where ground truth data is difficult to obtain. This is also true for our present problem as discussed earlier in Section III. Thus, the supervision signal for our model is derived using the input signals and the model predictions, resulting in a self-supervised learning scheme, whereby the model’s loss is computed between the input frame  $Q$ , which serves as a pseudo label, and its predictions  $\hat{Q}_{S \times S}^c$  corresponding to each reference frame and block size. The

predictions  $\hat{Q}_{S \times S}^c$  are generated by the spatial transformer module by using the model's MV predictions  $\mathcal{M}_{S \times S}$  to translate the blocks of the other two frames of the input triplet,  $R_P$  and  $R_F$ , which serve as the reference frames in either direction. We employed MS-SSIM [7] to derive this self-supervised loss when training the CBT-Net model as follows:

$$\mathcal{L} = \sum_{c=0}^1 \sum_{\hat{Q}_{S \times S}^c \in \mathcal{Q}^c} 10 \times \log(1 - \text{MS-SSIM}(Q, \hat{Q}_{S \times S}^c)), \quad (3)$$

where  $\text{MS-SSIM}(Q, \hat{Q}_{S \times S}^c)$  is the MS-SSIM score of the prediction  $\hat{Q}_{S \times S}^c$  with respect to the original frame  $Q$ . Thus, by minimizing the negative value of the MS-SSIM score (calculated in dB) between  $Q$  and its predictions obtained at each stage of the model as given by (3), the average objective visual quality of the inter frame predictions corresponding to all the block sizes was collectively improved by the training procedure.

#### D. Integration with SVT-AV1 Encoder

Block motion estimation takes place in two stages in the SVT-AV1 encoder: the full-pixel open loop ME performed before the main frame encoding process and a sub-pixel closed loop refinement performed alongside the mode decisions process. The open loop motion estimation process as implemented in the SVT-AV1 encoder itself consists of the following steps [8]:

- 1) *Hierarchical Motion Estimation (HME)*: first, a three stage HME is conducted at three different resolutions, corresponding to one-sixteenth resolution, one quarter resolution and the base resolution of the input frames, respectively. The goal of HME is to find the best candidate MV for each  $64 \times 64$  block. At each stage of HME, the frame is divided into multiple non-overlapping search areas, which are searched independently to produce multiple MV candidates for each  $64 \times 64$  block, using the locations given by the best candidate MVs from the previous stage as the starting points for the searches. The HME process then selects the best MV candidate for each  $64 \times 64$  block using the SAD metric as the block matching criterion.
- 2) *Integer full search*: the best MV candidates for the  $64 \times 64$  blocks given by HME are used in this step to find the MVs of all blocks of sizes  $64 \times 64, \dots, 8 \times 8$  that are encompassed by each  $64 \times 64$  block. The integer full search uses the displaced  $64 \times 64$  block locations given by the MVs computed at the previous HME step as the search center to estimate the MVs of all square blocks within each  $64 \times 64$  block using the SAD metric. This comprises the last stage of the open loop ME process, which refines the HME results to give integer pel MVs.

We replaced the above two sub-processes of the SVT-AV1 encoder's motion estimation process with the trained CBT-Net model. The MVs estimated by the integer full search process of the encoder are multiplied by a factor of four, in order to use the MVs in the quarter pel refinement process. Thus, in order to make the MV predictions of the CBT-Net suitable for

subsequent encoding processes, we scaled them in the same manner, and rounded the scaled MVs to the nearest integer values. The horizontal and vertical components of each MV were then combined to form a single 32 bit unsigned integer in accordance with the integer MV syntax used in SVT-AV1. Each block of size  $64 \times 64$  contains 85 blocks combining all four square prediction block sizes used in SVT-AV1, and the integer motion search process of the encoder collectively estimates the MVs of all 85 square blocks within each  $64 \times 64$  block as a single unit. So, we organized the four MV matrices corresponding to the four block sizes into a single matrix of size  $2 \times W/64 \times H/64 \times 85$ , where the first dimension corresponds to the number of reference frames. Finally, the HME and integer full search processes were disabled, and the combined MV matrices were used to provide the MV information for inter prediction while encoding.

## V. EXPERIMENTAL RESULTS

We trained the CBT-Net model to estimate MVs at two resolutions,  $1920 \times 1080$  and  $1280 \times 720$ . The training data for the smaller of the two resolutions were obtained by down-sampling the entire database of frame triplets originally created at  $1920 \times 1080$  to  $1280 \times 720$  using Lanczos resampling. At each resolution, four different instances of the CBT-Net model were trained for each of the four temporal layers shown in Fig. 3, using the frame triplets from the training partitions of the corresponding sets  $\mathcal{S}_1, \dots, \mathcal{S}_4$ . The CBT-Net model was comprehensively evaluated in terms of its prediction error, AV1 encoding performance achieved using the predicted MVs for inter-prediction, and computational complexity, as explained in Sections V-B V-C, and V-D, respectively. Although the MV prediction was carried out by the CBT-Net model using luma (Y) channels only, the predicted MVs were used to perform motion compensation of all three channels (Y, U, and V) of the test videos, which were encoded in YUV420p format.

#### A. System Settings and Hyperparameters

The CBT-Net model as shown in Fig. 4 has 1,914,832 trainable parameters. Each instance of the model was trained on four NVidia 1080-TI GPUs, while the prediction and encoding performance evaluation was conducted on a 64 bit, eight core Intel(R) Core(TM) i7-6700K CPU @ 4.00GHz system, running Ubuntu 18.04 (without GPU). All the evaluations conducted utilized all cores of the system. We used the Adam optimizer [52] with a learning rate of  $10^{-4}$  to optimize the parameters of each model instance. The batch size used was 8 in the case of 1080p triplets, and 16 in the case of 720p triplets, respectively.

#### B. Prediction Performance

In the absence of ground truth MVs, the pixelwise error between a B-frame  $Q$  and its prediction  $\hat{Q}_{S \times S}^c$  acts as an approximate measure of the precision of block MVs. We computed this prediction error as the mean absolute difference (MAD) value between  $Q$  and  $\hat{Q}_{S \times S}^c$  for all four block sizes, on the validation partition of each of the sets  $\mathcal{S}_1, \dots, \mathcal{S}_4$ .

TABLE II  
AVERAGE PREDICTION ERRORS IN TERMS OF MAD VALUES ON THE VALIDATION SET USING HME [8], ARPS [3] AND CBT-NET.

Temporal layer #	64 × 64 blocks			32 × 32 blocks			16 × 16 blocks			8 × 8 blocks		
	HME	ARPS	CBT-Net	HME	ARPS	CBT-Net	HME	ARPS	CBT-Net	HME	ARPS	CBT-Net
1	15.60	<b>10.48</b>	11.99	13.44	<b>7.07</b>	10.88	13.00	<b>5.20</b>	7.09	11.67	<b>4.08</b>	6.17
2	13.68	<b>9.29</b>	10.54	11.50	<b>6.02</b>	7.42	11.18	<b>4.25</b>	5.80	9.95	<b>3.23</b>	4.97
3	12.16	<b>8.27</b>	9.26	10.00	<b>5.13</b>	6.20	9.78	<b>3.47</b>	4.63	8.69	<b>2.55</b>	3.85
4	11.60	<b>7.49</b>	8.15	9.42	<b>4.47</b>	5.17	9.22	<b>2.91</b>	3.67	7.93	<b>2.07</b>	2.92

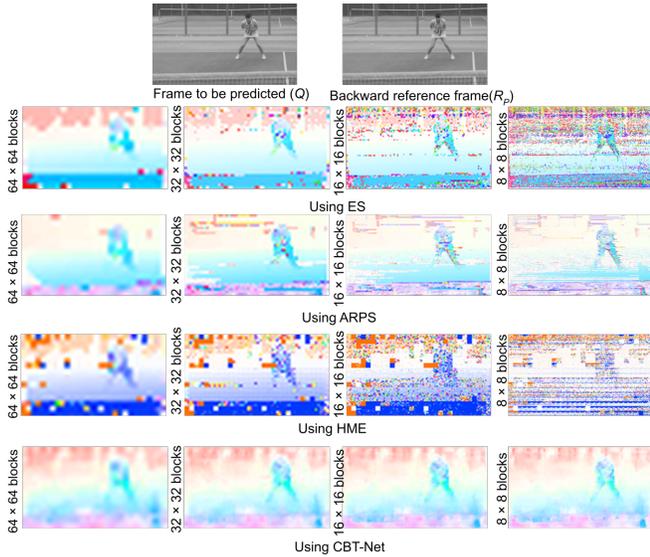
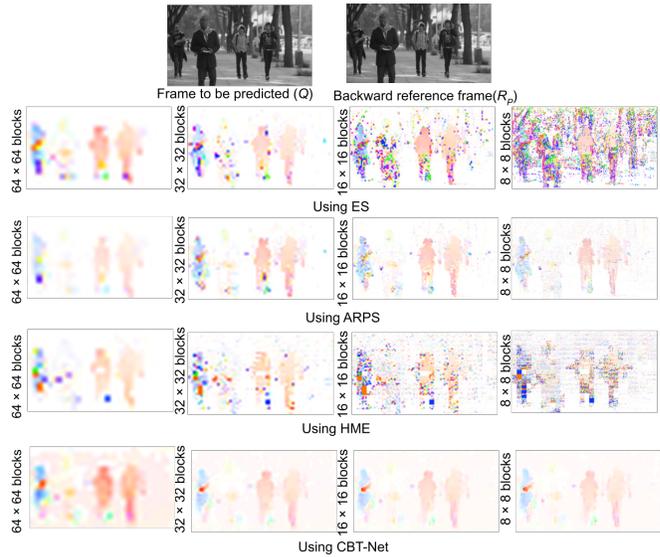
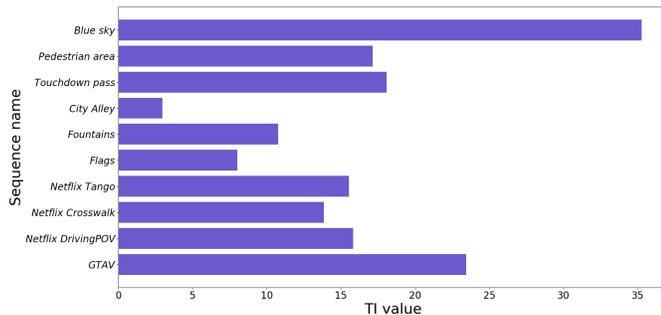
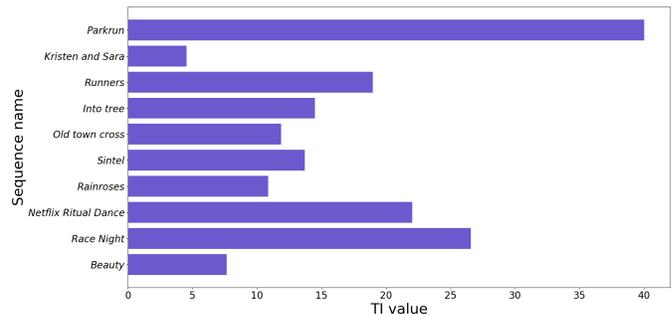
(a) Frames and color coded MVs for *Tennis* sequence.(b) Frames and color coded MVs for *Netflix Narrator* sequence.

Fig. 7. Visualization of the MVs computed between two frames with the ES, ARPS, HME and CBT-Net model from top to bottom; for each sequence and each method, the MVs computed using block sizes of 64 × 64, 32 × 32, 16 × 16, and 8 × 8 shown from left to right.



(a) 1080p sequences



(b) 720p sequences

Fig. 8. Distribution of TI values on test sequences at two resolutions.

We also compared our prediction performance to the HME implementation of the SVT-AV1 encoder [8] as well as the ARPS [3] method as presented in Table II<sup>1</sup>. While HME acts as the baseline for our work, the ARPS method was selected for comparison since it provides the best prediction error and computational complexity tradeoff among the widely used search pattern based block matching algorithms [53].

The average prediction errors reported in Table II show that

<sup>1</sup>We used the python implementations of ARPS available in the scikit-video library at <http://www.scikit-video.org/stable/motion.html>, and translated to optimized machine code using Numba (<https://numba.pydata.org/>), although actual encoders may use more optimized implementations.

HME has the largest prediction errors at all block sizes and temporal layers, while ARPS achieves the smallest values. The prediction errors of the CBT-Net model are much smaller than HME and are mostly comparable to that of ARPS.

The MVs obtained using the CBT-Net model are depicted in Fig. 7 for two sequences from the validation set, and compared against the corresponding MVs obtained using exhaustive search (ES), ARPS and HME based block matching procedures. We included the MVs given by ES in this comparison, as it produces the lowest pixelwise prediction errors among all block matching based motion estimation algorithms, although it is too computationally intensive to be

TABLE III  
PERFORMANCE OF CBT-NET MODEL AND ARPS [3] WITH RESPECT TO SVT-AV1 BASELINE.

Sequence	Resolution	# of frames	MS-SSIM BD-rates (%)		VMAF BD-rates(%)		PSNR BD-rates (%)		$\Delta T$ (%)		
			CBT-Net	ARPS	CBT-Net	ARPS	CBT-Net	ARPS	CBT-Net	ARPS	
<i>Blue sky</i>	1080p	209	<b>-0.95</b>	0.21	<b>-2.00</b>	-0.13	<b>-1.20</b>	-0.2	<b>10.2</b>	8.4	
<i>Pedestrian area</i>		369	<b>-3.22</b>	-2.10	<b>-3.92</b>	-2.94	<b>-2.62</b>	-1.48	<b>10.4</b>	8.3	
<i>Touchdown pass</i>		561	-6.62	<b>-6.84</b>	-5.10	<b>-5.37</b>	<b>-5.19</b>	-4.85	<b>10.2</b>	8.3	
<i>City Alley</i>		593	-3.39	<b>-5.08</b>	-2.42	<b>-3.49</b>	-1.73	<b>-2.43</b>	<b>5.0</b>	1.2	
<i>Fountains</i>		289	0.30	<b>-0.04</b>	<b>0.12</b>	0.13	0.24	<b>0.13</b>	<b>9.2</b>	8.4	
<i>Flags</i>		321	<b>0.41</b>	0.44	0.49	0.49	<b>0.03</b>	0.19	<b>9.7</b>	7.7	
<i>Netflix Tango</i>		289	<b>-1.48</b>	-1.00	<b>-1.03</b>	-0.81	<b>-1.22</b>	-0.79	<b>11.2</b>	9.6	
<i>Netflix Crosswalk</i>		289	<b>-3.79</b>	-1.77	<b>-3.13</b>	-1.49	<b>-2.72</b>	-1.51	<b>10.9</b>	8.8	
<i>Netflix DrivingPOV</i>		289	-1.84	<b>-1.93</b>	<b>0.30</b>	1.95	<b>-0.65</b>	0.01	<b>10.8</b>	9.9	
<i>GTA V</i>		689	<b>-2.31</b>	-1.71	<b>-2.27</b>	-1.08	<b>-1.96</b>	-1.14	<b>14.5</b>	13.4	
Average				<b>-2.29</b>	-1.98	<b>-1.90</b>	-1.27	<b>-1.70</b>	-1.21	<b>10.2</b>	8.4
<i>Parkrun</i>		720p	481	<b>-1.49</b>	-0.73	<b>-0.86</b>	-0.44	<b>-0.03</b>	0.77	<b>9.3</b>	8.4
<i>Kristen and Sara</i>	593		<b>2.57</b>	3.74	<b>0.72</b>	1.20	<b>0.73</b>	1.19	<b>6.1</b>	3.3	
<i>Runners</i>	289		<b>2.09</b>	2.87	<b>0.9</b>	1.36	<b>1.75</b>	1.99	<b>10.9</b>	9.8	
<i>Into tree</i>	497		<b>-1.89</b>	-0.73	<b>-0.84</b>	0.26	<b>-3.03</b>	-0.14	<b>8.4</b>	6.7	
<i>Old town cross</i>	497		0.82	<b>-3.59</b>	<b>1.75</b>	2.15	0.92	<b>-0.07</b>	<b>7.5</b>	5.3	
<i>Sintel</i>	337		-2.22	<b>-4.16</b>	<b>-1.38</b>	-1.20	<b>-1.59</b>	-0.79	<b>10.9</b>	8.5	
<i>Rainroses</i>	497		<b>0.08</b>	0.52	<b>0.07</b>	0.56	<b>-0.41</b>	0.62	<b>12.0</b>	10.4	
<i>Netflix Ritual dance</i>	417		<b>-2.64</b>	-1.60	<b>-1.73</b>	-1.45	<b>-2.12</b>	-1.64	<b>9.6</b>	8.2	
<i>Race Night</i>	593		<b>-7.13</b>	-5.91	<b>-5.83</b>	-4.77	<b>-5.74</b>	-4.34	<b>12.6</b>	11.0	
<i>Beauty</i>	593		<b>-1.98</b>	0.99	<b>0.05</b>	0.46	<b>-0.19</b>	0.54	<b>7.1</b>	4.7	
Average				<b>-1.18</b>	-0.86	<b>-0.72</b>	-0.19	<b>-0.97</b>	-0.19	<b>9.4</b>	7.6
Overall average				<b>-1.73</b>	-1.42	<b>-1.31</b>	-0.73	<b>-1.34</b>	-0.70	<b>9.8</b>	8.0

practicable for evaluation on the entire validation set as given in Table II. The MVs in Fig. 7 were color coded such that the hue and saturation represent the orientation and magnitude of the vectors, respectively [54]. From the first three rows of Figs. 7a and 7b, it is evident that for all the three block matching based approaches that were considered, the estimated motion fields become progressively noisier as the block size is reduced, thereby revealing that the aperture problem is more pronounced for smaller blocks, as discussed earlier in Section IV-A. In fact, in Figs. 7a and 7b, the ES algorithm that gives the lowest prediction errors also produces the most irregular motion fields at the two smallest block sizes as compared to the other three methods. By using the predicted motion information from larger blocks to guide estimation of the motion of smaller blocks, the CBT-Net model was able to avoid the aperture problem, yielding more coherent motion fields than the three other methods as shown in the bottom rows of Figs. 7a and 7b. The visual comparison in Fig. 7 shows that minimizing the pixelwise prediction error, as done by the block matching approaches such as ES, ARPS and HME was often unable to track the true motion of the objects for smaller block sizes, while our proposed multi-stage model accomplished consistent motion estimation even when block sizes were small. Since noisy motion fields are more expensive to encode in terms of the number of bits required to represent them, lower absolute pixelwise prediction errors do not necessarily ensure a better RD performance as we further demonstrate in Section V-C.

### C. RD Performance

We evaluated the RD performance of our block MV prediction approach on two resolutions (1080p and 720p), using ten test video sequences at each resolution. The contents of the test video sequences are distinctly different from the ones chosen for training and validation and they cover a wide range of

motion types, as indicated by the distribution of their temporal information (TI) values (computed as the maximum over time of the standard deviation of frame differences) [55] shown in Fig. 8.

The SVT-AV1 encoder operating at speed level 1 and using HME followed by integer full search and quarter pel refinement for inter-prediction was used as the baseline for this experiment. We also configured the encoder to use the temporal prediction structure shown in Fig. 3. The test sequences were encoded at four QP values (30, 40, 50 and 60), and the B-frames were predicted using the MVs obtained in three different ways as follows:

- 1) using HME (baseline configuration of the encoder)
- 2) using CBT-Net
- 3) using ARPS

The MVs obtained using ARPS were scaled and formatted in the same manner as described in the last paragraph of Section IV-D. The integer MVs of the P-frames from temporal layer 0 of Fig. 3 were computed using the integer motion search procedure of the encoder. As the mini-GOP size of the hierarchical prediction structure is 16, we encoded the first  $(\lfloor K/16 \rfloor \times 16) + 1$  frames of a sequence having  $K$  frames, to avoid an incomplete mini-GOP at the end of the sequence. The number of frames encoded for each sequence is listed in the third column of Table III. The qualities of the encodes thus obtained were evaluated using MS-SSIM [7], VMAF [11] and PSNR, where the first two are perceptual quality metrics. The RD performance on the test set, as measured by the Bjøntegaard delta bitrates (BD-rates) [56] with respect to each quality metric, are summarized in Table III.

Replacing the MVs estimated by the SVT-AV1 encoder with the CBT-Net model's MV predictions resulted in notable RD performance gains for all three metrics and at both resolutions tested, as shown in Table III, while the RD plots of two test sequences are presented in Fig. 9. The average BD-rates across

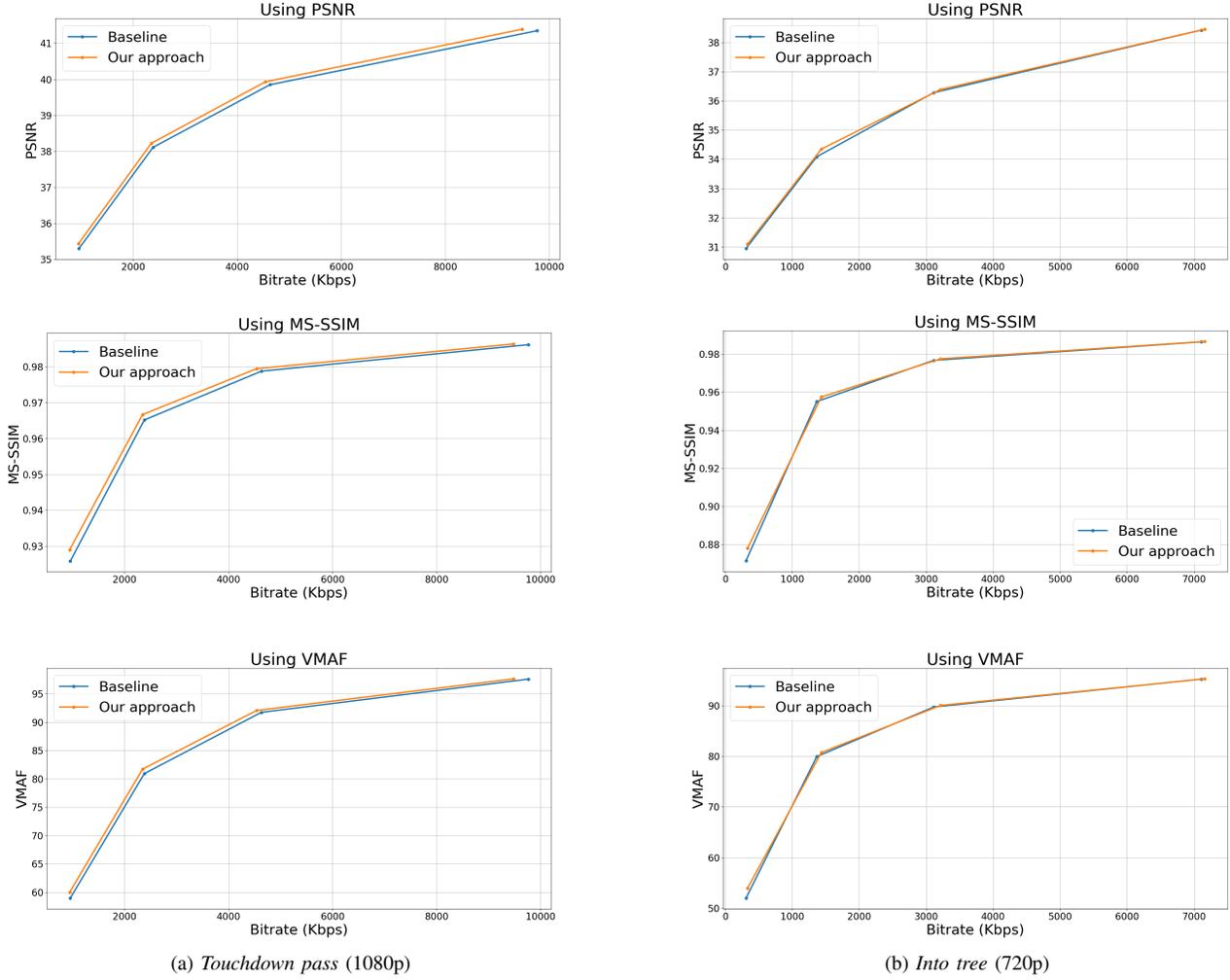


Fig. 9. RD plots of two test sequences.

all test sequences with respect to MS-SSIM VMAF and PSNR were -1.73%, -1.31% and -1.34%, respectively, showing consistent gains across the different metrics. The RD performance gain achieved against the VMAF metric shows that the model’s applicability to more general perceptual metrics, even though it was trained using MS-SSIM. This suggests that the CBT-Net learned to make MV predictions that contributed favorably towards optimizing the visual quality of the predicted B-frames. However, on some of the test sequences such as *Flags*, *Fountains*, *Kristen and Sara*, *Runners* etc. no RD gains were attained. As discussed in Section III-A, the dataset contains both sequences with camera motion and sequences with little to no camera motion. Our method is advantageous, and performs especially well on sequences that contain camera motion, while little or no performance advantage is attained on sequences with primarily static backgrounds, that have few moving regions.

Table III also compares the RD performance of our method with ARPS, which achieves average BD-rates of -1.42%, -0.73% and -0.70% with respect to MS-SSIM, VMAF and PSNR, respectively. Thus, although the ARPS method achieved slightly lower prediction errors than CBT-Net as

shown in Table II, the BD-rates it attains on the majority of the test sequences are higher than our method, substantiating our earlier claim that smoother motion fields are beneficial to attaining a better RD tradeoff at reasonably similar prediction errors. Fig. 10 compares the visual quality of the frames of two test sequences encoded with the baseline encoder against our MV estimation framework. Fewer compression artifacts are visible in the encoded frames obtained using our proposed framework, as shown in Fig. 10.

#### D. Encoding Complexity

The proposed block motion estimation approach improves the RD performance of the baseline without incurring an increase in computational complexity. In fact, it reduces the computational complexity of the original SVT-AV1 encoder to some extent due to the elimination of the searches needed by HME for block matching. For each test sequence, we computed the difference in the encoding time of our method or ARPS with respect to the baseline as  $\Delta T = \frac{T_0 - T}{T_0} \times 100$ , where  $T_0$  is the total encoding time of the baseline AV1 encoder for all QPs and  $T$  is the corresponding value when the motion estimation module is substituted with our model or

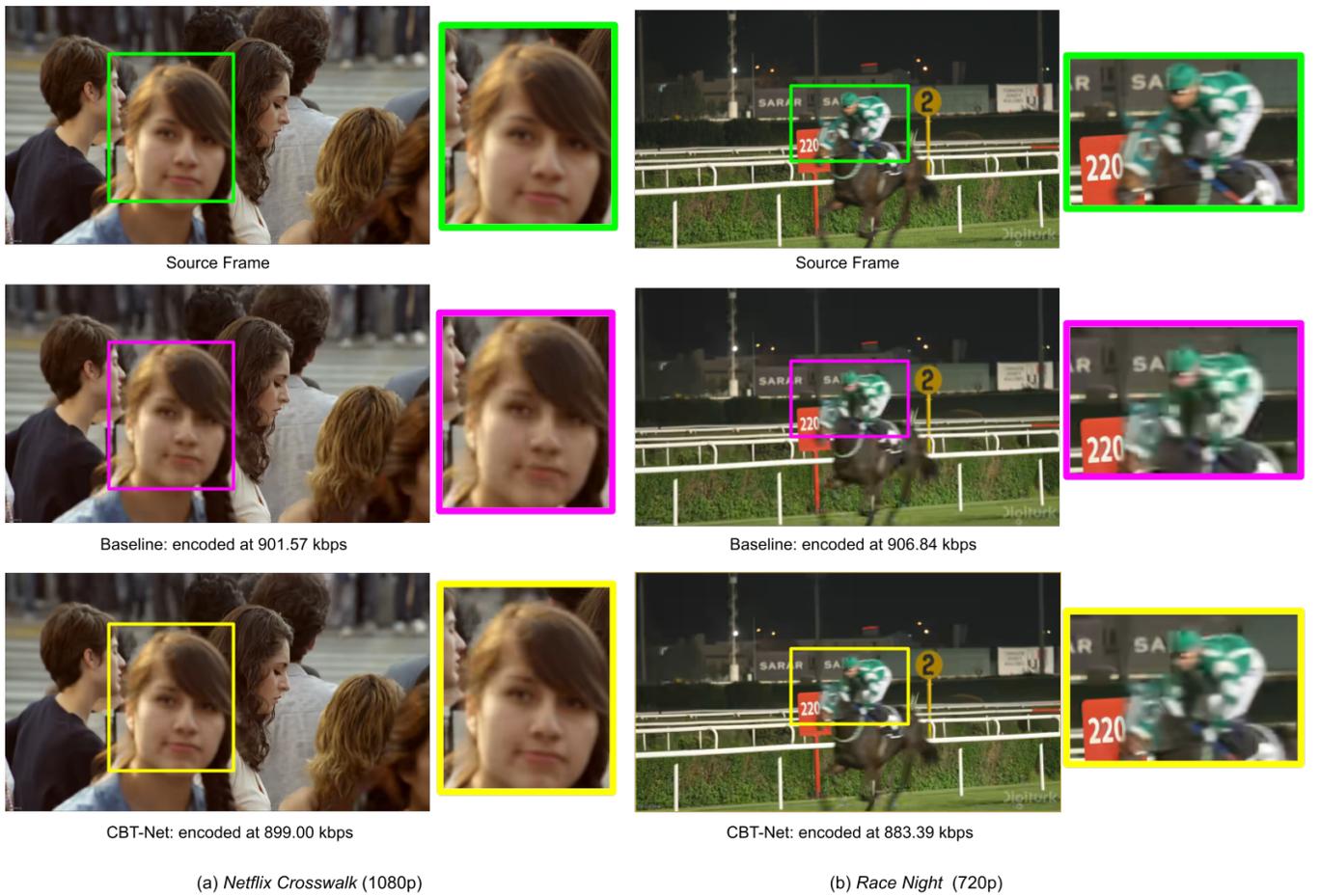


Fig. 10. Comparison of visual quality of the encoded frames obtained using the MVs computed by the SVT-AV1 baseline and the CBT-Net model.

ARPS. The time  $T$  includes the inference time of the CBT-Net model or ARPS on CPU.

The MV prediction time of the CBT-Net model is 0.14 seconds per frame, while the corresponding time for ARPS is 0.66 seconds per frame on the same CPU. ARPS takes longer to estimate the MVs as the underlying search process needs to be invoked eight times to estimate the MVs for four block sizes and two reference frames, while our method predicts all the required MVs simultaneously in a single forward pass. The  $\Delta T$  values of our approach as reported in Table III shows that our approach is faster than the baseline encoder on all the test sequences, achieving an average speedup of 10.2% and 9.4% on the 1080p and 720p test sequences, respectively. The corresponding average  $\Delta T$  values obtained using ARPS are 8.4% and 7.6% as also reported in Table III. As the total encoding time at all QP values is much larger than the time it takes to estimate MVs using either method, we get comparable speedups using the MVs derived from CBT-Net and ARPS, with our approach being slightly faster. Thus, using the composite learned model, we were able to improve the RD performance with respect to both the baseline encoder and ARPS without having a detrimental impact on the encoding speed.

#### E. Comparison with Deep Motion Estimation Frameworks

To the best of our knowledge there is no published work on deep learning based block motion estimation. Hence, to conduct a comparison with existing deep motion estimation schemes, we instead studied deep optical flow estimation methods. The following data-driven dense optical flow estimation models were selected for comparison:

- PWC-Net [43] - a supervised framework
- Unflow [45] - a self-supervised framework.

Since the PWC-Net and Unflow models were designed to estimate dense optical flow, we derived the per block MVs by downsampling the dense optical flow field corresponding to each of the four block sizes using bicubic resampling. The resulting motion fields were used as MVs to encode the videos, following the framework described in Section IV-D. Table IV compares the average RD performance of PWC-Net [43] and Unflow [45] with that of CBT-Net.

By contrast with our method, Table IV shows that RD losses were incurred on average when the dense optical flow fields predicted by PWC-Net [43] and Unflow [45] were used to derive block MVs. This experimental result suggests that the dense optical flow fields predicted by these methods do not provide an efficacious alternative to derive

TABLE IV  
AVERAGE RD PERFORMANCE OF CBT-NET MODEL AGAINST PWC-NET [43] AND UNFLOW [45].

Resolution	MS-SSIM BD-rates (%)			VMAF BD-rates(%)			PSNR BD-rates (%)		
	PWC-Net	Unflow	CBT-Net	PWC-Net	Unflow	CBT-Net	PWC-Net	Unflow	CBT-Net
1080p	0.44	0.28	-2.29	0.66	0.55	-1.90	-0.19	-0.40	-1.70
720p	1.26	1.37	-1.18	1.38	1.39	-0.72	1.13	1.18	-0.97

TABLE V  
COMPARISON OF PWC-NET [43], UNFLOW [45] AND CBT-NET IN TERMS OF COMPUTATIONAL COMPLEXITY.

Model	PWC-Net	Unflow	CBT-Net
# of trainable parameters	9.37M	116.58M	1.91M
# of FLOPS	771.14G	1894.8G	41.29G
Inference time on GPU	0.26s	0.48s	0.0025s
Inference time on CPU	5.56s	8.46s	0.14s

block MVs. Furthermore, the estimation of dense optical flow is generally much more computationally intensive than block motion estimation, which makes it uneconomical to derive block MVs from dense optical flow fields. This is also demonstrated with the help of Table V, that reports the number of trainable parameters, the number of floating point operations per second (FLOPS) as well as the inference times (for 1080p input frames) of the three models being compared. The data presented in Table V reveal that the CBT-Net model has significantly fewer trainable parameters, FLOPs and lower inference times as compared to [43] and [45], which were designed for dense optical flow estimation.

## VI. CONCLUSION

We developed a composite model to collectively estimate block MVs for multiple block sizes using a multi-stage CNN. The resulting CBT-Net model was trained on a database of frame-triplets created from publicly available video sources to support a hierarchical bidirectional inter prediction structure commonly used in hybrid codecs. We trained the CBT-Net model instances using a MS-SSIM based loss function in order to favor the perceptual quality of motion compensated frame predictions over pixelwise prediction errors traditionally used as block matching criteria. The proposed framework was applied to AV1 encoding, where substituting the integer motion search of the SVT-AV1 encoder with the trained CBT-Net model resulted in average BD-rates of -1.73% and -1.31% with respect to MS-SSIM and VMAF, respectively, outperforming the corresponding gains obtained using the ARPS based block matching algorithm for motion estimation. Our approach also attained faster encoding speeds as compared to the HME baseline and ARPS. The RD performance gain establishes the efficacy of our approach in improving the perceptual quality of the encoded videos, which is further substantiated by visual comparisons. As a pertinent future direction, the extension of our current framework to capture more complicated motion types that cannot be accounted for by simple block translations might be considered.

## APPENDIX A COMPUTING THE EFFECTIVE RECEPTIVE FIELD OF CBT-NET

The effective receptive field (ERF) of a CNN is defined as the area of the input image that affects a single feature in the output feature map. The ERF progressively enlarges as the depth of the network increases. The ERF of a layer at depth  $k$  can be computed recursively using the following formula [57]:

$$R_k = R_{k-1} + (f_k - 1) \prod_{i=1}^{k-1} s_i \quad (4)$$

where  $R_k$  is the ERF of layer  $k$ ,  $f_k$  is the size of the convolutional kernel at layer  $k$ , and  $s_i$  is the stride of the  $i^{\text{th}}$  layer.  $R_0 = 1$  at beginning, i.e. at the level of the input image. Applying this formula to the CBT-Net layer by layer, gives us an ERF of 255 at the final feature extraction layer as computed step by step in Table VI.

TABLE VI  
ERFS OF CBT-NET'S FEATURE EXTRACTION LAYERS.

Layer # ( $k$ )	filter size ( $f_k$ )	stride ( $s_k$ )	$\prod_{i=1}^{k-1} s_i$	ERF ( $R_k$ )
1	7	2	-	7
2	5	2	2	15
3	5	2	4	31
4	3	1	8	47
5	3	2	8	63
6	3	1	16	95
7	3	2	16	127
8	3	1	32	191
9	3	2	32	255

## REFERENCES

- [1] Z. Wang and A. C. Bovik, "Mean squared error: Love it or leave it? a new look at signal fidelity measures," *IEEE Signal Process. Mag.*, vol. 26, no. 1, pp. 98–117, Jan. 2009.
- [2] S. Zhu and K.-K. Ma, "A new diamond search algorithm for fast block-matching motion estimation," *IEEE Trans. Image Process.*, vol. 9, no. 2, pp. 287–290, Feb. 2000.
- [3] Y. Nie and K.-K. Ma, "Adaptive rood pattern search for fast block-matching motion estimation," *IEEE Trans. Image Process.*, vol. 11, no. 12, pp. 1442–1449, Dec. 2002.
- [4] Y. Chen *et al.*, "An overview of core coding tools in the AV1 video codec," in *Proc. IEEE Picture Coding Symp.*, 2018, pp. 41–45.
- [5] Y. Chen *et al.*, "An overview of coding tools in AV1: the first video codec from the alliance for open media," *APSIPA Trans. Signal Inf. Process.*, vol. 9, Feb. 2020.
- [6] M. Jaderberg, K. Simonyan, A. Zisserman *et al.*, "Spatial transformer networks," in *Advances Neural Inf. Process. Syst.*, 2015, pp. 2017–2025.
- [7] Z. Wang, E. P. Simoncelli, and A. C. Bovik, "Multiscale structural similarity for image quality assessment," in *Proc. 37th Asilomar Conf. Signals, Syst. Comput.*, vol. 2, 2003, pp. 1398–1402.
- [8] Scalable video technology for AV1 (SVT-AV1 encoder and decoder). Accessed May 09, 2021. [Online]. Available: <https://gitlab.com/AOMediaCodec/SVT-AV1>

- [9] F. Kossentini *et al.*, “The SVT-AV1 encoder: overview, features and speed-quality tradeoffs,” in *Applications Digital Image Process. XLIII*, vol. 11510, 2020, p. 1151021.
- [10] A. Norkin *et al.* SVT-AV1: open-source AV1 encoder and decoder. Accessed May 09, 2021. [Online]. Available: <https://netflixtechblog.com/svt-av1-an-open-source-av1-encoder-and-decoder-ad295d9b5ca2>
- [11] Z. Li, A. Aaron, I. Katsavounidis, A. Moorthy, and M. Manohara. Toward a practical perceptual video quality metric. Accessed: May 27, 2019. [Online]. Available: <https://medium.com/netflix-techblog/toward-a-practical-perceptual-video-quality-metric-653f208b9652>
- [12] Scalable video technology for HEVC encoder (SVT-HEVC Encoder). Accessed Oct. 09, 2022. [Online]. Available: <https://github.com/OpenVisualCloud/SVT-HEVC>
- [13] Scalable video technology for VP9 encoder (svt-vp9 encoder). Accessed Oct. 09, 2022. [Online]. Available: <https://github.com/OpenVisualCloud/SVT-VP9>
- [14] HM reference software for HEVC. Accessed Oct. 09, 2022. [Online]. Available: <https://vcgit.hhi.fraunhofer.de/jvet/HM>
- [15] VTM reference software for VVC. Accessed Oct. 09, 2021. [Online]. Available: [https://vcgit.hhi.fraunhofer.de/jvet/VVCSsoftware\\_VTM](https://vcgit.hhi.fraunhofer.de/jvet/VVCSsoftware_VTM)
- [16] G. Lu, W. Ouyang, D. Xu, X. Zhang, C. Cai, and Z. Gao, “DVC: An end-to-end deep video compression framework,” in *Proc. IEEE Conf. Comput. Vision Pattern Recog.*, 2019, pp. 11 006–11 015.
- [17] O. Rippel, S. Nair, C. Lew, S. Branson, A. G. Anderson, and L. Bourdev, “Learned video compression,” in *Proc. Int. Conf. Comput. Vision*, 2019, pp. 3454–3463.
- [18] A. Djelouah, J. Campos, S. Schaub-Meyer, and C. Schroers, “Neural inter-frame compression for video coding,” in *Proc. Int. Conf. Comput. Vision*, 2019, pp. 6421–6429.
- [19] M. Chen, T. Goodall, A. Patney, and A. C. Bovik, “Learning to compress videos without computing motion,” *arXiv preprint arXiv:2009.14110*, 2020.
- [20] J. Liu, S. Xia, W. Yang, M. Li, and D. Liu, “One-for-all: Grouped variation network-based fractional interpolation in video coding,” *IEEE Trans. Image Process.*, vol. 28, no. 5, pp. 2140–2151, May 2019.
- [21] C. Jia *et al.*, “Content-aware convolutional neural network for in-loop filtering in high efficiency video coding,” *IEEE Trans. Image Process.*, vol. 28, no. 7, pp. 3343 – 3356, Jul. 2019.
- [22] Y. Zhang, T. Shen, X. Ji, Y. Zhang, R. Xiong, and Q. Dai, “Residual highway convolutional neural networks for in-loop filtering in HEVC,” *IEEE Trans. Image Process.*, vol. 27, no. 8, pp. 3827–3841, Mar. 2018.
- [23] M. Xu, T. Li, Z. Wang, X. Deng, R. Yang, and Z. Guan, “Reducing complexity of HEVC: A deep learning approach,” *IEEE Trans. Image Process.*, vol. 27, no. 10, pp. 5044–5059, Jun. 2018.
- [24] S. Paul, A. Norkin, and A. C. Bovik, “Speeding up VP9 intra encoder with hierarchical deep learning-based partition prediction,” *IEEE Trans. Image Processing*, vol. 29, pp. 8134–8148, Jul. 2020.
- [25] I. Schiopu, H. Huang, and A. Munteanu, “CNN-based intra-prediction for lossless HEVC,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 30, no. 7, pp. 1816–1828, Sep. 2019.
- [26] Y. Wang, X. Fan, C. Jia, D. Zhao, and W. Gao, “Neural network based inter prediction for HEVC,” in *Proc. IEEE Int. Conf. Multimedia Expo*, 2018, pp. 1–6.
- [27] Z. Zhao, S. Wang, S. Wang, X. Zhang, S. Ma, and J. Yang, “CNN-based bi-directional motion compensation for high efficiency video coding,” in *Proc. IEEE Int. Symp. Circuits Syst.*, 2018, pp. 1–4.
- [28] S. Huo, D. Liu, F. Wu, and H. Li, “Convolutional neural network-based motion compensation refinement for video coding,” in *Proc. IEEE Int. Symp. Circuits Syst.*, 2018, pp. 1–4.
- [29] J.-S. Lee and T. Ebrahimi, “Perceptual video compression: A survey,” *IEEE J. Sel. Topics Signal Process.*, vol. 6, no. 6, pp. 684–697, Aug. 2012.
- [30] Z. Chen, W. Lin, and K. N. Ngan, “Perceptual video coding: Challenges and approaches,” in *IEEE Int. Conf. Multimedia Expo*, 2010, pp. 784–789.
- [31] S. Ki, S.-H. Bae, M. Kim, and H. Ko, “Learning-based just-noticeable-quantization-distortion modeling for perceptual video coding,” *IEEE Trans. Image Process.*, vol. 27, no. 7, pp. 3178–3193, Mar. 2018.
- [32] S. Paul, A. Norkin, and A. C. Bovik, “On visual masking estimation for adaptive quantization using steerable filters,” *Signal Process.: Image Communication*, vol. 96, p. 116290, Aug. 2021.
- [33] I. Marzuki and D. Sim, “Perceptual adaptive quantization parameter selection using deep convolutional features for HEVC encoder,” *IEEE Access*, vol. 8, pp. 37 052–37 065, Feb. 2020.
- [34] H. Hadizadeh and I. V. Bajić, “Saliency-aware video compression,” *IEEE Trans. Image Processing*, vol. 23, no. 1, pp. 19–33, Sep. 2013.
- [35] S. Zhu and Z. Xu, “Spatiotemporal visual saliency guided perceptual high efficiency video coding with neural network,” *Neurocomputing*, vol. 275, pp. 511–522, Jan. 2018.
- [36] Z. Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, “Image quality assessment: from error visibility to structural similarity,” *IEEE Trans. Image Process.*, vol. 13, no. 4, pp. 600–612, Apr. 2004.
- [37] S. Wang, A. Rehman, Z. Wang, S. Ma, and W. Gao, “Perceptual video coding based on SSIM-inspired divisive normalization,” *IEEE Trans. Image Processing*, vol. 22, no. 4, pp. 1418–1429, Dec. 2012.
- [38] L.-H. Chen, C. G. Bampis, Z. Li, A. Norkin, and A. C. Bovik, “ProxIQ: A proxy approach to perceptual optimization of learned image compression,” *IEEE Trans. Image Process.*, vol. 30, pp. 360–373, Nov. 2020.
- [39] J. Lee, K. Kong, G. Bae, and W.-J. Song, “BlockNet: A deep neural network for block-based motion estimation using representative matching,” *Symmetry*, vol. 12, no. 5, p. 840, May 2020.
- [40] G. Choi, P. Heo, S. R. Oh, and H. Park, “A new motion estimation method for motion-compensated frame interpolation using a convolutional neural network,” in *Proc. IEEE Int. Conf. Image Process.*, 2017, pp. 800–804.
- [41] A. Dosovitskiy *et al.*, “Flownet: Learning optical flow with convolutional networks,” in *Proc. IEEE Int. Conf. Comput. Vision*, 2015, pp. 2758–2766.
- [42] A. Ranjan and M. J. Black, “Optical flow estimation using a spatial pyramid network,” in *Proc. IEEE Conf. Comput. Vision Pattern Recog.*, 2017, pp. 4161–4170.
- [43] D. Sun, X. Yang, M.-Y. Liu, and J. Kautz, “PWC-Net: CNNs for optical flow using pyramid, warping, and cost volume,” in *Proc. IEEE Conf. Comput. Vision Pattern Recog.*, 2018, pp. 8934–8943.
- [44] J. Y. Jason, A. W. Harley, and K. G. Derpanis, “Back to basics: Unsupervised learning of optical flow via brightness constancy and motion smoothness,” in *Proc. Eur. Conf. on Comput. Vision*, 2016, pp. 3–10.
- [45] S. Meister, J. Hur, and S. Roth, “Unflow: Unsupervised learning of optical flow with a bidirectional census loss,” in *Proc. AAAI conf. artificial intelligence*, vol. 32, no. 1, 2018.
- [46] M. Cheon and J.-S. Lee, “Subjective and objective quality assessment of compressed 4K UHD videos for immersive experience,” *IEEE Trans. Circuits Syst. Video Technol.*, vol. 28, no. 7, pp. 1467–1480, Mar. 2017.
- [47] L. Song, X. Tang, W. Zhang, X. Yang, and P. Xia, “The SJTU 4K video sequence dataset,” in *Proc. Fifth Int. Workshop Qual. Multimedia Experience (QoMEX)*, 2013, pp. 34–35.
- [48] A. Mercat, M. Viitanen, and J. Vanne, “UVG dataset: 50/120fps 4K sequences for video codec analysis and development,” in *Proc. ACM Multimedia Syst. Conf.*, 2020, pp. 297–302.
- [49] Netflix open content. Accessed: Mar 27, 2022. [Online]. Available: <https://opencontent.netflix.com/>
- [50] Summary of frame triplet database contents. Accessed: April 24, 2021. [Online]. Available: <https://drive.google.com/file/d/1RUZu58jVAN4eHjQE164n1tr2MC1nrmnn/view?usp=sharing>
- [51] H. Choi and I. V. Bajić, “Affine transformation-based deep frame prediction,” *IEEE Trans. Image Process.*, vol. 30, pp. 3321–3334, Feb. 2021.
- [52] D. P. Kingma and J. Ba, “Adam: A method for stochastic optimization,” *arXiv preprint arXiv:1412.6980*, 2014.
- [53] R. Mariş, R. Brad, and B. Remus, “A comparative study of block matching optical flow algorithms,” *TEM Journal*, vol. 6, no. 4, p. 760, 2017.
- [54] S. Baker, D. Scharstein, J. Lewis, S. Roth, M. J. Black, and R. Szeliski, “A database and evaluation methodology for optical flow,” *Int. J. Comput. Vision*, vol. 92, no. 1, pp. 1–31, Nov. 2010.
- [55] P.910 : Subjective video quality assessment methods for multimedia applications. Accessed: April 24, 2021. [Online]. Available: <https://www.itu.int/rec/T-REC-P.910-200804-I>
- [56] G. Bjøntegaard, “Calculation of average PSNR differences between RD-curves,” *Proc. of the ITU-T Video Coding Experts Group (VCEG) Thirteenth Meeting*, Apr. 2001.
- [57] H. Le and A. Borji, “What are the receptive, effective receptive, and projective fields of neurons in convolutional neural networks?” *arXiv preprint arXiv:1705.07049*, 2017.