# Deep Animation Video Interpolation in the Wild

Li Siyao[1*]     Shiyu Zhao[1,2*]     Weijiang Yu[3]     Wenxiu Sun[1,4]

Dimitris Metaxas[2]     Chen Change Loy[5]     Ziwei Liu[5✉]

[1]SenseTime Research and Tetras.AI     [2]Rutgers University     [3]Sun Yat-sen University
[4]Shanghai AI Laboratory     [5]S-Lab, Nanyang Technological University

lisiyao1@sensetime.com   sz553@rutgers.edu   weijiangyu8@gmail.com
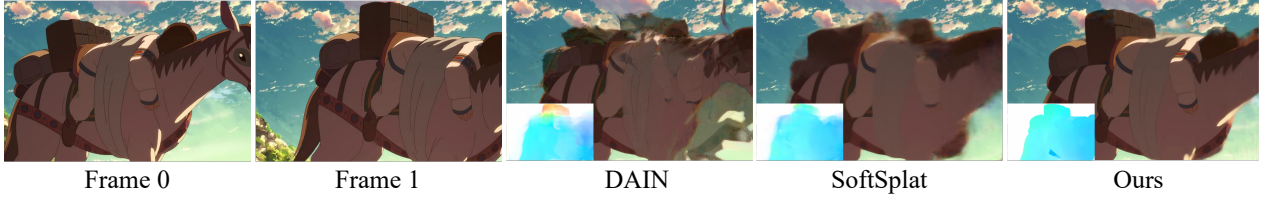sunwenxiu@sensetime.com   dnm@cs.rutgers.edu   {ccloy, ziwei.liu}@ntu.edu.sg

| Frame 0 | Frame 1 | DAIN | SoftSplat | Ours |

Figure 1: **A typical example for animation video interpolation.** Our approach is capable of estimating optical flows for large motions correctly and restore the content, while competing methods fail to handle such motions.

## Abstract

*In the animation industry, cartoon videos are usually produced at low frame rate since hand drawing of such frames is costly and time-consuming. Therefore, it is desirable to develop computational models that can automatically interpolate the in-between animation frames. However, existing video interpolation methods fail to produce satisfying results on animation data. Compared to natural videos, animation videos possess two unique characteristics that make frame interpolation difficult: 1) cartoons comprise lines and smooth color pieces. The smooth areas lack textures and make it difficult to estimate accurate motions on animation videos. 2) cartoons express stories via exaggeration. Some of the motions are non-linear and extremely large. In this work, we formally define and study the animation video interpolation problem for the first time. To address the aforementioned challenges, we propose an effective framework, **AnimeInterp**, with two dedicated modules in a coarse-to-fine manner. Specifically, 1) Segment-Guided Matching resolves the "lack of textures" challenge by exploiting global matching among color pieces that are piece-wise coherent. 2) Recurrent Flow Refinement resolves the "non-linear and extremely large motion" challenge by recurrent predictions using a transformer-like architecture. To facilitate comprehensive training and evaluations, we build a large-scale animation triplet dataset, **ATD-12K**, which comprises 12,000 triplets with rich annotations. Extensive experiments demonstrate that our approach outperforms ex-*

*isting state-of-the-art interpolation methods for animation videos. Notably, AnimeInterp shows favorable perceptual quality and robustness for animation scenarios in the wild. The proposed dataset and code are available at* `https://github.com/lisiyao21/AnimeInterp/`.

## 1. Introduction

In the animation industry, cartoon videos are produced from hand drawings of expert animators using a complex and precise procedure. To draw each frame of an animation video manually would consume tremendous time, thus leading to a prohibitively high cost. In practice, the animation producers usually replicate one drawing two or three times to reduce the cost, which results in the actual low frame rate of animation videos. Therefore, it is highly desirable to develop computational algorithms to interpolate the intermediate animation frames automatically.

In recent years, video interpolation has made great progress on natural videos. However, in animations, existing video interpolation methods are not able to produce satisfying in-between frames. An example from the film *Children Who Chase Lost Voices* is illustrated in Figure 1, where the current state-of-the-art methods fail to generate a piece of complete luggage due to the incorrect motion estimation, which is shown in the lower-left corner of the image. The challenges here stem from the two unique characteristics of animation videos: 1) First, cartoon images consist of explicit sketches and lines, which split the image into segments of smooth color pieces. Pixels in one segment are similar, which yields insufficient textures to match the corresponding

Figure 2: **Two challenges in animation video interpolation.** (a) Piece-wise smooth animations lack of textures. (b) Non-linear and extremely large motions.

pixels between two frames and hence increases the difficulty to predict accurate motions. 2) Second, cartoon animations use exaggerated expressions in pursuit of artistic effects, which result in non-linear and extremely large motions between adjacent frames. Two typical cases are depicted in Figure 2 (a) and (b) which illustrate these challenges respectively. Due to these difficulties mentioned above, video interpolation in animations remains a challenging task.

In this work, we develop an effective and principled novel method for video interpolation in animations. In particular, we propose an effective framework, **AnimeInterp**, to address the two aforementioned challenges. AnimeInterp consists of two dedicated modules: a Segment-Guided Matching (SGM) module and a Recurrent Flow Refinement (RFR) module, which are designed to predict accurate motions for animations in a coarse-to-fine manner. More specifically, the proposed SGM module computes a coarse piece-wise optical flow using global semantic matching among color pieces split by contours. Since the similar pixels belonging to one segment are treated as a whole, SGM can avoid the local minimum caused by mismatching on smooth areas, which resolves the "lack of textures" problem. To tackle the "non-linear and extremely large motion" challenge in animation, the piece-wise flow estimated by SGM is further enhanced by a Transformer-like network named Recurrent Flow Refinement. As shown in Figure 1, our proposed approach can better estimate the flow of the luggage in large displacement and generate a complete in-between frame.

A large-scale animation triplet dataset, **ATD-12K**, is built to facilitate comprehensive training and evaluations of video interpolation methods on cartoon videos. Unlike other animation datasets, which consists of only single images, ATD-12K contains 12,000 frame triplets selected from 30 animation movies in different styles with a total length over 25 hours. Apart from the diversity, our test set is divided into three difficulty levels according to the magnitude of the motion and occlusion. We also provide annotations on movement categories for further analysis.

The contributions of this work can be summarized as follows: **1)** We formally define and study the animation video

interpolation problem for the first time. This problem has significance to both academia and industry. **2)** We propose an effective animation interpolation framework named AnimeInterp with two dedicated modules to resolve the "lack of textures" and "non-linear and extremely large motion" challenges. Extensive experiments demonstrate that AnimeInterp outperforms existing state-of-the-art methods both quantitatively and qualitatively. **3)** We build a large-scale cartoon triplet dataset called ATD-12K with large content diversity representing many types of animations to test animation video interpolation methods. Given its data size and rich annotations, ATD-12K would pave the way for future research in animation study.

## 2. Related Work

**Video Interpolation.** Video interpolation has been widely studied in recent years. In [16], Meyer *et al.* propose a phase-based video interpolation scheme, which performs impressively on videos with small displacements. In [19, 20], Niklaus *et al.* design a kernel-based framework to sample interpolated pixels via convolutions on corresponding patches of adjacent frames. However, the kernel-based framework is still not capable of processing large movements due to the limitation of the kernel size. To tackle the large motions in videos, many studies use optical flows for video interpolation. Liu *et al.* [15] predict a 3D voxel flow to sample the inputs into the middle frame. Similarly, Jiang *et al.* [9] propose to jointly estimate bidirectional flows and an occlusion mask for multi-frame interpolation. Besides, some studies are dedicated to improving warping and synthesis with given bidirectional flows [2, 1, 17, 18], and higher-order motion information is used to approximate real-world videos [33, 5]. In addition to using pixel-wise flows on images, "feature flows" on deep features are also explored for video interpolation [8]. Although existing methods achieve great success in interpolating real-world videos, they fail to handle the large and non-linear motions of animations. Thus, the animation video interpolation still remains unsolved. In this paper, we propose a segment-guided matching module based on the color piece matching, boosting the flow prediction.

**Vision for Animation.** In vision and graphics community, there are many works on processing and enhancing the animation contents, *e.g.*, manga sketch simplification [24, 23], vectorization of cartoon images [35, 34], stereoscopy of animation videos [14], and colorization of black-and-white manga [22, 27, 11]. In recent years, with the surge of deep learning, researchers attempt to generate favorable animation contents. For example, generative adversarial networks are used to generate vivid cartoon characters [10, 32] and style-transfer models are developed to transfer real-world images to cartoons [4, 29, 3]. However, generating animation video contents is still a challenging task due to the difficulty of estimating the motion between frames.
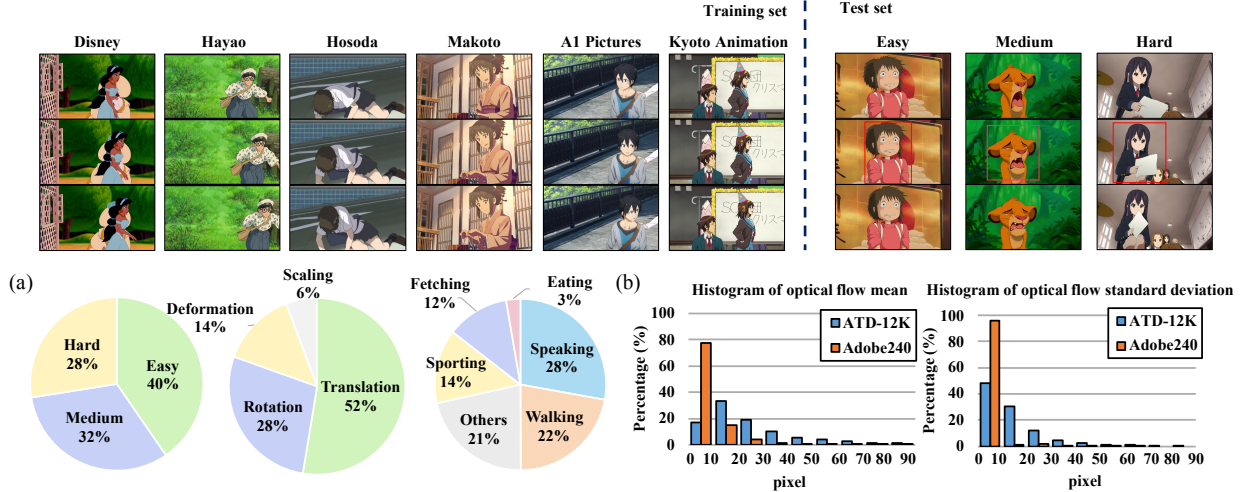
Figure 3: **Triplet samples and statistics of ATD-12K.** Top: typical triplets in the training and test sets. (a) The percentage of difficulty levels and motion tags in different categories. (b) Histograms of mean motion values and standard deviations calculated in each image.

## 3. ATD-12K Dataset

To facilitate the training and evaluation of animation video interpolation methods, we build a large-scale dataset named ATD-12K, which comprises a training set containing 10,000 animation frame triplets and a test set containing 2,000 triplets, collected from a variety of cartoon movies. To keep the objectivity and fairness, test data are sampled from different sources that are exclusive from the training set. Besides, to evaluate video interpolation methods from multiple aspects, we provide rich annotations to the test set, including levels of difficulty, the Regions of Interest (RoIs) on movements, and tags on motion categories. Some typical examples of ATD-12K and the annotation information are displayed in Figure 3.

### 3.1. Dataset Construction

In the first step to build ATD-12K, we download a large number of animation videos from the Internet. To keep the diversity and representativeness of ATD-12K, a total of 101 cartoon clips are collected from 30 series of movies made by diversified producers, including Disney, Hayao, Makoto, Hosoda, Kyoto Animation, and A1 pictures, with a total duration of more than 25 hours. The collected videos have a high visual resolution of $1920 \times 1080$ or $1280 \times 720$. Second, we extract sufficient triplets from the collected videos. In this step, we not only sample adjacent consecutive frames into the triplet but also extract those with one or two frames as an interval to enlarge the inter-frame displacement. Since the drawing of an animation frame may be duplicated several times, it is common to sample similar contents in the triplets extracted from one video. To avoid high affinity among our data, we filter out the triplets that contains two frames with a structural similarity (SSIM) [30] larger than 0.95.

Meanwhile, triplets with SSIM lower than 0.75 are also dropped off to get rid of scene transitions in our data. After that, we manually inspect the remaining 131,606 triplets for further quality improvement. Triplets that are marked as eligible by at least two individuals are kept. Frames that contain inconsistent captions, simple or similar scenes or untypical movements are removed. After the strict selection, 12,000 representative animation triplets are used to construct our final ATD-12K.

**Statistics.** To explore the gap between natural and cartoon videos, we compare the motion statistics of the ATD-12K and a real-world dataset Adobe240 [26], which is often used as an evaluation set for video interpolation. We estimate the optical flow between frames of every input pair in the two datasets, and compute the mean and standard deviation of displacements for each flow. The normalized histograms of means and the standard deviations in the two datasets are shown in Figure 3. They reveal that the cartoon dataset ATD-12K has a higher percentage of large and diverse movements than the real-world Adobe240 dataset.

### 3.2. Annotation

**Difficulty Levels.** We divided the test set of ATD-12K into three difficulty levels, *i.e.*, "Easy", "Medium", and "Hard" based on the average magnitude of motions and the area of occlusion in each triplet. For more details on the definition of each level, please refer to the supplementary file. Sample triplets of different levels are displayed in Figure 3.

**Motion RoIs.** Animation videos are made of moving foreground objects and motion-less background scenes. Since the foreground characters are more visually appealing, the motions of those regions have more impact on audiences' visual experience. To better reveal the performance of interpolation methods on motion areas, we provide a bounding
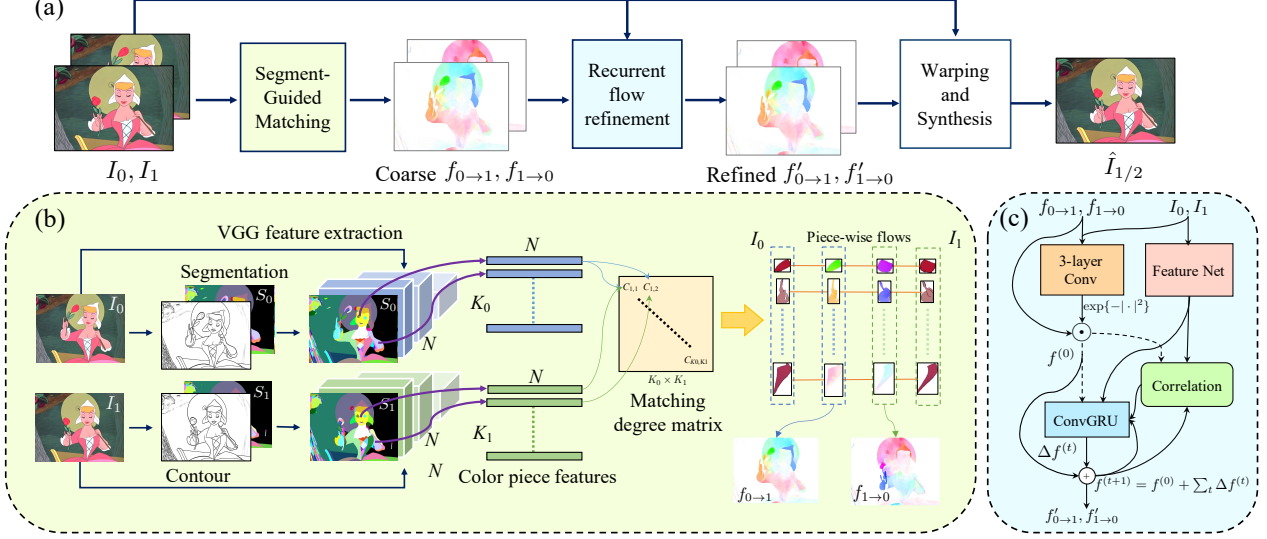
Figure 4: **(a) The overall pipeline of AnimeInterp.** The input $I_0$ and $I_1$ are fed to the SGM module to generate the coarse flows, *i.e.* $f_{0\to1}$ and $f_{1\to0}$. Then, $f_{0\to1}$ and $f_{1\to0}$ are refined by the RFR module. The final interpolation result is produced by the warping and synthesis network borrowed from SoftSplat [18]. **(b) The inner structure of the SGM module. (c) The workflow of the RFR module.**

box for the $2^{nd}$ image (as shown in Figure 3 test set) of each triplet to delineate the motion RoI.

**Motion Tags.** We also provide tags to describe the major motion of a triplet. Our motion tags can be classified into two categories, namely, 1) general motion types including translation, rotation, scaling, and deformation; 2) character behaviors containing speaking, walking, eating, sporting, fetching, and others. The percentage for tags in each category is shown in Figure 3.

## 4. Our Approach

The overview of our framework is shown in Figure 4. Given the input images $I_0$ and $I_1$, we first estimate the coarse flows $f_{0\to1}$ and $f_{1\to0}$ between $I_0$ and $I_1$ in both directions via the proposed SGM module in Section 4.1. Then, we set the coarse flow as the initialization of a recurrent neural network and gradually refine it to obtain the fine flows $f'_{0\to1}$ and $f'_{1\to0}$ in Section 4.2. Finally, we warp $I_0$ and $I_1$ using $f'_{0\to1}$ and $f'_{1\to0}$ and synthesize the final output $\hat{I}_{1/2}$ in Section 4.3.

### 4.1. Segment-Guided Matching

For typical 2D animations, objects are usually outlined with explicit lines, and each enclosed area is filled with a single color. The color of the moving object in one frame remains stable in the next frame despite the large displacement, which could be regarded as a strong clue to find appropriate semantic matching for the color pieces. In this paper, we leverage this clue to tackle smooth piece-wise motions by generating coarse optical flows. The procedure is illustrated in Figure 4(b), and we elaborate it in the following.

**Color Piece Segmentation.** Following Zhang *et al.*'s work [35], we adopt the Laplacian filter to extract contours of animation frames. Then, the contours are filled by the "trapped-ball" algorithm [35] to generate color pieces. After this step, we obtain a segmentation map $S \in \mathbb{N}^{H \times W}$, where pixels of each color piece is labeled by an identity number. In the rest of this section, we notate the segmentation of the input $I_0$ and $I_1$ as $S_0$ and $S_1$, respectively. $S_0(i)$ represents the pixels in the $i^{th}$ color piece of $I_0$, and $S_1(i)$ for $I_1$ is similar.

**Feature Collection.** We feed the input $I_0$ and $I_1$ into a pretrained VGG-19 model [25] and extract features of `relu1_2`, `relu2_2`, `relu3_4` and `relu4_4` layers. Then, we assemble the features belonging to one segment by the super-pixel pooling proposed in [13]. Features of smaller scales are pooled by the down-sampled segmentation maps and are concatenated together. After the pooling, each color piece is represented by an $N$-dimensional feature vector, and the whole image is reflected to a $K \times N$ feature matrix, where $K$ is the number of the color pieces and each row of the matrix represents the feature of the corresponding color piece. The feature matrices $I_0$ and $I_1$ are denoted as $F_0$ and $F_1$, respectively.

**Color Piece Matching.** We now use $F_0$ and $F_1$ to estimate a consistent mapping between color pieces of $I_0$ and $I_1$. Specifically, we predict a forward map $\mathcal{M}_{0\to1}$ and a backward map $\mathcal{M}_{1\to0}$. For the $i^{th}$ color piece in the first frame, the forward map $\mathcal{M}_{0\to1}(i)$ indicates the maximum-likelihood corresponding piece in the second frame and the backward map does the same from $I_1$ to $I_0$. To quantify the likelihood of a candidate pair $(i \in S_0, j \in S_1)$, we compute an affinity metric $\mathcal{A}$ using $F_0$ and $F_1$ as

$$\mathcal{A}(i,j) = \sum_n^N \min\left(\tilde{F}_0(i,n), \tilde{F}_1(j,n)\right), \quad (1)$$

where $\tilde{F}_0(i,n) = F_0(i,n)/\sum_n F_0(i,n)$ is the normalized feature of $F_0$, and $\tilde{F}_1$ is similar to $\tilde{F}_0$. This affinity metric measures the similarities of all pairs globally. Besides, to avoid potential outliers, we also exploit two constraint penalties, *i.e.*, the *distance penalty* and the *size penalty*. First, we assume that the displacement between two corresponding pieces is unlikely to be overly large. The *distance penalty* is defined as the ratio of the distance between the centroids of two color pieces and the diagonal length of the image,

$$\mathcal{L}_{dist}(i,j) = \frac{\|P_0(i) - P_1(j)\|}{\sqrt{H^2 + W^2}}, \quad (2)$$

where $P_0(i)$ and $P_1(j)$ represent the positions of the centroids of $S_0(i)$ and $S_1(j)$, respectively. Note that this penalty is only performed to the matching with the displacement larger than 15% of the diagonal length of the image. Second, we suggest that the sizes of matched pieces should be similar. The *size penalty* is designed as,

$$\mathcal{L}_{size}(i,j) = \left| \frac{|S_0(i)| - |S_1(j)|}{HW} \right|, \quad (3)$$

where $|\cdot|$ denotes the number of pixels in a cluster.

Combining all items above, a matching degree matrix $\mathcal{C}$ is then computed as,

$$\mathcal{C} = \mathcal{A} - \lambda_{dist}\mathcal{L}_{dist} - \lambda_{size}\mathcal{L}_{size}, \quad (4)$$

where $\lambda_{dist}$ and $\lambda_{size}$ are set to 0.2 and 0.05 in our implementation, respectively. For each pair ($i \in S_0, j \in S_1$), $\mathcal{C}(i,j)$ indicates the likelihood. Hence, for the $i^{th}$ color piece in $S_0$, the most likely matching piece in $S_1$ is the one with the maximum matching degree, and vice versa. The mapping of this matching is formulated as,

$$\begin{aligned}
\mathcal{M}_{0\rightarrow1}(i) &= \arg\max_j(\mathcal{C}(i,j)), \\
\mathcal{M}_{1\rightarrow0}(j) &= \arg\max_i(\mathcal{C}(i,j)).
\end{aligned} \quad (5)$$

**Flow Generation.** In the last step of SGM module, we predict dense bidirectional optical flows $f_{0\rightarrow1}$ and $f_{1\rightarrow0}$ using $\mathcal{M}_{0\rightarrow1}$ and $\mathcal{M}_{1\rightarrow0}$. We only describe the procedure to compute $f_{0\rightarrow1}$, for $f_{1\rightarrow0}$ can be obtained by reversing the mapping order. For each matched color-piece pair $(i,j)$ where $j = \mathcal{M}_{0\rightarrow1}(i)$, we first compute the displacement of the centroids as a shift base $f_c^i = P_1(j) - P_0(i)$, and then we compute the local deformation $f_d^i := (u,v)$ by variational optimization,

$$\begin{aligned}
E\left(f_d^i(\mathbf{x})\right) = &\int |I_1^j\left(\mathbf{x} + f_c^i(\boldsymbol{x}) + f_d^i(\mathbf{x})\right) - I_0^i(\mathbf{x})|d\mathbf{x} \\
&+ \int \left(|\nabla u(\mathbf{x})|^2 + \nabla|v(\mathbf{x})|^2\right)|d\mathbf{x}.
\end{aligned}$$
$$(6)$$

Here, $I_0^i$ represents a masked $I_0$ where pixels not belonging to $i^{th}$ color piece are set to 0. $I_1^j$ is the similar to $I_0^i$. The optical flow for pixels in $i^{th}$ color piece is then $f_{0\rightarrow1}^i = f_d^i + f_c^i$. Since the color pieces are disjoint, the final flow for the whole image is calculated by simply adding all piece-wise flows together as $f_{0\rightarrow1} = \sum_i f_{0\rightarrow1}^i$.

To further avoid outliers, we mask the flow of the $i^{th}$ piece to zero if it does not satisfy the mutual consistency, *i.e.*, $\mathcal{M}_{1\rightarrow0}(\mathcal{M}_{0\rightarrow1}(i)) \neq i$ . This operation prevents the subsequent flow refinement step to be misled by the low-confidence matching.

### 4.2. Recurrent Flow Refinement Network

In this section, we refine the coarse optical flows $f_{0\rightarrow1}$ and $f_{1\rightarrow0}$ to the finer views $f'_{0\rightarrow1}$ and $f'_{1\rightarrow0}$ via a deep Recurrent Flow Refinement (RFR) network. There are two main motivations for introducing the RFR module. First, since a strict mutual consistency constraint is adopted in the color-piece-matching step, non-robust pairs are masked off, leaving null flow values in some locations. RFR is able to produce valid flows for those locations. Second, the SGM module is beneficial for large displacements but is less effective to predict precise deformation for the non-linear and exaggerated motion in animation videos. Here, RFR complements the previous step by refining the coarse and piece-wise flows.

Inspired by the state-of-the-art optical flow network RAFT [28], we design a transformer-like architecture as shown in Figure 4(c) to recurrently refine the piece-wise flow $f_{0\rightarrow1}$. For the sake of brevity, we only illustrate the process to compute $f'_{0\rightarrow1}$. First, to further avoid potential errors in the coarse flow $f_{0\rightarrow1}$, a pixel-wise confidence map $g$ is learned to mask the unreliable values via a three-layer CNN, which takes concatenated $|I_1(\mathbf{x} + f_{0\rightarrow1}(\mathbf{x})) - I_0(\mathbf{x})|$, $I_0$ and $f_{0\rightarrow1}$ as input. Then, the coarse flow $f_{0\rightarrow1}$ is multiplied with $\exp\{-g^2\}$ to be the initialization $f_{0\rightarrow1}^{(0)}$ for the following refinement. Next, a series of residues $\{\Delta f_{0\rightarrow1}^{(t)}\}$ are learnt via a convolutional GRU [6]:

$$\Delta f_{0\rightarrow1}^{(t)} = \text{ConvGRU}\left(f_{0\rightarrow1}^{(t)}, \mathcal{F}_0, corr(\mathcal{F}_0, \mathcal{F}_{1\rightarrow0}^{(t)})\right), \quad (7)$$

where $corr(\cdot,\cdot)$ calculates the correlation between two tensors, $\mathcal{F}_0$ and $\mathcal{F}_1$ are frame features extracted by a CNN, and $\mathcal{F}_{1\rightarrow0}^{(t)}$ is bilinearly sampled from $\mathcal{F}_1$ with optical flow $f_{0\rightarrow1}^{(t)}$. The learnt residues are recurrently accumulated to update the initialization. The optical flow refined after $T$ iterations is computed as

$$f_{0\rightarrow1}^{(T)} = f_{0\rightarrow1}^{(0)} + \sum_{t=0}^{T-1} \Delta f_{0\rightarrow1}^{(t)}. \quad (8)$$

And the fine flow $f'_{0\rightarrow1}$ is the output of the last iteration.

Table 1: **Quantitative results on the test set of ATD-12K.** The best and runner-up values are bold and underlined, respectively.

| Method | Whole | | RoI | | Easy | | Medium | | Hard | |
|---|---|---|---|---|---|---|---|---|---|---|
| | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM | PSNR | SSIM |
| Super SloMo w/o. ft. | 27.88 | 0.946 | 24.56 | 0.886 | 30.66 | 0.966 | 27.29 | 0.948 | 24.63 | 0.917 |
| Super SloMo [9] | 28.19 | 0.949 | 24.83 | 0.892 | 30.86 | 0.967 | 27.63 | 0.950 | 25.02 | 0.922 |
| DAIN w/o. ft. | 28.84 | 0.953 | 25.43 | 0.897 | 31.40 | 0.970 | 28.38 | 0.955 | 25.77 | 0.927 |
| DAIN [1] | 29.19 | 0.956 | 25.78 | 0.902 | 31.67 | **0.971** | 28.74 | 0.957 | 26.22 | 0.932 |
| QVI w/o. ft. | 28.80 | 0.953 | 25.54 | 0.900 | 31.14 | 0.969 | 28.44 | 0.955 | 25.93 | 0.929 |
| QVI [33] | 29.04 | 0.955 | 25.65 | 0.901 | 31.46 | 0.970 | 28.63 | 0.956 | 26.11 | 0.931 |
| AdaCoF w/o. ft. | 28.10 | 0.947 | 24.72 | 0.886 | 31.09 | 0.968 | 27.43 | 0.948 | 24.65 | 0.916 |
| AdaCoF [12] | 28.29 | 0.951 | 24.89 | 0.894 | 31.10 | 0.969 | 27.63 | 0.951 | 25.10 | 0.925 |
| SoftSplat w/o. ft. | 29.15 | 0.955 | 25.75 | 0.904 | 31.50 | 0.970 | 28.75 | 0.956 | 26.29 | 0.934 |
| SoftSplat [18] | 29.34 | 0.957 | 25.95 | 0.907 | 31.60 | 0.970 | 28.96 | 0.958 | 26.59 | 0.938 |
| Ours w/o. SGM | 29.54 | **0.958** | 26.15 | **0.910** | 31.80 | **0.971** | 29.15 | **0.959** | 26.78 | **0.939** |
| Ours w/o. RFR | 27.62 | 0.944 | 24.43 | 0.887 | 29.78 | 0.959 | 27.29 | 0.946 | 24.94 | 0.920 |
| Ours | **29.68** | **0.958** | **26.27** | **0.910** | **31.86** | **0.971** | **29.26** | **0.959** | **27.07** | **0.939** |

## 4.3. Frame Warping and Synthesis

To generate the intermediate frame using flow $f'_{0\to1}$ and $f'_{1\to0}$, we adopt the splatting and synthesis strategy of Soft-Splat [18]. In short, a bunch of features are extracted from $I_0$ and $I_1$ using a multi-scale CNN, and then, all features and input frames are splatted via forward warping to the center position, *e.g.*, $I_0$ is splatted to $I_{0\to1/2}$ as

$$I_{0\to1/2}(\mathbf{x} + f_{0\to1}(\mathbf{x})/2) = I_0(\mathbf{x}). \qquad (9)$$

Finally, all warped frames and features are fed into a GridNet [7] with three scale levels to synthesize the target frame $\hat{I}_{1/2}$.

## 4.4. Learning

To supervise the proposed network, we adopt the $L_1$ distance $\|\hat{I}_{1/2} - I_{1/2}\|_1$ between the prediction $\hat{I}_{1/2}$ and the ground truth $I_{1/2}$ as the training loss.

We train this network in two phases: the *training phase* and the *fine-tuning phase*. In the *training phase*, we first pre-train the recurrent flow refinement (RFR) network following [28], and then fix the weights of RFR to train the rest parts of proposed network on a real-world dataset proposed in [33] for 200 epochs. During this phase, we do not use the coarse flows predicted by SGM module. The learning rate is initialized as $10^{-4}$ and decreases with a factor $0.1$ at the $100^{th}$ and the $150^{th}$ epochs. In the second phase, we fine-tune the whole system for another 50 epochs on the training set of our ATD-12K with the learning rate of $10^{-6}$. During fine-tuning, images are rescaled into $960 \times 540$ and are randomly cropped into $380 \times 380$ with batch size 16. We also stochastically flip the images and reverse the triplet order as data augmentation.

## 5. Experiments

**Compared Methods.** AnimeInterp is compared with five most recent state-of-the-art methods including Super SloMo [9], DAIN [1], QVI [33], AdaCoF [12] and SoftSplat [18]. We use the original implementation of DAIN, QVI and Ada-CoF, and the implementation from [21] for Super SloMo. Since no model is released for SoftSplat, we implement this method following the description in [18], and train it with the same strategy as our proposed model.

**Datasets.** We fine-tune these baseline models on the training set of ATD-12K with the hyper-parameters described in Section 4.4. Models before and after fine-tuning are separately evaluated on our proposed ATD-12K test set.

**Evaluation Metrics.** We use two distant images in every test triplet as input frames and use the middle one as ground truth. The PSNR and SSIM [31] between the predicted intermediate results and the ground truth are adopted as evaluation metrics. Using the annotation in Section 3.2, we evaluate the interpolation results not only on the whole image (denoted as "Whole"), but also on the regions of interest (denoted as "RoI"). Meanwhile, we show the evaluations for the three levels ("Easy", "Medium" and "Hard").

### 5.1. Comparative Results

**Quantitative Evaluation.** The quantitative results are shown in Table 1. According to the comparison, our proposed model consistently performs favorably against all the other existing methods on all evaluations. The PSNR score of our method improves 0.34dB compared to the best comparative method (SoftSplat) on the whole image evaluation. For the evaluation on RoI, which indicates the significant movement in a triplet, our method can also achieve an improvement of 0.32dB. Since animation videos contain many frames with static background, the evaluation on RoI is more precise to reflect the effectiveness on large motions.

**Qualitative Evaluation.** To further demonstrate the effectiveness of our proposed method, we show two examples for visual comparison in Figure 5. In each example, we display
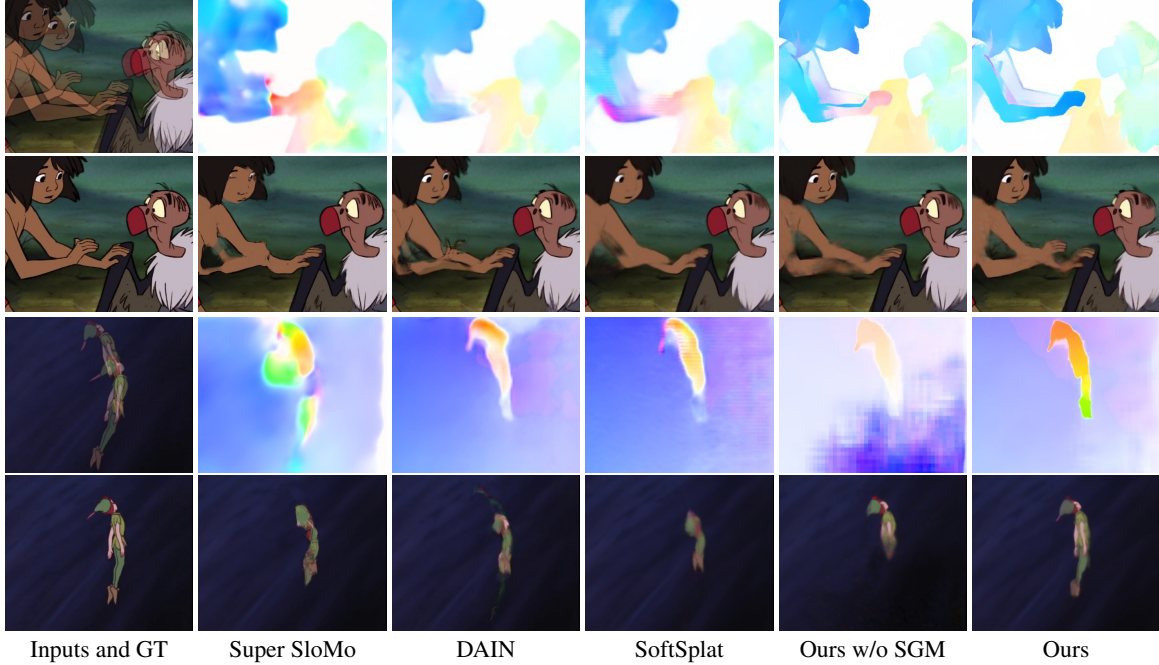
Figure 5: **Qualitative results on ATD-12K test set.** The first row of each example shows the overlapped input frames and the optical flow predicted by corresponding method, while the second row shows the ground truth and interpolation results.

the interpolation results and corresponding optical flows predicted by different methods. In the first example, *Mowgli* is reaching out a vulture using one hand, while his other hand is already on the wing of the vulture. Interpolation on this case is very challenging because the displacement of the boy's arm is locally discontinuous and extremely large. Moreover, the local patterns of the boy's two hands are very similar, which result in a local minimum for the optical flow prediction. Consequently, all existing methods fail to predict correct flow values on the boy's hand. In the interpolation results of these methods, the moving hand either disappears (*e.g.*, SoftSplat) or splits into two deformed ones (*e.g.*, Super SloMo, DAIN). In contrast, our method can estimate a more accurate optical flow and generate a relatively complete arm in the intermediate position. In the second case, *Peter Pan* flies fast in the air while the background moves with a large displacement. Similar to the first example, the compared methods estimate wrong flow values on the character, and hence fail to synthesize a complete body in the accurate position, while our method can produce the whole character, which is almost identical to the ground truth and looks correct.

## 5.2. Ablation Study

**Quantitative Evaluation.** To explain the effectiveness of the SGM and RFR modules in AnimeInterp, we evaluate two variants of the proposed method, where the SGM module and the RFR module are removed, respectively (denote

as "Ours w/o. SGM" and "Ours w/o. RFR"). The "w/o. SGM" variant directly predicts optical flows using the recurrent refinement network, without the guided initialization of global-context-aware piece-wise flow, while the "w/o. RFR" model uses the piece-wise flow predicted by SGM to warp and synthesize the interpolate results. For a fair comparison, we re-fine-tune the weights of these two variants after removing the corresponding modules. As shown in Table 1 and Figure 6(a), the PSNR of " w/o. SGM" variant on the whole-image evaluation is lowered about 0.14dB , while the score of "w/o. RFR" model drops sharply to 2.06dB, because the piece-wise flow are coarse and contain zero values on mutually inconsistent color pieces. The results suggest that both of the proposed SGM module and RFR module play a critical role in our method. To have a more precise view on different difficulty levels, the SGM module can improve about 0.3dB on the performance of triplets marked as "Hard", from 26.78dB to 27.07dB, but merely improves the "Easy" evaluation by 0.06dB, suggesting that the proposed SGM module is more useful on improving the performance of large-motion videos.

**Qualitative Evaluation.** As for the quality of visual results, if the SGM module is removed, the predicted flow of our methods becomes incorrect since it falls into a local minimum. For instance, in the first sample of Figure 5, the moving right hand of the boy in the first image is matched to the left hand in the second image, which is spatially closer but wrong for the matching. Interestingly, with the piece-
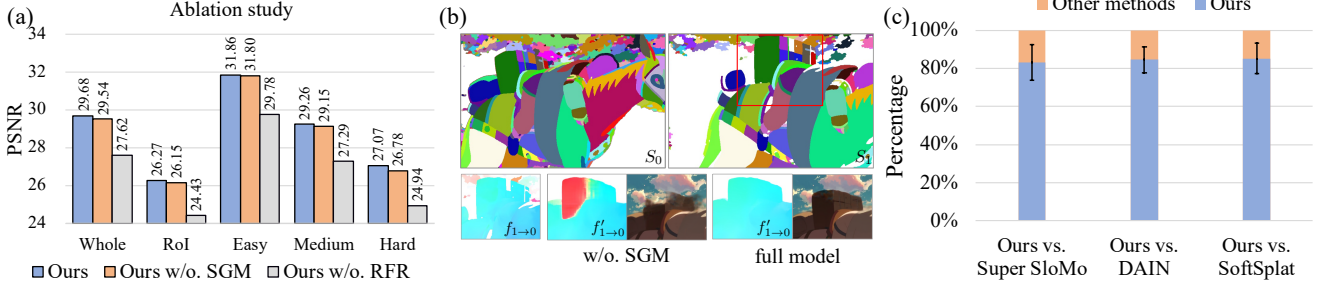
Figure 6: (a) **Ablation study for the SGM module.** Quantitative results of PSNR are exhibited. (b) **Visualization for the effectiveness of SGM.** In the first row, the color piece segmentation of frame 0 and frame 1 in Figure 1 are visualized, where the matching pieces are filled with the same random color in $S_0$ and $S_1$. In the second row, flow estimations and final interpolation results are displayed. (c) **User study results.** Our method outperforms others with a large margin.

wise flow as guidance, our method can avoid this local minimum and estimate correct flow values by taking advantage of global context information. A detailed illustration of the effectiveness of SGM is shown in Figure 6(b). In the first row, we show the color piece segmentation of the two input frames displayed in Figure 1. The matching pieces are filled with the same random color, and the color pieces disobeying the mutual consistency are masked as white. As can be seen in the red boxes, the segments of the luggage are correctly matched, even though it is segmented into several pieces. In the second row of Figure 6(b), we visually compare the optical flow and the interpolation results with and without SGM. Although the piece-wise flow $f_{1\to0}$ is coarse, it can be a good guidance to the refinement network to predict accurate $f'_{1\to0}$, which leads to a complete synthesis of the luggage. In conclusion, the proposed SGM module significantly improves the performance on local-discontinuous motions via global context matching.

### 5.3. Further Analysis

**Influence of Difficulty Levels and Motion RoIs.** As can be seen in Table 1, the difficulty levels have great impact on the performance. The PSNR scores of various methods will drop more than 2dB when the difficulty increases one level. The drop is more significant for those methods which are not designed to resolve "lack of textures" and "non-linear and extremely large motion" challenges in animation. For example, DAIN achieves the third-best score (31.67dB) in the evaluation on 'Easy' level, with the difference less than 0.2dB compared to our proposed method (31.86dB), but its performance drops more than ours for 'Medium' (28.74dB vs. 29.26dB) and 'Hard' (26.22dB vs. 27.07dB) levels. Meanwhile, the performance of various methods on RoIs are much lower than that on the whole images. Since RoI regions are the essential parts of animation videos, which have great impact on the visual experience, future studies focusing on restoring RoIs are encouraged.

**User Study.** To further evaluate the visual performance of our methods, we conduct a user study on the interpolation results of Super SloMo, DAIN, SoftSplat and our proposed method. We separately conduct subjective experiments with ten people. In every test, we randomly select 150 pairs, each of which contains the interpolation result of our method and the corresponding frame from one of the compared methods, and ask the subject to vote for the better one. We provide both the interpolated frames and the input frames to subjects so that they can take the temporal consistency in their decision. Figure 6 shows the percentages of the average votes of our method versus the compared methods. For each compared method, over 83% of the average votes account for ours as being better in terms of visual quality. These experimental results show our method produces more favorable results on both the quantitative evaluations and the visual quality.

## 6. Conclusion

This paper introduces the animation video interpolation task, which is an important step towards cartoon video generation. To benchmark animation video interpolation, we build a large-scale animation triplet dataset named ATD-12K comprising 12,000 triplets with rich annotations. We also propose an effective animation video interpolation framework called AnimeInterp, which employs the segmentation of color pieces as a match guidance to compute piece-wise optical flows between different frames and leverages a recurrent module to further refine the quality of optical flows. Comprehensive experiments demonstrate the effectiveness of the proposed modules as well as the generalization ability of AnimeInterp.

# References

[1] Wenbo Bao, Wei-Sheng Lai, Chao Ma, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. Depth-aware video frame interpolation. In *CVPR*, pages 3703–3712, 2019.

[2] Wenbo Bao, Wei-Sheng Lai, Xiaoyun Zhang, Zhiyong Gao, and Ming-Hsuan Yang. MEMC-Net: Motion estimation and motion compensation driven neural network for video interpolation and enhancement. *IEEE TPAMI*, 2019.

[3] Jie Chen, Gang Liu, and Xin Chen. AnimeGAN: A novel lightweight gan for photo animation. In *ISICA*, pages 242–256, 2019.

[4] Yang Chen, Yu-Kun Lai, and Yong-Jin Liu. CartoonGAN: Generative adversarial networks for photo cartoonization. In *CVPR*, pages 9465–9474, 2018.

[5] Zhixiang Chi, Rasoul Mohammadi Nasiri, Zheng Liu, Juwei Lu, Jin Tang, and Konstantinos N. Plataniotis. All at once: Temporally adaptive multi-frame interpolation with advanced motion modeling. In *ECCV*, 2020.

[6] Kyunghyun Cho, Bart Van Merriënboer, Dzmitry Bahdanau, and Yoshua Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.

[7] Damien Fourure, Rémi Emonet, Elisa Fromont, Damien Muselet, Alain Tremeau, and Christian Wolf. Residual conv-deconv grid network for semantic segmentation. In *BMVC*, pages 181.1–181.13, 2017.

[8] Shurui Gui, Chaoyue Wang, Qihua Chen, and Dacheng Tao. FeatureFlow: Robust video interpolation via structure-to-texture generation. In *CVPR*, pages 14004–14013, 2020.

[9] Huaizu Jiang, Deqing Sun, Varun Jampani, Ming-Hsuan Yang, Erik Learned-Miller, and Jan Kautz. Super SloMo: High quality estimation of multiple intermediate frames for video interpolation. In *CVPR*, pages 9000–9008, 2018.

[10] Yanghua Jin, Jiakai Zhang, Minjun Li, Yingtao Tian, Huachun Zhu, and Zhihao Fang. Towards the automatic anime characters creation with generative adversarial networks. *arXiv preprint arXiv:1708.05509*, 2017.

[11] Johannes Kopf and Dani Lischinski. Digital reconstruction of halftoned color comics. *ACM TOG*, 31(6):1–10, 2012.

[12] Hyeongmin Lee, Taeoh Kim, Tae-young Chung, Daehyun Pak, Yuseok Ban, and Sangyoun Lee. AdaCoF: Adaptive collaboration of flows for video frame interpolation. In *CVPR*, pages 5316–5325, 2020.

[13] Fayao Liu, Chunhua Shen, Guosheng Lin, and Ian Reid. Learning depth from single monocular images using deep convolutional neural fields. *IEEE TPAMI*, 38(10):2024–2039, 2015.

[14] Xueting Liu, Xiangyu Mao, Xuan Yang, Linling Zhang, and Tien-Tsin Wong. Stereoscopizing cel animations. *ACM TOG*, 32(6):1–10, 2013.

[15] Ziwei Liu, Raymond A Yeh, Xiaoou Tang, Yiming Liu, and Aseem Agarwala. Video frame synthesis using deep voxel flow. In *CVPR*, pages 4463–4471, 2017.

[16] Simone Meyer, Oliver Wang, Henning Zimmer, Max Grosse, and Alexander Sorkine-Hornung. Phase-based frame interpolation for video. In *CVPR*, pages 1410–1418, 2015.

[17] Simon Niklaus and Feng Liu. Context-aware synthesis for video frame interpolation. In *CVPR*, pages 1701–1710, 2018.

[18] Simon Niklaus and Feng Liu. Softmax splatting for video frame interpolation. In *CVPR*, pages 5437–5446, 2020.

[19] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive convolution. In *CVPR*, pages 670–679, 2017.

[20] Simon Niklaus, Long Mai, and Feng Liu. Video frame interpolation via adaptive separable convolution. In *ICCV*, pages 261–270, 2017.

[21] Avinash Paliwal. Pytorch implementation of Super SloMo. https://github.com/avinashpaliwal/Super-SloMo, 2018.

[22] Yingge Qu, Tien-Tsin Wong, and Pheng-Ann Heng. Manga colorization. *ACM TOG*, 25(3):1214–1220, 2006.

[23] Edgar Simo-Serra, Satoshi Iizuka, and Hiroshi Ishikawa. Mastering sketching: Adversarial augmentation for structured prediction. *ACM TOG*, 37(1):1–13, 2018.

[24] Edgar Simo-Serra, Satoshi Iizuka, Kazuma Sasaki, and Hiroshi Ishikawa. Learning to simplify: Fully convolutional networks for rough sketch cleanup. *ACM TOG*, 35(4):1–11, 2016.

[25] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *ICLR*, 2015.

[26] Shuochen Su, Mauricio Delbracio, Jue Wang, Guillermo Sapiro, Wolfgang Heidrich, and Oliver Wang. Deep video deblurring for hand-held cameras. In *CVPR*, pages 1279–1288, 2017.

[27] D. Sỳkora, J. Buriánek, and J. Žára. Unsupervised colorization of black-and-white cartoons. In *Int. Symp. NPAR*, pages 121–127, 2004.

[28] Zachary Teed and Jia Deng. RAFT: Recurrent all-pairs field transforms for optical flow. In *ECCV*, 2020.

[29] Xinrui Wang and Jinze Yu. Learning to cartoonize using white-box cartoon representations. In *CVPR*, pages 8090–8099, 2020.

[30] Zhou Wang, Alan C. Bovik, Hamid R Sheikh, and Eero P. Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 13(4):600–612, 2004.

[31] Zhou Wang, Alan C Bovik, Hamid R Sheikh, and Eero P Simoncelli. Image quality assessment: from error visibility to structural similarity. *IEEE TIP*, 13:600–612, 2004.

[32] Sitao Xiang and Hao Li. Disentangling style and content in anime illustrations. *arXiv preprint arXiv:1905.10742*, 2019.

[33] Xiangyu Xu, Li Siyao, Wenxiu Sun, Qian Yin, and Ming-Hsuan Yang. Quadratic video interpolation. In *NeurIPS*, pages 1647–1656, 2019.

[34] Chih-Yuan Yao, Shih-Hsuan Hung, Guo-Wei Li, I-Yu Chen, Reza Adhitya, and Yu-Chi Lai. Manga vectorization and manipulation with procedural simple screentone. *IEEE TVCG*, 23(2):1070–1084, 2016.

[35] Song-Hai Zhang, Tao Chen, Yi-Fei Zhang, Shi-Min Hu, and Ralph R. Martin. Vectorizing cartoon animations. *IEEE TVCG*, 15(4):618–629, 2009.