

# Wakey-Wakey: Animate Text by Mimicking Characters in a GIF

Liwenhan Xie<sup>\*†</sup>

liwenhan.xie@connect.ust.hk  
The Hong Kong University of Science  
and Technology  
Hong Kong SAR, China

Zhaoyu Zhou<sup>\*</sup>

21210980099@m.fudan.edu.cn  
Fudan University  
Shanghai, China

Kerun Yu

19307110540@fudan.edu.cn  
Fudan University  
Shanghai, China

Yun Wang<sup>‡</sup>

wangyun@microsoft.com  
Microsoft Research Asia  
Beijing, China

Huamin Qu

huamin@ust.hk  
The Hong Kong University of Science  
and Technology  
Hong Kong SAR, China

Siming Chen<sup>‡</sup>

simingchen@fudan.edu.cn  
Fudan University  
Shanghai Key Lab of Data Science  
Shanghai, China



Figure 1: We introduce an approach to revive static text by transferring the motion of characters in a driving GIF.

## ABSTRACT

With appealing visual effects, kinetic typography (animated text) has prevailed in movies, advertisements, and social media. However, it remains challenging and time-consuming to craft its animation scheme. We propose an automatic framework to transfer the animation scheme of a rigid body on a given meme GIF to text in vector format. First, the trajectories of key points on the GIF anchor are extracted and mapped to the text's control points based on local affine transformation. Then the temporal positions of the control points are optimized to maintain the text topology. We also develop an authoring tool that allows intuitive human control in the generation process. A questionnaire study provides evidence that the output results are aesthetically pleasing and well preserve the animation patterns in the original GIF, where participants were impressed by a similar emotional semantics of the original GIF. In

addition, we evaluate the utility and effectiveness of our approach through a workshop with general users and designers.

## CCS CONCEPTS

• **Human-centered computing** → **Graphical user interfaces**; • **Computing methodologies** → **Animation**.

## KEYWORDS

Kinetic typography, Animation, Motion transfer

## ACM Reference Format:

Liwenhan Xie, Zhaoyu Zhou, Kerun Yu, Yun Wang, Huamin Qu, and Siming Chen. 2023. Wakey-Wakey: Animate Text by Mimicking Characters in a GIF. In *The 36th Annual ACM Symposium on User Interface Software and Technology (UIST '23)*, October 29–November 1, 2023, San Francisco, CA, USA. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3586183.3606813>

## 1 INTRODUCTION

Nowadays, kinetic typography, *i.e.*, animated text or motion text, has become common in daily life. These vibrant artifacts can be observed in movies, website widgets, and online memes, such as the lyric video *Skyfall* [1] and the main title sequence of the movie *Spider-Man*. Kinetic typography is effective for expressing emotional content, creating characters, and capturing or directing attention [16, 31]. And there have been fruitful investigations of its

<sup>\*</sup>Both authors contributed equally to this research.

<sup>†</sup>This work is done during academic visit in Fudan University.

<sup>‡</sup>Siming Chen and Yun Wang are the corresponding authors

UIST '23, Oct 29–Nov 1, 2023, California Bay Area, USA

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM. This is the author's version of the work. It is posted here for your personal use. Not for redistribution. The definitive Version of Record was published in *The 36th Annual ACM Symposium on User Interface Software and Technology (UIST '23)*, October 29–November 1, 2023, San Francisco, CA, USA, <https://doi.org/10.1145/3586183.3606813>.

application scenarios, including animated visualization [60], instant messaging [19, 27], ambient displays [37], and captioning [30].

However, it remains non-trivial to craft the animation for text elements. Leveraging commercial animation software [2, 4] or programming toolkits [31] one may tweak the configuration of text in each animation keyframe, *e.g.*, color, positions of the anchor point, and the transition between keyframes like a slow-in easing function. Orchestrating these low-level parameters for a meaningful animation such as a melting scene requires careful considerations like which part of the text element to move, where to move, at what speed, *etc.* As such, this process remains challenging and time-consuming with the large design space. Previous studies [30, 66] tried to alleviate the authoring burden by designing a suite of templates. While categorizing animated effects allows a one-click or even automatic generation, this approach suffers from limited customizability. For instance, when one hopes to impress viewers with a refreshing presentation title, a pre-defined jumping effect may be inferior to a customized motion of breakdance.

Valuing uniqueness and personalization in digital communication [13, 53], we are motivated to find a sweet spot between automaticity and agency in kinetic typography tools. Inspired by the recent advances in artificial intelligence, where a head image can talk by mimicking the motion of a driving video, *e.g.*, [22, 68], we explore transferring existing animation designs to text. The flourishing of GIFs on the web offers myriad high-quality animation references that imply emotions and humor, which can enrich the expressiveness of kinetic typography and make it easy for creators to specify desired effects. However, existing approaches are not directly applicable to our goal. On the one hand, research in motion transfer hardly attends to the non-photorealistic domains [45, 62], especially for kinetic typography. On the other hand, relevant research in text stylization focus on static text (*e.g.*, [23, 63]), where the animation remains largely under-explored.

We propose a mixed-initiative framework for creating kinetic typography based on a driving GIF with a moving character. On the machine side, the motion of the driving GIF is represented as the trajectories of several key points, which are extracted and guide the positional changes in the control points of the target text. On the human side, people can steer the mapping process by directly manipulating these points to refine the automatically computed positions of each point, resulting in a more desirable output. Based on the proposed framework, we develop an interactive interface for creating kinetic typography. We perform a series of evaluation studies to evaluate the usefulness and effectiveness of our approach. First, we demonstrate how individual components of the proposed framework contribute to the final result and test several cases. Second, a questionnaire study shows evidence that the output is both aesthetically pleasing and similar to the driving GIF. Third, we organize a workshop with general users and expert designers to evaluate the utility of our approach.

In summary, our work contributes to the following three aspects.

- (Technique) An automatic approach to transfer the animation scheme from an anchor GIF to vector text.
- (Application) A prototype authoring tool for generating bespoke kinetic typography, which supports various scenarios, *e.g.*, design prototyping and instant messaging.
- (Evaluation) A questionnaire study validating our transfer approach and a workshop demonstrating the usefulness of the authoring tool.

## 2 BACKGROUND & RELATED WORK

In this section, we provide background information on typography and review existing research on kinetic typography, text stylization, and guided animation generation with an anchor.

### 2.1 Preliminaries on Digital Typography

Typography is defined as the art and technique of organizing text in a way that is easy to read, comprehend, and visually pleasing while presented. In general, the visual appearance of a digital letter is determined by its *font*, which is a particular size, weight, and style of a *typeface*. The typeface is a set of designed characters or letters, named glyphs, such as Courier New, Times New Roman, and Bookman Old. Internally, a typeface is represented in the raster domain or vector domain. As bitmap fonts may become distorted or blurred with mosaic-like jagged edges at high resolution, we chose to adopt a vector-based typeface—the TrueType [40] font, which describes glyphs with quadratic bezier curves.

### 2.2 Kinetic Typography

Kinetic typography enriches animated user interface [9] and digital media, which has received scholarly interest since the 1990s [16]. Most recently, Xie et al. [60] summarized a design space of kinetic typography concerning changes in style, shape, position, and scale. Compared with static text, kinetic typography is more competent in guiding attention [7, 37] and communicating emotions or semantics with the paralinguistic clues underlying animation [30, 34]. Accordingly, there has been a series of works seeking to lower the burden of creating kinetic typography. Kinetic Typography Engine [31] set the basis of modern animation software (*e.g.*, Adobe After Effects [2] and TypeMonkey [70]) with frame-based low-level specifications and a library of common effects. The specification concerns text properties like position, rotation, *etc.* And the library was composed of functional time filters like oscillation. TextAlive [24] featured kinetic typography synchronized with audio signals in video editing.

A stream of work investigated tools for average users rather than professional designers, where reducing efforts in animation configurations is a primary goal. These works normally predefined a suite of animated effects and support selection or automatic matching under various contexts. Instant messaging has been most studied, *e.g.*, [17, 20, 38, 66]. For instance, Kinedit [17] allowed users to integrate text animation into a line of words. Minakuchi and Tanaka [38] conceptualized an automatic composer that analyzes the semantics of text and queries suitable motions from a static repository to amplify its meanings. Other scenarios include emotional animation for lyric videos [54] and dynamic display based on viewers' emotions [32]. These works suffered from the number of animated effects provided. For instance, there is hardly any consideration of transforming the text shape, which is common in animation [52] yet requires by-frame editing. In comparison, our work takes advantage of the ubiquitous online memes or stickers and can scrape their animation schemes to a random text with reliable transformation on its outlines.

## 2.3 Text Stylization

Our work closely relates to the task of text style transfer and semantic typography in text stylization, an area widely studied in computer vision/graphics to make a given text visually appealing.

Similar to our workflow, text style transfer concerns transferring the style of a given source (font samples, natural image/video) into text. Some works explored propagating the design of a few stylized letters to others, such as typeface geometry [41] and glyph decorations [56]. Other works followed the general workflow of neural style transfer [18] and viewed style as local neural patterns of the input image/video, e.g., [35, 36, 65]. In contrast, our work deforms the vectorized outline of the text to match the reference GIF. We propose to vivify text by animating it in the way of a cartoon character, which diverges from their focus on learning image patch-based features. Additionally, compared with kinetic typography, text style transfer emphasizes the artistic effect rather than an affective impact and typically produces static output.

Semantic typography amplifies the semantic meanings through visual cues in typography, which is also our goal. Xu and Kaplan [63] proposed calligraphic packing, which deforms letters in a word to fit a given shape, which was improved by Zou et al. [69]. In contrast to the intense deformation in letters, Word-As-Image [23] stroke the balance of transformation on both sides, preserving the original font’s style and legibility while ensuring the semantic implication, which was constrained by a pre-trained Stable Diffusion model [42]. Other approaches operate in the raster domain and leveraged external icons to replace parts of a text [50, 67]. Our work differentiates from these works in that we imply semantics/emotion via animation of the text geometry rather than its static appearance, where the continuity between frames is considered. To the best of our knowledge, this work is the first attempt to incorporate semantics in generating kinetic typography.

## 2.4 Guided Animation Generation

As we aim to produce emotionally or semantically resonant kinetic typography based on a given text, relevant constraints need to be introduced in the animation generation process. Some works infer motions directly from a given still image, concerning features like texture [10, 25, 28], status in a motion cycle [64], periodic patterns [21], etc. These methods are unsuitable for our goal because a text usually appears with no background and is not equipped with equivalently rich properties for motion inference.

Motion transfer has been a standard task in computer vision, which is to generate a video based on a source image and a driven video by learning the motion from the driving video while preserving the appearance of the source image. Monkey-Net [44] was the first model-free approach to transfer motions of arbitrary objects by aligning key points between the source and target domain. FOMM [45] further enhanced it with local affine transformations on the extracted key points. It is one of the state-of-the-art models and we adapted it to fit the vector-based text. Specifically, we maintained the text legibility by regularizing motion anchors with the distance change in the Laplacian coordinate. Our method shares the same idea to preserve the structural information as DAM [49], which introduced a latent root anchor to model the structure of objects.

Different from our focus on the text, most existing datasets and models concern talking heads and human posture (e.g., [8, 22, 46, 47, 68]) and do not yield desired results on texts where legibility matters (see Section 4). Our exploration of kinetic typography contributes to a unique case of cross-domain motion transfer.

In addition to fully automatic approaches, mixed-initiative interfaces for animation authoring have been investigated. Users may specify the intended effect with sketch-based demonstration [25, 26, 58, 61], gestures [5], or examples [14]. Pose2Pose [59] supports creating cartoon character animation by minimizing the design efforts through clustering postures and automatically matching the stylized postures designed by the artists to the driving video. Most similar to our work, Live Sketch [48] leveraged motion transfer to let novice users create animated sketches, where users are required to define control points in both the source and target domain. Our approach also allows users to specify their desired animation effect through a GIF, which is easy to access online. However, the key points in the driving video are automatically extracted and automatically mapped to the target domain. For a finer-grain control, users can adjust the extracted key points and internal parameters.

## 3 DESIGN CONSIDERATIONS

Motivated to lower the barrier in creating kinetic typography, we explore motion transfer techniques. Instead of tweaking keyframe configurations from scratch, users may specify the desired animation effect based on a reference GIF. With numerous online GIF instances, users may derive more diverse animation effects compared with using template-based tools.

One major design consideration is to *support both direct generation and fine-grain refinement (C1)*. We expect our approach can generalize to various user requirements, including casual use as in online-messaging and professional editing like video-making. In addition to producing one-off results, the tool should allow refinement over fine-grain configurations of each frame. This is because motion transfer inherently introduces uncertainties in the generated result from motion transfer, which may violate user preference.

Additionally, we strive to *empower creators with interpretable algorithmic parameters (C2)*. We hope the system supports iterative refinement, which necessitates providing explanations for the generation process so that users can provide feedback and make adjustments at every step of the generation process. By combining human experience and supervision, we seek to achieve higher quality and more consistent generations in line with users’ expectations.

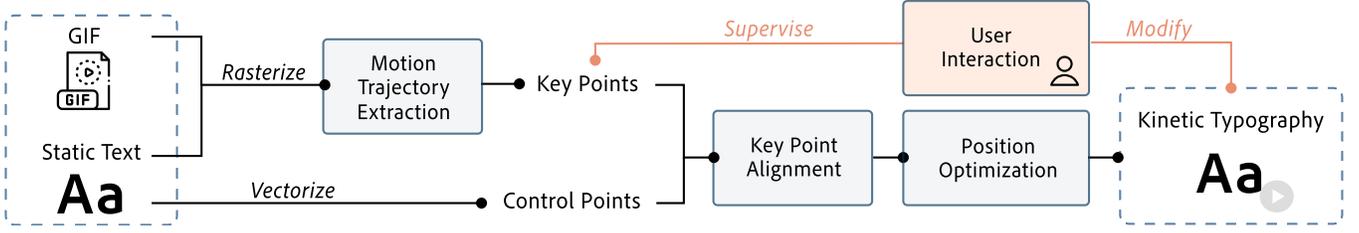
## 4 FRAMEWORK

In this section, we introduce a general framework for transferring the motion of a given GIF to a static text.

### 4.1 Overview

Figure 2 illustrates our framework. The computational pipeline takes in a driving GIF and static text as input and outputs the kinetic typography. Users can tweak the intermediate key points and control points in the generated result (C1).

Internally, the input text is represented in the TrueType format [40]. It is first converted into an image and fed into a FOMM



**Figure 2: Overview of our approach.** Inputs are a driving GIF and static text. The output is kinetic typography echoing the GIF’s animations. The motion trajectory extraction module captures the key points in the driving GIF. The key point alignment module aligns the control points of the vectorized text to the key points. The position optimization module regularizes the text outline. And the User Interaction module allows human intervention on the intermediate key points and final results.

model [45] together with the anchor GIF to obtain the trajectory of the motion key points at each frame  $X_i^f$ , where the model identifies  $N$  key points, and the GIF consists of  $F$  frames, *i.e.*,  $i = 1, \dots, N$ ,  $f = 1, \dots, F$ . The input text is also parsed to the initial control point set  $C_j^0$  of its glyphs, with a total of  $M$  control points, *i.e.*,  $j = 1, \dots, M$ . A local affine transformation is applied to both the initial control point set  $C^0$  and the key point set trajectory  $X^f$  to obtain the motion trajectory of the control point set  $C_j^f$ . The updated control point trajectory  $C_j^f$  is attained through position optimization. Subsequently, a vectorized glyph sequence is generated, culminating in the creation of animated text in vector form. Through the user interaction module,  $X_i^f$  and  $C_j^f$  can be directly manipulated, and users can control some hyperparameters (C2).

## 4.2 Motion Trajectory Extraction

We convert the static text to an image, and input it along with the anchor GIF to obtain the trajectory of motion key points. We adopted motion transfer to support the fast and flexible generation of kinetic typography. As the object in the GIF usually differs from the text in shape, we need to separate the appearance and extract the motion trajectories of key points from the source GIF.

FOMM [45] is applied for key points extraction in our task. It is a self-supervised method using a framework that decouples appearance and motion, which effectively enriches the possible transferable motions to support motion transfer within any object category. To address the problem of large differences in key points between the driving frame  $D$  and the source image  $S$ , the FOMM model introduces an abstract reference frame  $R$  and obtains  $\mathcal{T}_{S \leftarrow D}$  by separately calculating  $\mathcal{T}_{S \leftarrow R}$  and  $\mathcal{T}_{D \leftarrow R}^{-1}$ .

$$\mathcal{T}_{S \leftarrow D} = \mathcal{T}_{S \leftarrow R} \circ \mathcal{T}_{R \leftarrow D} = \mathcal{T}_{S \leftarrow R} \circ \mathcal{T}_{D \leftarrow R}^{-1},$$

where  $\mathcal{T}_{A \leftarrow B}$  denotes the mapping from the image  $B$  to  $A$ .

In the implementation,  $\mathcal{T}_{S \leftarrow R}$  and  $\mathcal{T}_{D \leftarrow R}$  are obtained by key points detection in  $S$  and  $D$ , respectively, which supports us to extract the key point trajectories from both the source and generated pixel-based text GIFs. Either of the two trajectories of the key points can be applied to drive the subsequent generation, and we use  $X_i^f$  to represent the selected key point trajectory for simplicity. The separate detection mode also supports the relative generation

( $\mathcal{T}_{S_i \leftarrow S_1}$  to deform from the source image) following a similar mindset, in addition to the absolute way ( $\mathcal{T}_{S_i \leftarrow D_i}$  to deform from the corresponding frame of the source GIF directly).

In our implementation, we utilized the pre-trained FOMM model on the MGif dataset [44], which has shown good performance in key point detection. Following the pre-trained model, the features extracted for each frame are estimated independently, and the number of key points is set to 10. However, to further enhance the integration of emotion into generation and analysis, we gather and create a dataset of Puppy Maltese [39] with 77 emotional-labeled GIFs and use it to fine-tune the model. This is done to better cater to the needs of the subsequent case studies and user surveys.

Due to the difficulty of FOMM in achieving good performance in motion transfer across different categories of objects, we only extract intermediate results from the key point detection module and redesign the subsequent generation steps based on our task. We compare our result with the rasterized output of FOMM in a crowdsourcing study introduced in Section 7.

## 4.3 Key Point Alignment

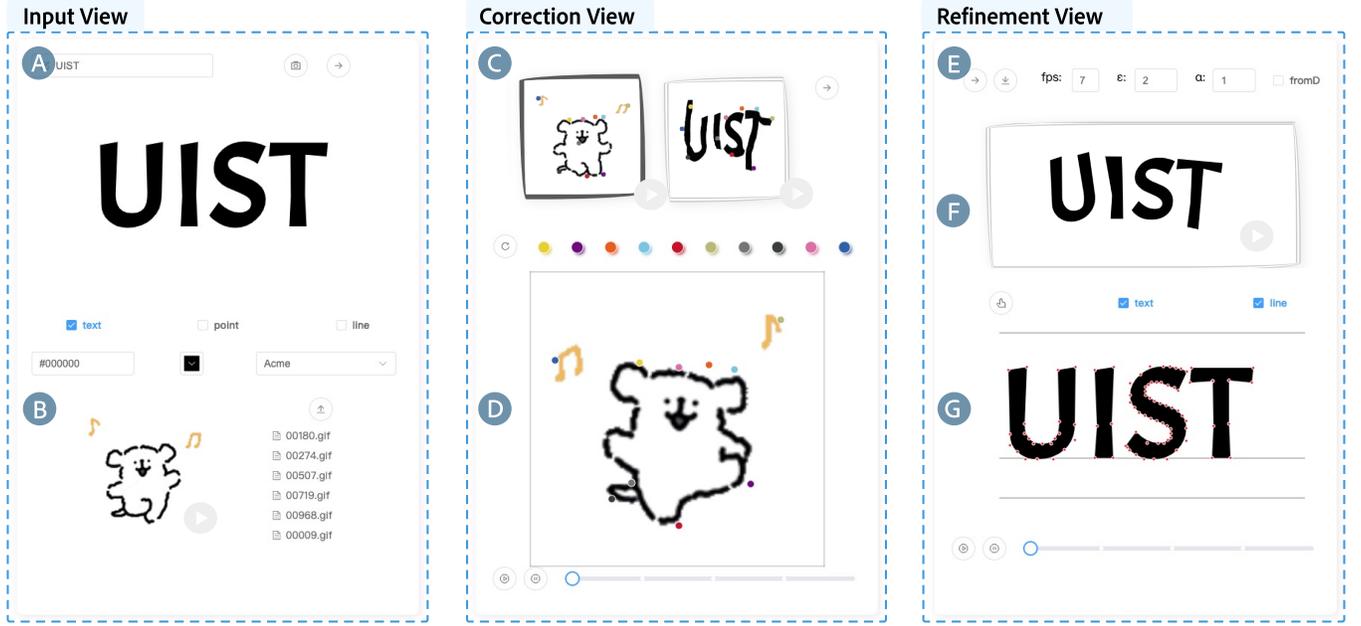
A local affine transformation is applied to align the motion trajectories of key points to the control points. It introduces non-linearity to preserve local information better and achieve richer deformation. In our task, each key point extracted from each frame from the GIF is considered a local region, and the local affine transformation matrix set is obtained by computing the translational transformation of each key point in adjacent frames. The global non-linear transformation is then calculated using a distance-weighted interpolation-based approach with the matrix set.

$$\begin{bmatrix} C_j^{f+1} \\ 1 \end{bmatrix} = \sum_{i=1}^N w_i(C_j) \cdot \begin{bmatrix} C_j^0 \\ 1 \end{bmatrix} \begin{bmatrix} \mathcal{I} & 0 \\ (X_i^f - X_i^1)^T & 1 \end{bmatrix},$$

$$w_i(C_j) = \frac{1/\|C_j^0 - X_i^1\|^e}{\sum_i 1/\|C_j^0 - X_i^1\|^e}.$$

$C_j^f$  and  $X_i^f$  denote the control point  $j$  and the key point  $i$  at frame  $f$ , respectively. The control point’s position at each frame is calculated in reference to the key point at the first frame to achieve global stability.  $\mathcal{I}$  is a 2nd-order identity matrix.  $w_i$  is a weight function for a control point with respect to the key point  $i$ .

The weight decays according to the inverse of the  $e$ -th power of the relative distance from  $X_i$  to  $C_j$ , where  $e$  controls the locality



**Figure 3: Wakey-Wakey: An authoring interface to interactively create anchor-based kinetic typography. There are three views: input view, correction view, and refinement view. (A) Input and preview the text, where font and color can be specified. (B) Upload a driving GIF. (C) Preview the matching of key points between the text and GIF at each frame. (D) Directly manipulate key points locations. (E) Fine-tune the hyperparameter. (F) Preview result GIF. (G) Refine the text control points at each frame.**

of the affine transformations, *i.e.*, the degree to which each affine transformation affects the target point.

#### 4.4 Position Optimization

To alleviate inappropriate deformation of glyphs caused by changes in the relative position of the control points, we optimize the positions of the control points by frame based on the Laplacian coordinate, which generally describes the relative positions on the surface using the neighbor information. For a control point  $j$  at frame  $f$ , its Laplacian coordinate  $L_j^f$  is calculated as

$$L_j^f = \sum_{k \in N_j} \omega_{jk}^f (C_k^f - C_j^f) = \sum_{k \in N_j} \omega_{jk}^f C_k^f - C_j^f,$$

where  $C_j^f$  and  $C_k^f$  denote the Cartesian coordinate of the control point  $j$  and  $k$ , respectively.  $N_j$  is the set of  $K$ -nearest neighboring control points with the smallest Euclidean distance to the control point  $j$ , which is calculated based on the initial control points set  $C^0$ , invariant to changes in  $f$ .  $\omega_{jk}^f$  denotes the weight of the neighbor point  $k$  in the Laplacian representation of the control point  $j$ , where

$$\omega_{jk}^f = \frac{1/\|C_k^f - C_j^f\|^2}{\sum_{k \in N_j} 1/\|C_k^f - C_j^f\|^2}.$$

Considering the inhomogeneity of the distribution of discrete sampling points, we use the aforementioned distance-based weights to describe the detailed location information better. Further, we

design the following objective function  $\mathcal{L}_{\text{total}}$  to optimize the coordinates of the sequence of control points obtained by frame.

$$\begin{aligned} \mathcal{L}_{\text{total}} &= \alpha \cdot \mathcal{L}_{\text{glyph}} + \mathcal{L}_{\text{motion}}, \quad \alpha \in [0, +\infty), \\ \mathcal{L}_{\text{glyph}} &= \sum_{j=1}^M \|L_j^f - L_j^0\|^e, \quad \mathcal{L}_{\text{motion}} = \sum_{j=1}^M \|C_j^f - C_j^{f'}\|^e. \end{aligned}$$

$L_j^f$  and  $L_j^0$  denote the Laplacian coordinates of the control point  $j$  in frame  $f$  and 0, respectively.  $C_j^f$  and  $C_j^{f'}$  denote the coordinates of the control points before and after optimization.  $\mathcal{L}_{\text{glyph}}$  measures how much the local shape details are preserved, which is computed as the sum of the distance of Laplacian coordinates between the optimized and initial control points.  $\mathcal{L}_{\text{motion}}$  measures how much of the motions are preserved, as the sum of the distance of the control points before and after optimization, *i.e.*, minimizing edit distance.  $\alpha$  is a hyperparameter representing the trade-off between the two loss functions. The larger  $\alpha$  is, the more details of the initial glyph and the less motion are preserved. As for the norm  $e$ , the larger it is, the locality is more regulated, which leads to stronger deformation. We use a  $K$ -dimensional tree to accelerate the nearest-neighbor search, where the parameter is empirically set:  $K = 3$ . While we employ frame-by-frame optimization, we note it is worth introducing global temporal regularization terms in the loss function to promote smoothness and consistency.

#### 4.5 User Interaction

The user interaction module allows direct manipulation of the computed positions of the key points and control points at each frame,

	Time →					Time →					Time →				
$\alpha = 0$	sleepy	sleepy	sleepy	sleepy	sleepy	thanks	thanks	thanks	thanks	thanks	thanks	wakey	wakey	wakey	wakey
$\alpha = 2$	sleepy	sleepy	sleep	sleepy	sleepy	thanks	thanks	thanks	thanks	thanks	thanks	wakey	wakey	wakey	wakey
$\alpha = 4$	sleepy	sleepy	sleepy	sleepy	sleepy	thanks	thanks	thanks	thanks	thanks	thanks	wakey	wakey	wakey	wakey

Figure 4: Comparison of the generation results with  $\alpha = 0, 2, 4$ . Increasing  $\alpha$  enhances the smoothness of the glyph, but an excessive value of  $\alpha$  may negatively impact the amplitude of the motion.

i.e.,  $\{X_i^f\}$  and  $\{C_j^f\}$ ,  $\forall i \in [1, N] \cap \mathbb{N}, j \in [1, M] \cap \mathbb{N}, f \in [1, F] \cap \mathbb{N}$ . In this way, creators of kinetic typography can participate in the motion transfer process and adjust the final results according to their needs. The hyperparameter  $\alpha$  in the position optimization stage can also be adjusted for different texts, as illustrated in Section 7.1.

## 5 AUTHORIZING INTERFACE

Based on the proposed framework, we implement a mixed-initiative authoring tool called Wakey-Wakey<sup>1</sup> that allows fine-grain adjustment for more natural and aesthetic results (see Figure 3). This section offers a step-by-step guide showing how to interactively generate animated text with our tool referring to the interface.

The user-oriented process mainly consists of three steps: text and GIF input, key point correction, and glyph refinement, each supported by a view: *Input View*, *Correction View*, and *Refinement View*. While the input step is mandatory, the correction and refinement stages are optional. This allows for simple end-to-end personalized generation, as well as interactive improvement.

*Input View*. Users first input the text and customize its static appearance with the global typeface and color (Figure 3 A). They can upload and preview the driving GIF through a button (Figure 3 B). The section will record and list the recent upload history. After clicking “Next”, the system will process the input with our method, and both intermediate and final results will be displayed in the other two views. Users can easily obtain the generated animated text here without any additional effort.

*Correction View*. Users can then drag the displayed the key points from the motion trajectory extraction module to a suitable location at a specific frame, as shown in Figure 3 D. Special attention can be paid to the key point trajectory with a corresponding colored button above. Two thumbnails (Figure 3 C) enable switching between the anchor GIF and the extracted key points for correction. As the corresponding key points share the same color, users can learn how the mapping is. Users are supported to ensure better quality by maintaining reasonable motion trajectory to drive the generation of the vector animated text.

*Refinement View*. Users may configure the parameters  $\alpha$  and  $e$ , and select whether the kinetic typography is generated by aligning to key points from the anchor GIF or the extracted key points (Figure 3 E). The bottom panel (Figure 3 G) enables the users to adjust the glyph by dragging the control points. And the final GIF of kinetic typography is displayed in the middle panel (Figure 3 F).

<sup>1</sup>The name suggests that the authoring tool awakes static text and makes it lively.

## 6 IMPLEMENTATION

Wakey-Wakey<sup>2</sup> was implemented as a client/server web application. The front end was built with Vue for user interactions. The computational framework for generating kinetic typography was implemented in Python. An automatic generation takes around 300ms/frame (CPU: Intel i7 4.9 GHz). The Flask framework is used to handle the messaging between the front end and the back end.

## 7 METHOD ANALYSIS

Due to the absence of pre-defined “ground truth” and lack of standard metrics in the nascent area of motion transfer for quantitative assessment, we empirically evaluated our approach by (1) analyzing the impacts introduced by each component, (2) comparing the automatically generated result from different styles of GIFs and typefaces, and (3) conducting questionnaire studies to understand how general people perceive the outputs based on several cases.

### 7.1 Effects of Components

We evaluated the effect of each component in the workflow to analyze how our adaption to FOMM and the introduced human interventions can improve the generated result, including local position optimization, vectorized text representation, key point correction, and glyph refinement.

*Local Position Optimization*. The position optimization module is introduced to preserve the local shape of each glyph better. Figure 4 demonstrates three motion transfer results with  $\alpha$  set to 0, 2, 4.  $\alpha$  is the weight of  $\mathcal{L}_{\text{glyph}}$ , which controls the degree of preservation of local shape. As can be seen, when alpha is set to 0, i.e., without local position optimizing, some local parts of the glyphs are unsatisfactory, such as the “p” in “sleep”, the “t” and “k” in “thanks”, and the “a” and “k” in “wakey”. With the increment of  $\alpha$ , the glyph becomes smoother. However, when alpha is too large, it may cause too much preservation of the original glyph and result in a loss of motion, for example, the “w” and “k” of “wakey” when  $\alpha = 4$ . Through experiments, we find a suitable default value of 2.

*Vectorized Text Representation*. Instead of directly employing existing image-based motion transfer models, our approach operates on the control points of text glyphs. As shown in Figure 5, the output results based on pixels are not stable enough. For example, in the pixel-based glyph generated by FOMM, the letter “y” of “angry” has breaks and extra noisy strokes. As the anchor GIF is hardly the targeted category of kinetic typography, using FOMM for motion transfer does not produce satisfactory results. In contrast, our

<sup>2</sup>Source code available at <https://github.com/KeriYuu/Wakey-Wakey>.



Figure 5: A comparison with the FOMM model [45]. Our approach operates on the control points of vectorized text, which improves legibility.

method better preserves the integrity and legibility of the glyph and produces a more stable frame sequence.

*Key Point Correction.* The key points  $\{X_i^f\}$  detected by the model may not always be accurate, which can result in unexpected deformations in the generated animated text that rely on these key points. Our approach allows users to interactively correct the key points, thereby obtaining a more desirable generation that aligns with their expectations. As shown in Figure 6, by analyzing the preceding and following frames, we can find that in the fourth frame of the pixel-based animated text image sequence generated by FOMM, the key point marked in red noticeably shifts towards the right. This caused an excessive deformation towards the right in the lower right part of the letter “W” in the vector-based animated text generated with this key point. By dragging the key point towards the left to an area consistent with the preceding and following frames, the deformation of the generated glyph appears more reasonable and smoother.

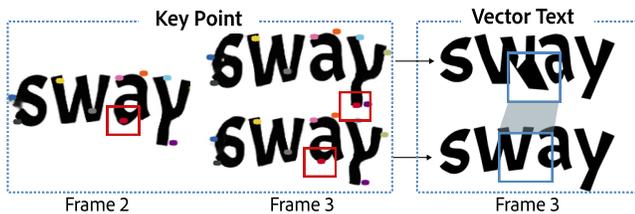


Figure 6: Comparison of the output before (top) and after (bottom) key point correction. The automatic mapping fails when the highlighted red key point shifts to another location in two consecutive frames, which yields distortion in the text.

*Glyph Refinement.* The control point sequence  $\{C_j^f\}$  can be manually updated for fine-grain refinement. Through the authoring interface, users are supported to drag the control points and preview the result immediately. As shown in Figure 7, the sharp corners inside the first letter “p” affect the glyph aesthetics, where the highlighted left line segment is tilted to the left and needs to be adjusted. By moving the three control points in the sharp corner area to the right and adjusting the relative positions of the three points, the refinement process is done. It is evident that the updated glyph achieves a better effect through simple and immediate dragging.

## 7.2 Generalizability

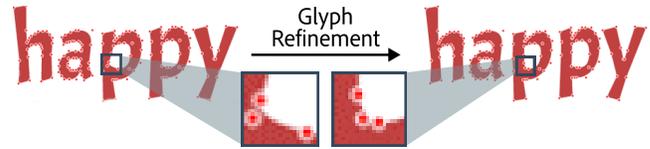


Figure 7: Manual refinement of the text control points. By dragging the distorted control points, one may intuitively refine the output kinetic typography at a fine-grained level.

Drawing from our experience, we reflect on the generalizability of our approach in terms of the driving GIF and the input typeface.

In general, Wakey-Wakey can accommodate input GIFs with a clean background and a simple-shape moving rigid body, such as instances from the Puppy Maltese dataset. This is because our implementation adopts the pre-trained FOMM model based on the MGif dataset, which features a white background and one cartoon animal. Seen from Figure 8, the automatic key point extraction may fail and cause large distortion when there are multiple moving objects, or the moving object exhibits complex patterns. A complex background also threatens the reliability of extracted key points in the driving GIFs, such as a clip from natural videos. While these issues can be addressed by manual correction, we also note that the motion trajectory extraction module can be improved by unsupervised training on a larger dataset with representative cases or using a large universal model.

As for the input fonts, our approach empirically performs well for typefaces with more than 5 control points in a glyph. The more control points encapsulated in the typeface, the more likely that the position optimization can maintain its legibility. Figure 9 showcases the automatic generation results for ten typefaces of common classes. Typefaces with the most control point number also yield the most smooth results, including Fredericka and Cedarville. However, there might be strong deformation for handwriting-styled typefaces, potentially due to their high flexibility.

## 7.3 Questionnaire Study

We conducted two questionnaire studies to evaluate the effectiveness of our approach. Specifically, we seek to understand (1) whether our approach convincingly transfers the motion, and (2) to what extent the semantics of the original GIF can be preserved.

*7.3.1 Setup.* The questionnaires are distributed on Qualtrics. Participants are required to complete Study I before Study II. And the questions appear in a random order in each study. We used meaningless pseudo-words from the Lorem Ipsum corpus [33] as input text in order to minimize the influence of text content. For driving GIFs, we used the Puppy Maltese dataset to generate cases. To avoid confounding effects, the driving GIFs are non-repetitive. And we employed the typeface “Akronim” for it has over 300 control points, which may lead to satisfying results without human intervention and therefore suitable for our scenario requiring batch generation.

*Study I: Motion Transfer.* The first study aimed to evaluate the overall quality of the output kinetic typography. As no quantitative metric is available in our task, we obtained subjective assessments by asking the participants to rate the similarity between the driving

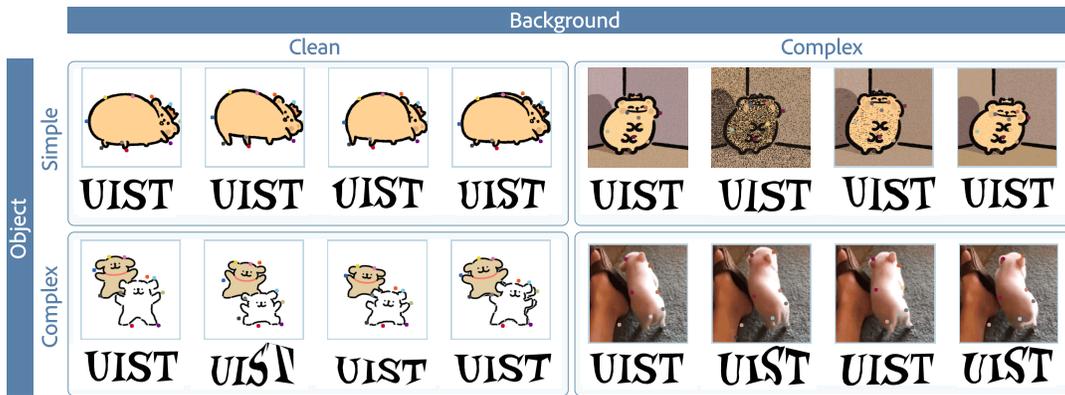


Figure 8: A comparison of results from driving GIFs of different complexities in background and target object(s).

Typeface	serif		sans serif		monospace	display	handwriting			
	Merriweather	PT Serif	Open Sans	Oswald			Cutive Mono	Cousine	Fredericka	Acme
Avg. #Control	62	55	31	47	63	32	1588	30	144	67
	wakey	wakey	wakey	<b>wakey</b>	wakey	wakey	<b>wakey</b>	<b>wakey</b>	wakey	wakey
	wakey	wakey	wakey	<b>wakey</b>	wakey	wakey	<b>wakey</b>	<b>wakey</b>	wakey	wakey
	wakey	wakey	wakey	<b>wakey</b>	wakey	wakey	<b>wakey</b>	<b>wakey</b>	wakey	wakey
	wakey	wakey	wakey	<b>wakey</b>	wakey	wakey	<b>wakey</b>	<b>wakey</b>	wakey	wakey

Figure 9: A comparison of results from typefaces of different categories and average number of control points for the 26 English alphabets (Avg. #Control). The first column shows four key frames of the driving GIF. Each column in the rest shows the font information and the corresponding motion transfer result.

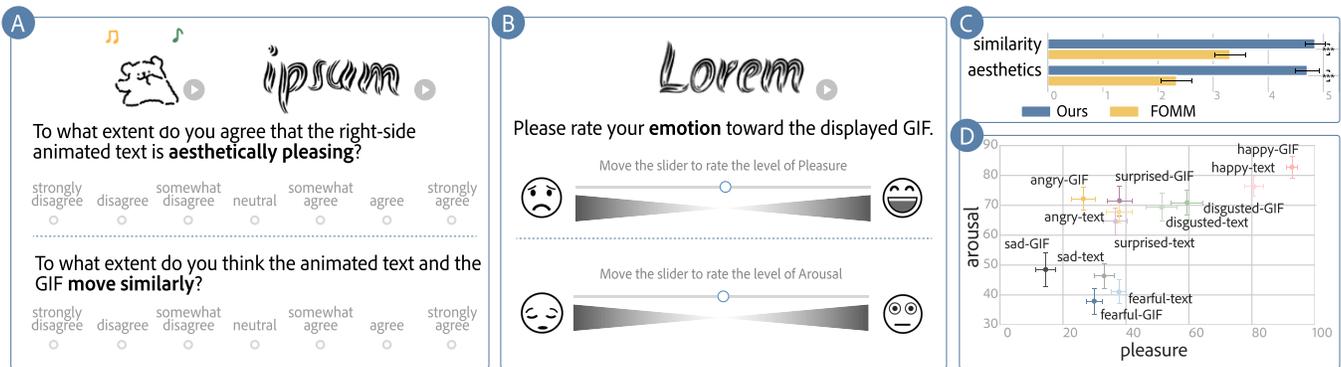
GIF and the output kinetic typography and their aesthetics. On the one hand, the similarity between the source and target is the primary goal in motion transfer. On the other hand, aesthetics is a common pursuit in animation design.

A sample question is shown in Figure 10 A. When designing the questionnaires, we tried to familiarize participants with simple and concrete questions. For instance, we asked whether a GIF is “aesthetically pleasing” to align participants’ appraisal of the aesthetic property to their feelings [11]. For each question, the animated text and the corresponding anchor GIF are displayed, and participants are asked to rate them on a 7-point Likert scale (0–strongly disagree to 6–strongly agree) for aesthetics and motion similarity, respectively. 20 driving GIFs were randomly selected. For each driving GIF, we set two questions, one for the baseline–rasterized kinetic typography generated with FOMM, and one for the experimental group–vectorized kinetic typography with our approach. Therefore, a questionnaire consists of 40 questions. Considering the influence of the font, participants are asked to rate the aesthetics of the static font before viewing the main body of the questionnaire.

*Study II: Semantic Preservation.* The second study aimed to verify whether the learned animation can preserve the semantics of the driving GIF. We tackled the problem from the perspective of emotion, which is an integral part of semantics. Specifically, we focused on Ekman’s six basic emotions [15], i.e. sadness, happiness, fear, anger, surprise, and disgust.

Figure 10 B illustrates a sample question. We adopted the Affective Slider [6] for participants to self-report their emotions, which consist of two dimensions: pleasure and arousal from the extent 1 to 100. Pleasure means the degree of positivity or negativity of an individual’s emotional state. Arousal corresponds to the level of physiological activation or stimulation in an individual’s emotional state. For each basic emotion, we selected two GIFs as anchors according to the pre-defined labels in the dataset. A question comprises one kinetic typography, where we required participants to assess the perceived emotions. Hence, there were 12 questions.

*7.3.2 Participants.* We recruited participants from a local university by posting advertisements on social media. Each participant is paid £3.5 for completing the questionnaire. A total of 33 people



**Figure 10: Example questions and results in two questionnaire studies (N=33).** (A) Study I: Subjective ratings for the aesthetic property of the output kinetic typography and similarity between the target and source GIF. (B) Study II: Perception of emotions underlying kinetic typography. (C) The average ratings and standard errors of cases in Study I, where our approach outperforms the baseline FOMM in motion similarity and result aesthetics. (D) The average ratings and standard errors of the pairwise (pleasure, arousal) ratings for sample GIFs, where the driving GIF and corresponding kinetic typography posit in adjacent areas. Ratings of the same emotion are encoded with colors of a similar hue.

signed up for the questionnaire study. Most participants were between 18–24 years old, with 15 females and 18 males. In addition, all participants reported using emojis frequently in daily communication, where 14 (42%) reported to use emojis *very often*.

### 7.3.3 Result Analysis.

**Study I. Motion Transfer.** Participants spent an average of 10.4 minutes in completing the 20 questions (std=6.2, ranging from 3.5 to 28). We deemed all the responses valid. Seen from Figure 10 C, participants generally recognized the aesthetics and the similarity of movement between the driving GIFs and the animated text in our approach. For our approach, the average score on the aesthetics was 4.71 (std=1.21), and the motion similarity was 4.85 (std=1.06), where both scores exceeded 4, *i.e.*, somewhat agree, on the 7-point scale. For the baseline, the average aesthetic score was 2.34 (std=1.60) and the average similarity score was 3.31 (std=1.59). Compared to the baseline method, our approach obtained significantly higher scores (significance level  $\alpha = 0.001$ , Student’s t-test), which suggests that our method outperforms the FOMM in terms of motion transfer, making the generated animations more visually appealing and more similar to the source motion. This finding echoes our ablation study on the vectorized text representation, where we identified certain glitches in the pixel-based methods. Moreover, after adding the dynamic motion, the aesthetics of the text improves compared to static text with an average aesthetics score of 3.79, further validating the effectiveness of our method.

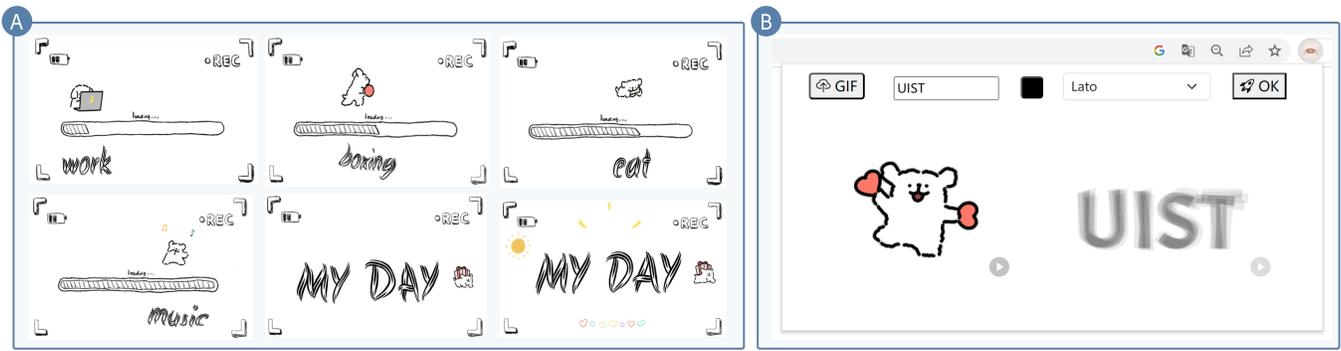
In summary, Study I verified that our model could achieve better motion transfer compared to the baseline. The improved aesthetics and motion similarity scores, along with the user comments, demonstrated the effectiveness and applicability of our method in generating visually appealing and motion-consistent animated text.

**Study II. Emotion Preservation.** Participants spent an average of 8.5 minutes to complete the questionnaire, with a maximum of 30 minutes and a minimum of 2 minutes (std=8.9). The results

of study II are shown in Figure 10 B, where the scatterplot maps the average scores of the (pleasure, arousal) pairs. The attached error bars indicate the standard error of each dimension with their lengths, where vertical for arousal and horizontal for pleasure.

Inspecting the diagram, we could see that the data points for the same emotion under both the driving GIF and Text are projected on adjacent areas, revealing that the expression of emotion has also been successfully transferred through the motion transfer. Furthermore, different emotion classes largely varied. For example, there is a significant difference between happy and sad emotions, demonstrating the effectiveness of our method in preserving the emotional semantics of the anchor GIFs. One could see that the emotions of disgust and anger are very close, and the demarcation is not obvious, with only the ordering of pleasure being altered. This suggested that the emotional expressions in these two emotions might be similar, making it harder for users to differentiate them clearly. In addition, the variances in the arousal dimension were generally greater, possibly due to the users’ perception of arousal being more ambiguous or subjective, leading to a wider range of responses. In contrast, the arousal and pleasure scores for the text are more neutral (around 50 points). Although the emotional expression has been learned and transferred, it is not as strong as in the source motion pictures. It might result from the limitations of the method or the inherent difference between text and the figure-like domain in representing emotions.

In summary, results from Study II showed that our method could effectively preserve the emotional semantics when transferring animations from the driving GIF to a text. However, some emotions may not be as distinct as they are in the original anchor GIFs. And users’ perceptions of arousal might be more ambiguous. While the questionnaire shows the success of emotion transfer on the particular typeface being used, more studies are needed to validate similar mechanisms for other fonts, as we did not eliminate the influence on emotion perception from the typeface.



**Figure 11: Demonstrations in the workshop. (A) A video opening featuring synchronous animation of textual descriptions and the cartoon character. (B) A browser plugin for the automatic generation of kinetic typography based on our approach.**

## 8 WORKSHOP

We organized a workshop to evaluate the utility of our method.

### 8.1 Demonstrations

To elicit in-depth discussions in the workshop, we designed and implemented several demonstrations of potential application scenarios, including a video opening, an online messaging widget, and an emotional word cloud.

*Video opening.* Vlogs (Video blogs) are becoming a prevalent form to share personal experiences creatively. To make a vlog stand out, an engaging opening animation can help grab viewers’ attention and set the tone for the rest of the video. Leveraging results generated by our method, we created a vlog opening as an example of its practical application in daily content creation. As shown in Figure 11 A, D1 is a vlog opening with the Puppy Maltese theme, showcasing different states of the character and introducing “My Day”, which suggests the video topic. We envision that integrating animated text with consistent motions enhances both explanatory and entertaining values. Similarly, users can use our method to create and personalize animated text in their video creations across different themes and contents.

*Browser Widget for Online Chatting.* Informed by the previous efforts in enhancing emotion communication in online messaging [3, 34, 55], we implemented a light-weighted Chrome extension to facilitate the real-time creation of kinetic typography (see Figure 11 B). It has a simplified interface compared with Wakey-wakey, which features real-time generation and removes the human interaction module. Users may upload an anchor GIF, type down the text, configure the color and font, and then directly obtain the generated kinetic typography in the GIF format.

*Emotional Animated Word Cloud.* The word cloud is a common visualization technique to summarize text data, where the text size represents the word frequency. Xie et al. [60] coined the word “emordle” representing animated word clouds that suggest underlying emotions. Based on our approach, we generated an “emordle” by transferring the animated scheme of a bumpy cartoon pig from the MGif dataset. Instead of using the parsed control points, we transferred the text anchors for each text element that constitute the



**Figure 12: Application of animated word cloud that delivers a certain emotion with the animation.**

word cloud. In other words, we replace the vector control points with the anchor points of each word in the proposed approach. Note that one word has one anchor point at its central position. The generated word cloud example is displayed in Figure 12.

### 8.2 Protocol

The workshop proceeded in the following four stages.

*Briefing.* The briefing session took 10 minutes, during which we introduced the background of our work and briefly explained our method. We then used a demo video to demonstrate the functions of the tool and illustrate the generation process.

*Demonstration.* We spent another 10 minutes displaying the video opening, the online messaging widget, and the animated word cloud to demonstrate potential application scenarios. We also introduced the installation and use of the plugin through a demo.

*Self-creation.* We then allow users to freely use and explore the system and widget to create their kinetic typography for 20 minutes.

*Post-interview.* After trying Wakey-Wakey, participants are asked to complete a questionnaire with six questions on a 7-point Likert scale about its usability, including covering *Practicality*, *Customization*, *Pleasure*, *Efficiency*, and the *Intuitiveness* for widget and interface respectively. We also raised open-ended questions following a structured template in order to understand their perceptions of the authoring process as well as the generated kinetic typography.

### 8.3 Participants

Both designers and general users are invited to obtain feedback from different perspectives in the evaluation. We recruited 20 participants through our personal network and advertisements on social media, with 7 females and 13 males. There are 3 professional designers: P1 works on user experience design; P2 engages in self-media creation as a blogger and vlogger; P3 is a digital painter (P3). The rest are graduate students majoring in data science at a local university (denoted as P4–P20 in ascending time order of their interviews). Among all participants, 13 have seen kinetic typography before, such as in online memes, short videos, slides, and advertisements. Only two participants have experience in creating kinetic typography with other software (P2, P3).

### 8.4 Results

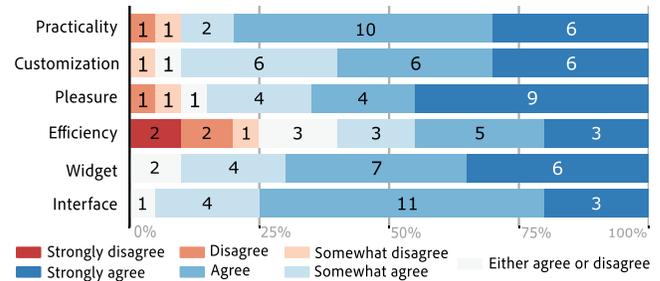
Here we present both the quantitative and qualitative results.

**8.4.1 Observations.** Nonetheless, when it comes to the time required for creation, there is substantial dissent, as our concurrency fell short, leading to extended completion times for some users. In spite of this setback, the system garners favorable acknowledgment for its expressiveness, user-friendliness, intuitiveness, and real-world applicability. To elevate the overall user experience, refinements in creation duration and concurrency should be considered. As for the usage of  $\alpha$  or  $e$ , most users were generally satisfied with the empirical default value. However, two users with a design background occasionally fine-tuned this parameter to derive better results. In terms of the manual adjustments on control points, users without a design background barely adjusted the control points. Three users with design backgrounds occasionally make manual adjustments, about 8/14 frames, with an average duration of 2.1/5.8 min for one kinetic typography. Users tended to adjust key points in GIFs when there were noticeable detection errors. Otherwise, they normally increased  $\alpha$  to mitigate unexpected deformations, though less motion preserving.

**8.4.2 Usability Ratings.** Figure 13 shows the distribution of users' subjective ratings of the six usability questions. More than half of users concurred that the tool was intuitive, tailored, and delivers an enjoyable experience during usage. To be more specific, 45% users strongly agreed that our work is pleasant and interesting. They valued the innovative animated text design, the straightforward comprehension of emotions portrayed, and the uncomplicated tool configurations for selecting animation approaches. They also observed that blending textual meanings with animations facilitated a beneficial expression of the content.

**8.4.3 User Feedback.** We summarize the following insights and implications for future improvements from the users' feedback.

◊ On mixed-initiative authoring. Participants expressed their agreement with the balance we struck between user involvement and automatic generation. P3 said: “supporting quick automatic generation while also offering optional customization and improvement from users”. Participants were optimistic about the mixed-initiative way, as “intelligence improves efficiency while users enhance quality and creativity, as the imagination of users cannot be dismissed” (P10).



**Figure 13: Quantitative Results of the usefulness of our method. We measured the practicality, customization, pleasure, efficiency and the intuitiveness for the browser widget and authoring interface respectively.**

◊ On personalization. 19 participants (except P11) valued customization highly and believed it helps incorporate their own ideas and shape a unique personal style, while three designers showed an “innate aversion to preconceived solutions” (P1). Among them, 17 believed that customization is also crucial in the context of animated text, which helps to “express opinions and feelings more effectively and precisely” (P10), as well as “making the creations more specific and memorable” (P18).

◊ On motion transfer. Users expressed their agreement with the novelty, interest, and inspiration of our approach. Seven participants mentioned “novel”, where P9 commented “it is pleasantly surprising”. Four participants mentioned “interesting”, and 3 participants found it inspiring. P2 found “the generated results inspiring and insightful for my design progress”. P1 said: “I appreciate the smart use of motion transfer, as the interpretation of motion is subjective, but you use memes as an intermediate medium whose ability to convey emotions is validated through wide applications, making the generated results meet subjective expectations. In addition, using text as a vector container for transformation, where the container can be liquefied, allows for slight deformation, thus in support of more subtle emotional expression.”

◊ On application scenarios. Participants brainstormed multiple application scenarios in their personal life with kinetic typography, including online chatting, social media post, website banners, presentations, subtitle enhancement, and E-invitation. P13 also imagined that in order to attract interest, animated text may be used as a teaching instrument for introducing words to little children. P1 identified some challenges in the application of the animated text. “When used alone, the interpretation by users can be ambiguous. There are challenges in accessibility, readability, as well as efficiency of perception and recognition”. He suggested that we can “emphasize the combination with animated images, which can enhance contextual effects and create a synergistic interaction greater than the sum of the separate parts”.

◊ On future improvements. First, Participants suggested a possible enhancement in the guidance of interactions, preferably by introducing recommendations for interactive operations. The operations may be clearer “with the help of some icons and text” (P16), and “the generation efficiency can be improved by recommending interactive behaviors. In addition, it would be very helpful to support

*automatic modifications of other frames after modifying one frame in key point correction and glyph refinement*" (P2). Besides, integrating external resources may enrich the generated results. P11 suggested *"incorporating language models to generate using instructions enhances its convenience"*. P3 commented that *"integrating meme and artistic font libraries helps generate richer and more artistic results"*.

## 9 DISCUSSION

We summarize the implications of our investigation, reflect on our limitations, and discuss promising directions for future research.

### 9.1 Implication

*Create animated effects with model-free motion transfer.* We contribute a novel approach in the emerging area of AI-generated content to design motion graphics using prevalent cross-domain GIFs as references. Participants in the workshop acknowledge the ease of guiding animation generation with reference GIFs. Despite various properties to coordinate in animation design, the underlying workflow of Wakey-Wakey helps users to author in a top-down manner instead of tweaking every details. Beyond template engines, model-free motion transfer helps create more diversified effects.

*Support human-AI collaboration with interpretable features.* According to the user feedback in the workshop, the extracted key point helps them understand the causes of misalignment. While most of our participants are unaware of the internal mechanism, they are able to correct the flaws introduced by the black-box model and intervene in the generation process to produce more desirable results. In designing AI-empowered authoring tools, interpretable features are practical entry points for humans to inject requirements into the content creation process.

*Consider design requirements of users at different levels.* Wakey-Wakey supports both one-off generation and fine-grain control. On the one hand, there are default values underlying algorithms to cater to the fast-generation need of causal users. On the other hand, configurable parameters and the vector representation of generated results are also exposed for further adjustment. When developing authoring support for prevailing artifacts, such as kinetic typography or data visualization, it is important to consider the requirements of different user profiles to make the authoring tool more useful.

### 9.2 Limitation

*Deformation stability.* When the motion amplitude of the character in the driving GIF is too large, excessive deformation may occur in the glyph, especially for fonts with only a few control points. Severe deformation can result in distorted glyphs with discontinued outlines and low legibility. This may be addressed by expanding the number of control points along the predefined glyph outline [23], using the triangulated mesh representation of glyphs [12], and introducing global penalty in the loss function. In addition, as discussed in Section 7.2, the automatic pipeline may fail for over-complicated GIF styles or typefaces, which requires more generalized models in motion trajectory extraction.

*Motion semantic perseverance.* As with other cross-domain motion transfer problems, when generated result may not preserve

the original semantics, as text generally lacks comparable internal structures with GIFs, where the length of a text strongly influences the success. With our approach, a "goodbye" may crawl like a snake but hardly a giraffe swinging its neck. In addition, our objective function prioritizes the overall deformation of the exterior and may neglect the local and independent deformation of the interior. For instance, the generated result may learn the waving gesture but miss the delicate eye blinking.

*Animation expressiveness.* We leverage motion transfer on text control points to deform the shape of text elements and mimic the general animated effects in the driving GIF. However, in addition to learning the shape-deformation patterns, kinetic typography also concerns other properties [60]. Future works may explore incorporating visual properties like colors and designing an integrated environment for a more flexible authoring experience.

### 9.3 Future Work

This work demonstrates a novice-friendly approach to creating text animation through cross-domain motion transfer. While we focus on texts, it is also interesting to explore arbitrary anthropomorphized shapes, such as the dancing mushrooms in Disney's *Fantasia* [51] or sketches of monsters [47, 48]. Unlike human postures constrained by bones and flesh, motion graphics enjoy higher flexibility for exaggerating effects, which share similarities with text. We note that the outline optimization should shift the focus from maintaining text legibility to shape semantics. Recent advances in large language-vision models like Stable Diffusion [42] may help to regularize undesired artifacts. In addition, animated data storytelling (e.g., [43, 57]) remains an exciting avenue for integrating expressive motions on visual marks. As illustrated in the case of an animated word cloud (see Figure 12), the positions of individual visual marks in a visualization can be regarded as the control points of a text. Using our approach, it is possible to transfer motion into visualizations, which may extend existing visual vocabularies and further facilitate the comprehension of abstract data and the expression of emotions [29, 60].

## 10 CONCLUSION

In this study, we explore the opportunity to create expressive kinetic typography based on a driving GIF with character motions. Based on the unique characteristics of text, we propose a framework that adapts existing motion transfer models to the vector domain. Specifically, we animate text based on their control points predefined in font specification. To mimic the motion while maintaining fair legibility, the by-frame positions of each control point are regularized by the extracted motion key point in the driving GIF and neighboring control points. We also introduce an interaction module that allows human-computer collaboration, where humans can steer the intermediate results and guide the generation of kinetic typography. Based on the framework, we developed a mixed-initiative authoring tool and a browser widget featuring automatic generation. A questionnaire study (N=33) initially validated the effectiveness of our approach. Participants generally recognized that the results were animated in a desirable manner, both aesthetically pleasing and semantically resonant. Moreover, we evaluated the novel transfer-based kinetic typography tools by organizing a

workshop (N=20) with both general users and professional designers. People were positive about the interactive system and showed interest in employing our tools for various scenarios.

## ACKNOWLEDGMENTS

This research was supported by the Natural Science Foundation of China (NSFC No.62202105), Shanghai Municipal Science and Technology (No. 21ZR1403300 and No. 21YF1402900), and Hong Kong Research Grants Council General Research Fund 16210722. We thank the anonymous reviewers, participants in the user studies, Zhan Wang, Ziyue Lin, Dr. Xinhuan Shu, Dr. Jiaxiong Hu, and Prof. Zhenjie Zhao for valuable feedback.

## REFERENCES

- [1] Adele. 2012. *Adele - Skyfall (Official Lyric Video)*. Retrieved May 30, 2023 from <https://youtu.be/DeumyOzKqgl> 1:37–1:53.
- [2] Adobe Inc. 2023. *After Effects*. Retrieved Mar 20, 2023 from <https://www.adobe.com/products/aftereffects.html>
- [3] Toshiki Aoki, Rintaro Chujo, Katsufumi Matsui, Saemi Choi, and Ari Hautasaari. 2022. EmoBalloon—Conveying Emotional Arousal in Text Chats with Speech Balloons. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*. ACM, New York, NY, USA, Article 527, 16 pages. <https://doi.org/10.1145/3491102.3501920>
- [4] Apple Inc. 2023. *Motion*. Retrieved Mar 20, 2023 from <https://www.apple.com/final-cut-pro/motion/>
- [5] Rahul Arora, Rubaiat Habib Kazi, Danny M Kaufman, Wilmot Li, and Karan Singh. 2019. Magicalhands: Mid-Air Hand Gestures for Animating in VR. In *Proceedings of the ACM Annual Symposium on User Interface Software and Technology (UIST)*. ACM, New York, NY, USA, 463–477. <https://doi.org/10.1145/3332165.3347942>
- [6] Alberto Betella and Paul FMJ Verschure. 2016. The Affective Slider: A Digital Self-assessment Scale for the Measurement of Human Emotions. *PLoS one* 11, 2, Article e0148037 (2016), 11 pages. <https://doi.org/10.1371/journal.pone.0148037>
- [7] George Borzyskowski. 2004. Animated Text: More than Meets the Eye?. In *Proceedings of the ASCILITE Conference*. 141–144. <https://www.ascilite.org/conferences/perth04/procs/pdf/borzyskowski.pdf>
- [8] Caroline Chan, Shiry Ginosar, Tinghui Zhou, and Alexei A Efros. 2019. Everybody Dance Now. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, Piscataway, NJ, USA, 5933–5942. <https://doi.org/10.1109/ICCV.2019.00603>
- [9] Bay-Wei Chang and David Ungar. 1993. Animation: From Cartoons to the User Interface. In *Proceedings of the ACM Annual Symposium on User Interface Software and Technology (UIST)*. ACM, New York, NY, USA, 45–55. <https://doi.org/10.1145/168642.168647>
- [10] Yung-Yu Chuang, Dan B Goldman, Ke Colin Zheng, Brian Curless, David H. Salesin, and Richard Szeliski. 2005. Animating Pictures with Stochastic Motion Textures. *ACM Transactions on Graphics* 24, 3 (2005), 853–860. <https://doi.org/10.1145/1073204.1073273>
- [11] O. Conolly. 2003. Aesthetic Principles. *The British Journal of Aesthetics* 43 (04 2003), 114–125. <https://doi.org/10.1093/bjaesthetics/43.2.114>
- [12] Mathieu Desbrun, Mark Meyer, Peter Schröder, and Alan H Barr. 1999. Implicit Fairing of Irregular Meshes Using Diffusion and Curvature Flow. In *Proceedings of the Annual Conference on Computer Graphics and Interactive Techniques (SIGGRAPH)*. ACM, New York, NY, USA, 317–324. <https://doi.org/10.1145/311535.311576>
- [13] Laura Devendorf and Kimiko Ryokai. 2013. AnyType: Provoking Reflection and Exploration with Aesthetic Interaction. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*. ACM, New York, NY, USA, 1041–1050. <https://doi.org/10.1145/2470654.2466133>
- [14] Marek Dvorožňák, Pierre Bénéard, Pascal Barla, Oliver Wang, and Daniel Sýkora. 2017. Example-Based Expressive Animation of 2D Rigid Bodies. *ACM Transactions on Graphics* 36, 4, Article 127 (2017), 10 pages. <https://doi.org/10.1145/3072959.3073611>
- [15] Paul Ekman. 1999. Basic Emotions. *Handbook of Cognition and Emotion* 98, 45–60 (1999), 16. <https://doi.org/10.1002/0470013494.ch3>
- [16] Shannon Ford, Jodi Forlizzi, and Soguro Ishizaki. 1997. Kinetic Typography: Issues in Time-Based Presentation of Text. In *Extended Abstracts on Human Factors in Computing Systems (CHI EA)*. ACM, New York, NY, USA, 269–270. <https://doi.org/10.1145/1120212.1120387>
- [17] Jodi Forlizzi, Johnny Lee, and Scott Hudson. 2003. The Kinedit System: Affective Messages Using Dynamic Texts. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*. ACM, New York, NY, USA, 377–384. <https://doi.org/10.1145/642611.642677>
- [18] Leon A Gatys, Alexander S Ecker, and Matthias Bethge. 2016. Image Style Transfer Using Convolutional Neural Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Piscataway, NJ, USA, 2414–2423. <https://doi.org/10.1109/CVPR.2016.265>
- [19] Weston Gaylord, Vivian Hare, and Ashley Ngu. 2015. Adding Body Motion and Intonation to Instant Messaging with Animation. In *Adjunct Proceedings of the ACM Symposium on User Interface Software and Technology*. ACM, New York, NY, USA, 105–106. <https://doi.org/10.1145/2815585.2815741>
- [20] Weston Gaylord, Vivian Hare, and Ashley Ngu. 2015. Adding Body Motion and Intonation to Instant Messaging with Animation. In *Adjunct Proceedings of the ACM Symposium on User Interface Software & Technology (UIST Adjunct)*. ACM, New York, NY, USA, 105–106. <https://doi.org/10.1145/2815585.2815741>
- [21] Tavi Halperin, Hanit Hakim, Orestis Vantzos, Gershon Hochman, Netai Benaim, Lior Sassy, Michael Kupchik, Ofir Bibi, and Ohad Fried. 2021. Endless Loops: Detecting and Animating Periodic Patterns in Still Images. *ACM Transactions on Graphics* 40, 4, Article 142 (2021), 12 pages. <https://doi.org/10.1145/3450626.3459935>
- [22] Fa-Ting Hong, Longhao Zhang, Li Shen, and Dan Xu. 2022. Depth-Aware Generative Adversarial Network for Talking Head Video Generation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Piscataway, NJ, USA, 3397–3406. <https://doi.org/10.1109/CVPR52688.2022.00339>
- [23] Shir Iluz, Yael Vinker, Amir Hertz, Daniel Berio, Daniel Cohen-Or, and Ariel Shamir. 2023. Word-As-Image for Semantic Typography. arXiv:2303.01818 Accepted in ACM SIGGRAPH 2023.
- [24] Jun Kato, Tomoyasu Nakano, and Masataka Goto. 2015. TextAlive: Integrated Design Environment for Kinetic Typography. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*. ACM, New York, NY, USA, 3403–3412. <https://doi.org/10.1145/2702123.2702140>
- [25] Rubaiat Habib Kazi, Fanny Chevalier, Tovi Grossman, Shengdong Zhao, and George Fitzmaurice. 2014. Draco: Bringing Life to Illustrations with Kinetic Textures. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*. ACM, New York, NY, USA, 351–360. <https://doi.org/10.1145/2556288.2556987>
- [26] Rubaiat Habib Kazi, Tovi Grossman, Nobuyuki Umetani, and George Fitzmaurice. 2016. Motion Amplifiers: Sketching Dynamic Illustrations Using the Principles of 2D Animation. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*. ACM, New York, NY, USA, 4599–4609. <https://doi.org/10.1145/2858036.2858386>
- [27] Minhwan Kim, Kyungah Choi, and Hyeon-Jeong Suk. 2016. Yo! Enriching Emotional Quality of Single-Button Messengers through Kinetic Typography. In *Proceedings of the ACM Conference on Designing Interactive Systems (DIS)*. ACM, New York, NY, USA, 276–280. <https://doi.org/10.1145/2901790.2901835>
- [28] Yu-Chi Lai, Bo-An Chen, Kuo-Wei Chen, Wei-Lin Si, Chih-Yuan Yao, and Eugene Zhang. 2016. Data-driven NPR Illustrations of Natural Flows in Chinese Painting. *IEEE Transactions on Visualization and Computer Graphics* 23, 12 (2016), 2535–2549. <https://doi.org/10.1109/TVCG.2016.2622269>
- [29] Xingyu Lan, Yang Shi, Yanqiu Wu, Xiaohan Jiao, and Nan Cao. 2021. Kineticcharts: Augmenting Affective Expressiveness of Charts in Data Stories with Animation Design. *IEEE Transactions on Visualization and Computer Graphics* 28 (2021), 933–943. <https://doi.org/10.1109/TVCG.2021.3114775>
- [30] Daniel G. Lee, Deborah I. Fels, and John Patrick Udo. 2007. Emotive Captioning. *Computers in Entertainment* 5, 2, Article 11 (2007), 15 pages. <https://doi.org/10.1145/1279540.1279551>
- [31] Johnny C. Lee, Jodi Forlizzi, and Scott E. Hudson. 2002. The Kinetic Typography Engine: An Extensible System for Animating Expressive Text. In *Proceedings of the ACM Annual Symposium on User Interface Software and Technology (UIST)*. ACM, New York, NY, USA, 81–90. <https://doi.org/10.1145/571985.571997>
- [32] Sooyeon Lim. 2022. A Study on the Interactive Expression of Human Emotions in Typography. *International Journal of Advanced Culture Technology* 10, 1 (2022), 122–130. <https://doi.org/10.17703/IJACT.2022.10.1.122>
- [33] lipsum.com. 1996. *Lorem Ipsum*. Retrieved Mar 20, 2023 from <https://www.lipsum.com/>
- [34] Sabrina Malik, Jonathan Aitken, and Judith Kelly Waalen. 2009. Communicating Emotion with Animated Text. *Visual Communication* 8, 4 (2009), 469–479. <https://doi.org/10.1177/1470357209343375>
- [35] Wendong Mao, Shuai Yang, Huihong Shi, Jiaying Liu, and Zhongfeng Wang. 2022. Intelligent Typography: Artistic Text Style Transfer for Complex Texture and Structure. *IEEE Transactions on Multimedia* (2022). <https://doi.org/10.1109/TMM.2022.3209870> Early Access.
- [36] Yifang Men, Zhouhui Lian, Yingmin Tang, and Jianguo Xiao. 2019. DynTypo: Example-based Dynamic Text Effects Transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. IEEE, Piscataway, NJ, USA, 5870–5879. <https://doi.org/10.1109/CVPR.2019.00602>
- [37] Mitsuru Minakuchi and Yutaka Kidawara. 2008. Kinetic Typography for Ambient Displays. In *Proceedings of the International Conference on Ubiquitous Information Management and Communication (ICUIMC)*. ACM, New York, NY, USA, 54–57. <https://doi.org/10.1145/1352793.1352805>

- [38] Mitsuru Minakuchi and Katsumi Tanaka. 2005. Automatic Kinetic Typography Composer. In *Proceedings of the ACM Conference on Advances in Computer Entertainment Technology (ACE)*. ACM, New York, NY, USA, 221–224. <https://doi.org/10.1145/1178477.1178512>
- [39] Moonlab Studio Co., Ltd. 2023. *Puppy Maltese*. Retrieved Mar 20, 2023 from <https://weibo.com/u/7776232700?tabtype=home>
- [40] Laurence Penny. 1996. *A History of TrueType*. Retrieved Mar 20, 2023 from <https://www.true-type.com>
- [41] Huy Quoc Phan, Hongbo Fu, and Antoni B Chan. 2015. Flexyfont: Learning Transferring Rules for Flexible Typeface Synthesis. *Computer Graphics Forum* 34, 7 (2015), 245–256. <https://doi.org/10.1111/cgf.12763>
- [42] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis with Latent Diffusion Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Piscataway, NJ, USA, 10684–10695. <https://doi.org/10.1109/CVPR52688.2022.01042>
- [43] Xinhuan Shu, Aoyu Wu, Junxiu Tang, Benjamin Bach, Yingcai Wu, and Huamin Qu. 2021. What Makes a Data-GIF Understandable? *IEEE Trans. Vis. Comput. Graph.* 27, 02 (2021), 1492–1502. <https://doi.org/10.1109/TVCG.2020.3030396>
- [44] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019. Animating Arbitrary Objects via Deep Motion Transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Piscataway, NJ, USA, 2377–2386. <https://doi.org/10.1109/CVPR.2019.00248>
- [45] Aliaksandr Siarohin, Stéphane Lathuilière, Sergey Tulyakov, Elisa Ricci, and Nicu Sebe. 2019. First Order Motion Model for Image Animation. In *Proceedings of the International Conference on Neural Information Processing Systems (NeurIPS)*. Curran Associates Inc., Red Hook, NY, USA, Article 641, 11 pages. <https://doi.org/10.5555/3454287.3454928>
- [46] Aliaksandr Siarohin, Oliver J Woodford, Jian Ren, Menglei Chai, and Sergey Tulyakov. 2021. Motion Representations for Articulated Animation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Piscataway, NJ, USA, 13653–13662. <https://doi.org/10.1109/CVPR46437.2021.01344>
- [47] Harrison Jesse Smith, Qingyuan Zheng, Yifei Li, Somya Jain, and Jessica K. Hodgins. 2023. A Method for Animating Children’s Drawings of the Human Figure. *ACM Transactions on Graphics* 42, 3, Article 32 (2023), 15 pages. <https://doi.org/10.1145/3592788>
- [48] Qingkun Su, Xue Bai, Hongbo Fu, Chiew-Lan Tai, and Jue Wang. 2018. Live Sketch: Video-Driven Dynamic Deformation of Static Drawings. In *Proceedings of the ACM Conference on Human Factors in Computing Systems (CHI)*. Association for Computing Machinery, New York, NY, USA, Article 662, 12 pages. <https://doi.org/10.1145/3173574.3174236>
- [49] Jiale Tao, Biao Wang, Borun Xu, Tiezheng Ge, Yuning Jiang, Wen Li, and Lixin Duan. 2022. Structure-Aware Motion Transfer with Deformable Anchor Model. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Piscataway, NJ, USA, 3637–3646. <https://doi.org/10.1109/CVPR52688.2022.00362>
- [50] Purva Tendulkar, Kalpesh Krishna, Ramprasaath R Selvaraju, and Devi Parikh. 2019. Trick or TReAT: Thematic Reinforcement for Artistic Typography. In *Proceedings of the International Conferences on Computational Creativity (ICCC)*. ACC, 9 pages. arXiv:1903.07820
- [51] The Walt Disney Company. 1940. *Fantasia*. Retrieved Mar 20, 2023 from <https://youtu.be/r7gLLv4it0> 17:03–18:07.
- [52] Frank Thomas, Ollie Johnston, and Frank Thomas. 1995. *The Illusion of Life: Disney Animation*. Hyperion New York.
- [53] Fernanda B. Viégas, Martin Wattenberg, and Jonathan Feinberg. 2009. Participatory Visualization with Wordle. *IEEE Transactions on Visualization and Computer Graphics* 15, 6 (2009), 1137–1144. <https://doi.org/10.1109/TVCG.2009.171>
- [54] Quoc V Vy, Jorge A Mori, David W Fourney, and Deborah I Fels. 2008. EnACT: A Software Tool for Creating Animated Text captions. In *Proceedings of the International Conference on Computers Helping People with Special Needs (ICHP)*. Springer, Berlin, Heidelberg, 609–616. [https://doi.org/10.1007/978-3-540-70540-6\\_87](https://doi.org/10.1007/978-3-540-70540-6_87)
- [55] Hua Wang, Helmut Prendinger, and Takeo Igarashi. 2004. Communicating Emotions in Online Chat Using Physiological Sensors and Animated Text. In *Extended Abstracts on Human Factors in Computing Systems (CHI EA)*. ACM, New York, NY, USA, 1171–1174. <https://doi.org/10.1145/985921.986016>
- [56] Wenjing Wang, Jiaying Liu, Shuai Yang, and Zongming Guo. 2019. Typography with Decor: Intelligent Text Style Transfer. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, Piscataway, NJ, USA, 5889–5897. <https://doi.org/10.1109/CVPR.2019.00604>
- [57] Yun Wang, Yi Gao, Ray Huang, Weiwei Cui, Haidong Zhang, and Dongmei Zhang. 2021. Animated Presentation of Static Infographics with InfoMotion. *Computer Graphics Forum* 40, 3 (2021), 507–518. <https://doi.org/10.1111/cgf.14325>
- [58] Nora S. Willett, Rubaiat Habib Kazi, Michael Chen, George Fitzmaurice, Adam Finkelstein, and Tovi Grossman. 2018. A Mixed-Initiative Interface for Animating Static Pictures. In *Proceedings of the ACM Annual Symposium on User Interface Software and Technology (UIST)*. ACM, New York, NY, USA, 649–661. <https://doi.org/10.1145/3242587.3242612>
- [59] Nora S Willett, Hijung Valentina Shin, Zeyu Jin, Wilmot Li, and Adam Finkelstein. 2020. Pose2Pose: Pose Selection and Transfer for 2D Character Animation. In *Proceedings of the 25th International Conference on Intelligent User Interfaces (IUI)*. ACM, New York, NY, USA, 88–99. <https://doi.org/10.1145/3377325.3377505>
- [60] Liwenhan Xie, Xinhuan Shu, Jeon Cheol Su, Yun Wang, Siming Chen, and Huamin Qu. 2023. Creating Emordle: Animating Word Cloud for Emotion Expression. *IEEE Transactions on Visualization and Computer Graphics* (2023). <https://doi.org/10.1109/TVCG.2023.3286392> Early Access.
- [61] Jun Xing, Rubaiat Habib Kazi, Tovi Grossman, Li-Yi Wei, Jos Stam, and George Fitzmaurice. 2016. Energy-Brushes: Interactive Tools for Illustrating Stylized Elemental Dynamics. In *Proceedings of the ACM Annual Symposium on User Interface Software and Technology (UIST)*. ACM, New York, NY, USA, 755–766. <https://doi.org/10.1145/2984511.2984585>
- [62] Borun Xu, Biao Wang, Jinhong Deng, Jiale Tao, Tiezheng Ge, Yuning Jiang, Wen Li, and Lixin Duan. 2022. Motion and Appearance Adaptation for Cross-Domain Motion Transfer. In *Proceedings of the European Conference on Computer Vision (ECCV), Part XVI*. Springer, Cham, Germany, 529–545. [https://doi.org/10.1007/978-3-031-19787-1\\_40](https://doi.org/10.1007/978-3-031-19787-1_40)
- [63] Jie Xu and Craig S Kaplan. 2007. Calligraphic Packing. In *Proceedings of Graphics Interface (GI)*. ACM, New York, NY, USA, 43–50. <https://doi.org/10.1145/1268517.1268527>
- [64] Xuemiao Xu, Liang Wan, Xiaopei Liu, Tien-Tsin Wong, Liansheng Wang, and Chi-Sing Leung. 2008. Animating Animal Motion from Still. In *Proceedings of the Special Interest Group on Computer Graphics and Interactive Techniques (SIGGRAPH Asia)*. ACM, New York, NY, USA, Article 117, 8 pages. <https://doi.org/10.1145/1457515.1409070>
- [65] Shuai Yang, Zhangyang Wang, and Jiaying Liu. 2021. Shape-Matching GAN++: Scale Controllable Dynamic Artistic Text Style Transfer. *IEEE Transactions on Pattern Analysis and Machine Intelligence* 44, 7 (2021), 3807–3820. <https://doi.org/10.1109/TPAMI.2021.3055211>
- [66] Zhiquan Ye. 2008. Emotional Instant Messaging with KIM. In *Extended Abstracts on Human Factors in Computing Systems (CHI EA)*. ACM, New York, NY, USA, 3729–3734. <https://doi.org/10.1145/1358628.1358921>
- [67] Junsong Zhang, Yu Wang, Weiyi Xiao, and Zhenshan Luo. 2017. Synthesizing Ornamental Typefaces. *Computer Graphics Forum* 36, 1 (2017), 64–75. <https://doi.org/10.1111/cgf.12785>
- [68] Yang Zhou, Xintong Han, Eli Shechtman, Jose Echevarria, Evangelos Kalogerakis, and Dingzeyu Li. 2020. MakeltTalk: Speaker-Aware Talking-Head Animation. *ACM Transactions on Graphics* 39, 6, Article 221 (2020), 15 pages. <https://doi.org/10.1145/3414685.3417774>
- [69] Changqing Zou, Junjie Cao, Warunika Ranaweera, Ibraheem Alhashim, Ping Tan, Alla Sheffer, and Hao Zhang. 2016. Legible Compact Calligrams. *ACM Transactions on Graphics* 35, 4, Article 122 (2016), 12 pages. <https://doi.org/10.1145/2897824.2925887>
- [70] Ebberts + Zucker. 2023. *TypeMonkey*. Retrieved Mar 20, 2023 from <http://aescrpts.com/typemonkey/>