

DIFFANT: Diffusion Models for Action Anticipation

Zeyun Zhong^{1,2} Chengzhi Wu¹ Manuel Martin² Michael Voit² Juergen Gall³ Jürgen Beyerer^{1,2}
¹Karlsruhe Institute of Technology ²Fraunhofer IOSB ³University of Bonn
{firstname.lastname}@{kit.edu, iosb.fraunhofer.de} gall@iai.uni-bonn.de

Abstract

Anticipating future actions is inherently uncertain. Given an observed video segment containing ongoing actions, multiple subsequent actions can plausibly follow. This uncertainty becomes even larger when predicting far into the future. However, the majority of existing action anticipation models adhere to a deterministic approach, neglecting to account for future uncertainties. In this work, we rethink action anticipation from a generative view, employing diffusion models to capture different possible future actions. In this framework, future actions are iteratively generated from standard Gaussian noise in the latent space, conditioned on the observed video, and subsequently transitioned into the action space. Extensive experiments on four benchmark datasets, i.e., Breakfast, 50Salads, EpicKitchens, and EGTEA Gaze+, are performed and the proposed method achieves superior or comparable results to state-of-the-art methods, showing the effectiveness of a generative approach for action anticipation. Our code and trained models will be published on GitHub.

1. Introduction

In contexts such as human-machine cooperation and robotic assistance, the anticipation of potential future daily-living actions is vital. For instances where timely assistance is needed or proactive dialogues are essential, the capability to accurately predict actions, even in the absence of direct observation, is of utmost importance. Yet, the ever-changing landscape of daily activities introduces a natural unpredictability to future actions, as depicted in Fig. 1. This inherent uncertainty becomes even larger if we are going to predict far into the future, which poses significant challenges to the precise anticipation of future actions. Consequently, modeling the underlying uncertainty may be beneficial, allowing to capture different possible futures.

Addressing the intricate challenges and uncertainties of action anticipation calls for a departure from conventional methods. Traditional approaches [13, 18, 27, 45] predominantly employ a discriminative and deterministic stance,

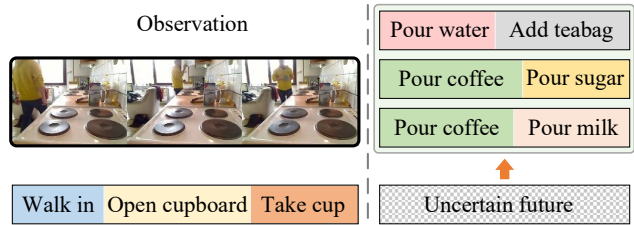


Figure 1. The inherent uncertainty in predicting future actions presents a complex scenario where, given a single observation, multiple feasible future action sequences may emerge. Example frames are from the Breakfast [30] dataset.

framing action prediction as a classification task with deterministic future outcomes. This approach, however, can inadvertently diminish the visibility of plausible alternative actions [60]. Although probabilistic models [1, 38, 60], including generative models like GANs [60] and VAEs [38], have been introduced to embrace the non-deterministic essence of future anticipations, they rely on action labels from observed frames, thereby constraining their immediate applicability. As such, a model that acknowledges inherent uncertainties and supports direct deployment is of paramount importance.

Within this context, we rethink action anticipation from a generative view, introducing an end-to-end probabilistic generative model, DIFFANT. Leveraging diffusion models [23, 47], which have emerged as a paradigm with significant promise across various domains, our approach iteratively generates future actions in the latent space from the standard Gaussian noise. To adapt the diffusion models for action anticipation, we extend the standard models by introducing a future action embedding function and an action predictor. Moreover, to enhance the precision of our predictions, we incorporate visual observations as conditional information and employ an encoder-decoder structure [5, 18, 53] for seamless integration. We evaluate DIFFANT on four standard benchmarks for long-term action anticipation, achieving state-of-the-art results on Breakfast, 50Salads, and EGTEA Gaze+, and comparable results on EpicKitchens. In summary, our main contributions are:

- A probabilistic generative approach, DIFFANT, for long-term action anticipation, employing diffusion models and leveraging their stochastic and continuous nature to iteratively refine predictions and navigate the intrinsic uncertainties that exist in predicting future actions.
- We conduct extensive experiments on four widely used benchmark datasets, *i.e.*, Breakfast, 50Salads, EpicKitchens, and EGTEA Gaze+. Our results consistently demonstrate superior or comparable performance to state-of-the-art deterministic methods.
- We present an in-depth analysis of our approach, demonstrating that our model surpasses the current state-of-the-art probabilistic method.

2. Related Work

Action Anticipation aims to predict future actions given a video clip of the past and present. Many approaches initially investigated different forms of action and activity anticipation from third person video [13, 15, 18, 27]. Recently, along with development of multiple challenge benchmarks [8, 9, 20, 34], the first-person (egocentric) vision has also gained popularity. To accurately predict future actions, the summarization of temporal progression of past actions is essential. To model the past action progression, earlier methods mainly used RNN [13, 14] or TCN [27]-based architectures, which have been shown to be inferior to the recent Transformer-based approaches [16, 18, 40, 62]. Based on the predicted time horizon, action anticipation approaches can be broadly grouped into two categories [61]: short-term anticipation approaches [8, 9] and long-term anticipation approaches [13, 20]. While short-term approaches predict actions a few seconds into the future, long-term approaches aim to predict a sequence of future actions (with their durations) up to several minutes into the future.

Long-term Anticipation. An initial work [13] introduced two distinct models for long-term action anticipation. While the RNN model performed in a recursive manner, the CNN model outputs a sequence of future actions in the form of a matrix in one single step. Ke *et al.* [27] developed a model targeting the prediction of a specific future action, bypassing the need for anticipating intermediate actions to avoid error accumulations. Farha *et al.* [2] introduced a cycle consistency module to predict past activities using the projected future, demonstrating improved outcomes in comparison to its counterpart lacking the consistency module. Sener *et al.* [45] suggested a multi-scale temporal aggregation model that aggregates past visual features in condensed vectors and then iteratively predicts future actions using an LSTM. Recently, Transformer-based approaches [18, 40] have been also employed for long-term anticipation. Different from these approaches, our method is non-deterministic and is capable to take the uncertain future into account.

Uncertainty-aware Anticipation. To address the inherent uncertainty involved in forecasting future actions, various researches have proposed non-deterministic approaches, including generating an array of potential outcomes with multiple rules [42, 55] and by sampling from the learned distribution [1, 37, 38, 60]. Vondrick *et al.* [55] proposed training a mixture of networks, each aiming to predict one potential future. Piergiovanni *et al.* [42] proposed a differentiable grammar model and applied adversarial techniques to enable efficient learning and avoid enumerating all possible rules. Farha and Gall [1] anticipated all subsequent actions and durations stochastically, employing an action model and a time model trained to predict the probability distribution of future action labels and durations. In the literature, probabilistic generative models like GANs [60] and VAEs [38] have been employed to capture the uncertain aspects of future predictions. Nonetheless, these methods typically require the action labels of observed frames as inputs, which are obtained either from ground truth labels or inferred through an action segmentation model [11]. Drawing inspiration from these works, we approach the anticipation task from a generative perspective, employing diffusion models to generate diverse and plausible future actions, while leveraging the rich visual features.

Diffusion Models. Initially unified with score-based models [48–50], diffusion models [23, 46, 47] are renowned for their stable training processes, which do not rely on adversarial mechanism for generative learning. These models have demonstrated remarkable achievements across various domains, including image generation [10, 44, 56], natural language generation [19, 33, 59], text-to-image synthesis [21, 28], and audio generation [31, 32]. Recent advances have extended the application of diffusion models to image and video comprehension tasks within computer vision, such as object detection [7], image segmentation [3, 4], video forecasting and infilling [24, 54, 57], and action segmentation [35]. In this work, we leverage the iterative refinement capabilities of diffusion models for future prediction. To the best of our knowledge, this work is the first one employing diffusion models for action anticipation.

3. Diffusion Action Anticipation

To address the inherently non-deterministic nature of future action anticipation, we present DIFFANT, a novel diffusion model for action anticipation. The model integrates an action embedding function and an action predictor into the original diffusion framework, enabling the incorporation of discrete future actions within both the forward and reverse diffusion processes (see Fig. 2). To refine the predictive capabilities of DIFFANT, an encoder-decoder architecture [5, 18, 53] is utilized to assimilate visual observations as conditional information. We acquaint the reader

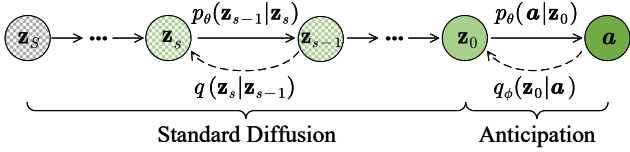


Figure 2. Concept of DIFFANT. Standard diffusion models reconstruct original data \mathbf{z}_0 from pure noise \mathbf{z}_S . Our approach integrates discrete future actions \mathbf{a} by introducing an embedding function and a predictor to facilitate the conversion between continuous \mathbf{z}_0 and discrete \mathbf{a} .

with the problem statement and the background of diffusion models in Section 3.1. Subsequently, the inference and training processes of the proposed method are introduced in Section 3.2 and 3.3, respectively.

3.1. Preliminaries

Problem Statement. Given the past observation containing observed video features $F \in \mathbb{R}^{L \times K}$ with K dimensions for L frames, the long-term anticipation task aims to predict following N future actions \mathbf{a} within a video, consisting of action classes $\mathbf{a}^c \in \mathbb{R}^{N \times C}$ and durations $\mathbf{a}^t \in \mathbb{R}^{N \times 1}$, where C is the total number of action classes.

Diffusion Models aim to approximate the data distribution $q(\mathbf{z}_0)$ with a model distribution $p_\theta(\mathbf{z}_0)$ [23, 47]. A diffusion model typically contains forward and reverse processes. The forward process or diffusion process corrupts the real data $\mathbf{z}_0 \sim q(\mathbf{z}_0)$ into a series of noisy data $\mathbf{z}_1, \mathbf{z}_2, \dots$ and finally into a standard Gaussian noise $\mathbf{z}_S \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$. For each forward step $s \in [1, 2, \dots, S]$, the perturbation is controlled by $q(\mathbf{z}_s | \mathbf{z}_{s-1}) = \mathcal{N}(\mathbf{z}_s; \sqrt{1 - \beta_s} \mathbf{z}_{s-1}, \beta_s \mathbf{I})$, with predefined $\beta_s \in (0, 1)$ as different variance scales. By denoting $\alpha_s = 1 - \beta_s$ and $\bar{\alpha}_s = \prod_{i=1}^s \alpha_i$, we can directly obtain \mathbf{z}_s from \mathbf{z}_0 in a closed form without recursion: $q(\mathbf{z}_s | \mathbf{z}_0) = \mathcal{N}(\mathbf{z}_s; \sqrt{\bar{\alpha}_s} \mathbf{z}_0, (1 - \bar{\alpha}_s) \mathbf{I})$, which can be further simplified using the reparameterization trick [29].

The reverse process or denoising process tries to gradually remove the noise from $\mathbf{z}_S \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ to reconstruct the original data \mathbf{z}_0 . Each reverse step is defined as $p_\theta(\mathbf{z}_{s-1} | \mathbf{z}_s) = \mathcal{N}(\mathbf{z}_{s-1}; \boldsymbol{\mu}_\theta(\mathbf{z}_s, s), \sigma_s^2 \mathbf{I})$, where σ_s^2 is controlled by β_s , and $\boldsymbol{\mu}_\theta(\mathbf{z}_s, s)$ is a predicted mean parameterized by a step-dependent diffusion model $f_\theta(\mathbf{z}_s, s)$. Several different ways [36] are possible to parameterize p_θ , including the prediction of mean $\boldsymbol{\mu}_\theta(\mathbf{z}_s, s)$, the prediction of the noise ϵ , and the prediction of \mathbf{z}_0 . We choose to predict \mathbf{z}_0 , as suggested in several related works [19, 33, 35].

3.2. Overall Inference Pipeline

To better harness observed video features F as conditional information for controlled future action generation, we propose an encoder-decoder structure depicted in Fig. 3b.

Encoder. The input features F , typically extracted per short

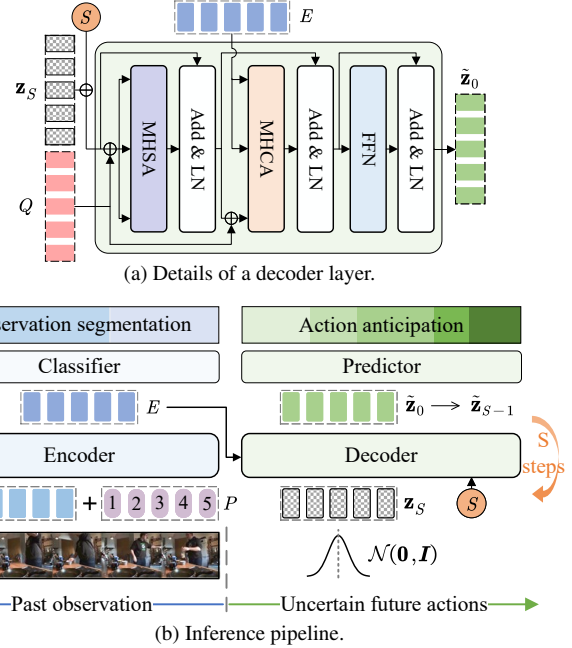


Figure 3. Inference pipeline of DIFFANT. (a) The decoder uses action queries Q and step information to refine the noisy futures \mathbf{z}_S , conditioned on encoded past observations E via cross-attention. (b) Future action embeddings \mathbf{z}_S are drawn from the standard Gaussian distribution and iteratively denoised to produce the final $\tilde{\mathbf{z}}_0$, which is then decoded into future action labels and durations.

clip by pre-trained models, are first processed through an encoder to incorporate long-range temporal context, making them task-oriented. Although our framework permits flexibility in the choice of the encoder, we opt for the standard transformer model used in DETR [5] and FUTR [18]. We begin by passing the video features F through a linear layer, which adjusts them to a suitable hidden dimension D , yielding transformed features F^* . These, combined with sinusoidal positional encodings $P \in \mathbb{R}^{L \times D}$, enable the encoder to generate refined representations $E \in \mathbb{R}^{L \times D}$. Such representations are conducive to action anticipation and can be mapped to the action space using a linear classifier, facilitating observation segmentation.

Decoder. Our decoder, as depicted in Fig. 3a, follows the query-based paradigm [5, 63], originally proposed to eliminate handcrafted components in object detection and recently applied in long-term action anticipation [18, 40]. It processes noisy action embeddings $\mathbf{z}_S \in \mathbb{R}^{M \times D'}$ and action queries $Q \in \mathbb{R}^{M \times D'}$, generating refined futures $\tilde{\mathbf{z}}_0$ in parallel, under the guidance of encoded past observations E via cross-attention. The action queries Q consist of M learnable tokens of dimension $\mathbb{R}^{1 \times D'}$, with their temporal order aligned with the sequence of future actions, *i.e.*, each query directly corresponds to a respective future action [18].

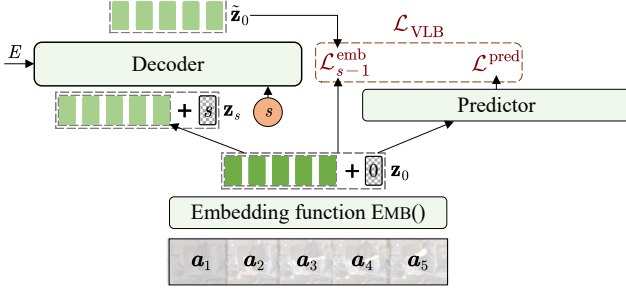


Figure 4. Training pipeline. By employing an embedding function and an action predictor, we integrate the discrete future actions \mathbf{a} into the standard diffusion models. The training optimizes a variational lower bound objective, which combines an embedding reconstruction loss $\mathcal{L}_{s-1}^{\text{emb}}$ and an action prediction loss $\mathcal{L}^{\text{pred}}$.

Predictor. The predictor decodes the refined future embeddings $\tilde{\mathbf{z}}_0$ into action labels $\tilde{\mathbf{a}}^c \in \mathbb{R}^{M \times C}$ and durations $\tilde{\mathbf{a}}^t \in \mathbb{R}^{M \times 1}$ through two dedicated linear layers. To guarantee the non-negativity of the predicted durations, an exponential activation function is employed.

Inference. During inference, future action embeddings \mathbf{z}_S are drawn from a standard Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$. Alternatively, we can initialize \mathbf{z}_S with zero vectors (mean of the standard Gaussian distribution) for deterministic results. The decoder is made step-aware by embedding the current diffusion step $s = S$ with sinusoidal positional encodings, further integrated via a multi-layer perceptron. The decoder inputs \mathbf{z}_S , Q , and step data, aiming to initially rectify the noisy futures and predict $\tilde{\mathbf{z}}_0$, taking cues from the encoded past observations E via cross-attention. For E that does not match the decoder’s dimension D' , a linear transformation adjusts its dimensionality. To enable the iterative reverse diffusion process, we apply the forward chain to get $\tilde{\mathbf{z}}_{S-1}$ and input it with the updated step $s = S - 1$ to the decoder. By iteratively removing noise, a final sample $\tilde{\mathbf{z}}_0$ is generated via a trajectory $\mathbf{z}_S, \tilde{\mathbf{z}}_{S-1}, \dots, \tilde{\mathbf{z}}_0$. To accelerate the reverse process, we adopt DDIM [47] to skip steps in the trajectory. The resulting action embeddings $\tilde{\mathbf{z}}_0$ are decoded into future action labels and durations by the predictor.

Given the preset number of action queries M may surpass the actual number of future actions N , an EOS class is introduced to signify sequence termination. In the inference stage, following [18], all predictions after the first prediction of EOS are discarded. Moreover, we apply Gaussian normalization to the predicted durations to make the sum of all durations equates to one, consistent with the previous work [13, 18].

3.3. Training

To integrate discrete actions \mathbf{a} into the continuous diffusion processes, we extend standard diffusion models with a future action embedding function and an action predictor to

facilitate the conversion between continuous embeddings \mathbf{z}_0 and discrete actions \mathbf{a} , as illustrated in Fig. 4. The following details DIFFANT’s learning process.

Forward Process with Future Actions as Input. DIFFANT takes future actions \mathbf{a} , which comprise action classes $\mathbf{a}^c \in \mathbb{R}^{M \times C}$ and durations $\mathbf{a}^t \in \mathbb{R}^{M \times 1}$, where C is the total number of action classes. It first linearly transforms the one-hot encoded classes and real-value durations into embeddings, $\text{EMB}_c(\mathbf{a}^c)$ and $\text{EMB}_t(\mathbf{a}^t)$, respectively. An additional linear layer then merges these embeddings to form a joint feature space $\text{EMB}(\mathbf{a}^c, \mathbf{a}^t) = f_a([\text{EMB}_c(\mathbf{a}^c), \text{EMB}_t(\mathbf{a}^t)])$. When only action classes are predicted, the model bypasses f_a and directly concatenates the class embeddings. This step facilitates the inclusion of discrete future actions into the diffusion model’s forward process by extending the original forward chain to a new Markov transition $q_\phi(\mathbf{z}_0|\mathbf{a}) = \mathcal{N}(\text{EMB}(\mathbf{a}), \beta_0 \mathbf{I})$, as shown in Fig. 2.

Reverse Process with Conditional Denoising. The reverse process aims to recover the initial state \mathbf{z}_0 from the noised version \mathbf{z}_S by applying the denoising probability $p_\theta(\mathbf{z}_{0:S}) := p(\mathbf{z}_S) \prod_{s=1}^S p_\theta(\mathbf{z}_{s-1}|\mathbf{z}_s)$. Incorporating encoded past observations E as conditional inputs transforms the learning process into $p_\theta(\mathbf{z}_{s-1}|\mathbf{z}_s, E)$, modeled with the proposed DIFFANT $f_\theta(\mathbf{z}_s, s, E)$. For training, in line with [19, 33], we streamline the variational lower bound (\mathcal{L}_{VLB}) to consist of an embedding reconstruction loss and an action prediction loss, $\mathcal{L}_{\text{VLB}} = \sum_{s=1}^S \mathcal{L}_{s-1}^{\text{emb}} + \mathcal{L}^{\text{pred}}$. The embedding loss $\mathcal{L}_{s-1}^{\text{emb}}$ is calculated by:

$$\mathcal{L}_{s-1}^{\text{emb}} = \begin{cases} \|\mathbf{z}_0 - f_\theta(\mathbf{z}_s, s, E)\|^2 & \text{if } 2 \leq s \leq S \\ \|\text{EMB}(\mathbf{a}) - f_\theta(\mathbf{z}_1, 1, E)\|^2 & \text{if } s = 1 \end{cases} \quad (1)$$

and the prediction loss $\mathcal{L}^{\text{pred}}$ is defined as: $\mathcal{L}^{\text{pred}} = -\log p_\theta(\mathbf{a}^c|\mathbf{z}_0, E) - \log p_\theta(\mathbf{a}^t|\mathbf{z}_0, E)$, under the assumption that action classes and durations are conditionally independent for mathematical convenience, following common practices [38]. The action class log-likelihood is computed using cross-entropy and the duration log-likelihood is computed by the mean squared error (MSE) loss. More details are provided in the supplementary material.

Training Objective. In addition to \mathcal{L}_{VLB} for the decoder outputs, we append a classification head to the encoder and apply a cross-entropy loss \mathcal{L}_{seg} and a temporal smoothness loss $\mathcal{L}_{\text{smooth}}$ as auxiliary supervision. The temporal smoothness loss [12] is computed as the mean squared error of the log-likelihoods between adjacent video frames to promote the local similarity along the temporal dimension. The final training objective is thus defined as

$$\mathcal{L} = \mathcal{L}_{\text{seg}} + \mathcal{L}_{\text{smooth}} + \mathcal{L}_{\text{VLB}}. \quad (2)$$

Type	Backbone	Methods	Breakfast β ($\alpha = 0.2$)				Breakfast β ($\alpha = 0.3$)				50Salads β ($\alpha = 0.2$)				50Salads β ($\alpha = 0.3$)			
			0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5
Pred. label	Fisher	RNN [13]	18.11	17.20	15.94	15.81	21.64	20.02	19.73	19.21	30.06	25.43	18.74	13.49	30.77	17.19	14.79	9.77
		CNN [13]	17.90	16.35	15.37	14.54	22.44	20.12	19.69	18.76	21.24	19.03	15.98	9.87	29.14	20.14	17.46	10.86
		UAAA [1]	16.71	15.40	14.47	14.20	20.73	18.27	18.42	16.86	24.86	22.37	19.88	12.82	29.10	20.50	15.28	12.31
		Timecond. [27]	18.41	17.21	16.42	15.84	22.75	20.44	19.64	19.75	32.51	27.61	21.26	15.99	35.12	27.05	22.05	15.59
Features	Fisher	CNN [13]	12.78	11.62	11.21	10.27	17.72	16.87	15.48	14.09	–	–	–	–	–	–	–	–
		TempAgg [45]	15.60	13.10	12.10	11.10	19.50	17.00	15.60	15.10	25.50	19.90	18.20	15.10	30.60	22.50	19.10	11.20
	I3D	TempAgg [45]	24.20	21.10	20.00	18.10	30.40	26.30	23.80	21.20	–	–	–	–	–	–	–	–
		Cycle Cons[2]	25.88	23.42	22.42	21.54	29.66	27.37	25.58	25.20	34.76	28.41	21.82	15.25	34.39	23.70	18.95	15.89
		A-ACT [22]	<u>26.70</u>	24.30	<u>23.20</u>	21.70	30.80	28.30	26.10	25.80	35.40	<u>29.60</u>	22.50	16.10	35.70	<u>25.30</u>	20.10	<u>16.30</u>
		FUTR [18]	27.70	<u>24.55</u>	22.83	<u>22.04</u>	32.27	<u>29.88</u>	<u>27.49</u>	<u>25.87</u>	39.55	27.54	<u>23.31</u>	<u>17.77</u>	<u>35.15</u>	24.86	<u>24.22</u>	15.26
		DIFFANT	25.33	24.59	24.39	22.74	<u>32.13</u>	31.83	31.18	30.77	<u>36.13</u>	34.00	30.46	25.29	34.09	30.14	26.34	20.23

Table 1. Comparison with state-of-the-art methods on Breakfast [30] and 50Salads [51] using MoC (%). Bold and underlined numbers indicate the highest and second highest accuracy, respectively. DIFFANT achieves the state-of-the-art performance in almost all settings.

4. Experimental Setup

Datasets. The Breakfast [30] dataset comprises 1,712 videos of 52 different individuals making breakfast in 18 different kitchens, totalling 77 hours. Every video is categorized into one of the 10 activities related to breakfast preparation. The videos are annotated by 48 fine-grained actions. The 50Salads [51] dataset comprises 50 top-view videos of 25 people preparing a salad. The dataset contains over 4 hours of RGB-D video data, annotated with 17 fine-grained action labels and 3 high-level activities. EpicKitchens [8] and EGTEA Gaze+ [34] are egocentric datasets. EpicKitchens contains 39,596 segments labeled with 125 verbs, 352 nouns, and 2,513 combinations (actions), totalling 55 hours. EGTEA Gaze+ contains 28 hours of videos including 10.3K action annotations, 19 verbs, 51 nouns, and 106 unique actions.

Metrics. For both Breakfast and 50Salads datasets, we calculate the mean accuracy over classes, averaged across all future timestamps within a specified anticipation duration. We observe the first portion (α) of a video, adhering to benchmarks from [13] that set α values at 0.2 or 0.3. Subsequent predictions cover segments β of the entire video, where β can be one of $\{0.1, 0.2, 0.3, 0.5\}$. For our evaluation metrics, we average results across 4 splits for the Breakfast dataset and 5 splits for the 50Salads dataset. For the EpicKitchens and EGTEA Gaze+ datasets, we employ a multi-label classification metric (mAP) targeting specific action classes. A portion α of each untrimmed video serves as input to forecast all subsequent action classes, representing the remaining $(1 - \alpha)$ duration of the video. As in [39], we use $\alpha = \{0.25, 0.50, 0.75\}$ for evaluation. In alignment with [39], we also report the mAP metrics in both low-shot (rare) and many-shot (freq) scenarios.

Architecture Details. The encoder has four layers with a

hidden dimension $D = 256$. The decoder has 8 layers for 50Salads and 4 layers for the other datasets. The decoder dimensions D' are set at 1024 for Breakfast and EpicKitchens, 512 for EGTEA Gaze+, and 256 for 50Salads. We set the number of action queries M to 8 for Breakfast and 16 for 50Salads, since 50Salads includes more actions than Breakfast in a video. For EpicKitchens and EGTEA Gaze+, we set M to 1, as the task is treated as a multi-label task defined in [39].

Training Details. To allow a fair comparison to other state-of-the-art long-term action anticipation methods, we use the pre-extracted I3D features [6] as input visual features F for all datasets, provided by [12] and [39]. We sample the I3D features with a stride of 3 for Breakfast and 50Salads and 1 for EpicKitchens and EGTEA Gaze+. In training, we set the observation rate $\alpha \in \{0.2, 0.3, 0.4, 0.5\}$ for Breakfast and 50Salads and additionally use $\{0.6, 0.7, 0.8\}$ for EpicKitchens and EGTEA. We use AdamW optimizer and train our model for 50, 30, 50, 100 epochs for Breakfast, 50Salads, EpicKitchens, and EGTEA Gaze+, respectively. We employ a cosine annealing warm-up scheduler with warm-up stages of 10 epochs. The total number of steps is set as $S = 1000$. To sample diffusion steps s , we adopt importance sampling [19, 41] in our experiments. For all trained models, we apply gradient clipping at 1. More training details can be found in the supplementary material.

5. Results

To compare DIFFANT with the state-of-the-art methods in Section 5.1 and conduct ablation studies in Section 5.2, we initialize the action embeddings \mathbf{z}_S as zeros for deterministic predictions. In Section 5.3 and Section 5.4, we sample \mathbf{z}_S from the standard Gaussian distribution, and present an in-depth analysis of our approach.

Method	EpicKitchens			EGTEA Gaze+		
	All	Freq	Rare	All	Freq	Rare
I3D [6]	32.7	53.3	23.0	72.1	79.3	53.3
ActionVLAD [17]	29.8	53.5	18.6	73.3	79.0	58.6
Timeception [25]	35.6	55.9	26.1	74.1	79.7	<u>59.7</u>
VideoGraph [26]	22.5	49.4	14.0	67.7	77.1	47.2
EGO-TOPO [39]	38.0	<u>56.9</u>	<u>29.2</u>	73.5	80.7	54.7
ANTICIPATR [40]	39.1	58.1	29.1	<u>76.8</u>	<u>83.3</u>	55.1
DIFFANT (Ours)	<u>38.7</u>	55.0	31.0	77.3	83.5	61.4

Table 2. Comparison with state-of-the-art methods on EpicKitchens [8] and EGTEA Gaze+ [34] in mAP. Bold and underlined numbers indicate the highest and second highest accuracy, respectively. DIFFANT achieves competitive results on EpicKitchens and sets a new state-of-the-art on EGTEA Gaze+.

5.1. Comparison with the State-of-the-art

In Table 1, we compare our methods with the state-of-the-art methods on Breakfast and 50Salads. While we list the methods that take the action labels predicted by an action segmentation method (e.g., [43]) at the top of the table for completeness, we mainly compare our method with methods that take visual features as input. Note that we do not list the results of ANTICIPATR [40] in Table 1, as they used a different evaluation protocol [61]. DIFFANT achieves the state-of-the-art performance in almost all experimental settings on Breakfast and 50Salads. Notably, our method outperforms other methods with a large margin for long time horizons for anticipation, e.g., for $\alpha = 0.3$ and $\beta = 0.5$, accuracy improvement on Breakfast and 50Salads is 18.9% (25.87 \rightarrow 30.77) and 24.1% (16.30 \rightarrow 20.23), respectively.

Next we evaluate our method on EpicKitchens and EGTEA Gaze+, as presented in Table 2. The results indicate that our approach achieves competitive results when compared to the state-of-the-art method [40] on EpicKitchens. While our method ranks second for the *All* set, it achieves the best results for the more challenging *Rare* set. Notably, our approach does not require separate training for the encoder and decoder, a step that is necessary in [40]. For EGTEA Gaze+, our method establishes a new benchmark, setting the state-of-the-art across all sets.

5.2. Ablation Study

Extensive ablation studies are performed to validate the design choices in our method. The experiments on Breakfast are conducted across all splits, and the results are averaged across these splits. Default settings are marked in gray.

Encoder Architecture. We first evaluate the impact of different encoder architectures in Table 3. In our default setting, we follow [5, 18] to use a transformer encoder with global attention across all past frames. As locality of fea-

Encoder	β ($\alpha = 0.3$)			
	0.1	0.2	0.3	0.5
DETR Encoder [5]	32.13	31.83	31.18	30.77
DETR Encoder [5] (local)	31.01	30.58	30.38	29.16
RoFormer [52]	22.26	21.94	22.42	21.20

Table 3. Impact of the encoder architecture on Breakfast in MoC.

\mathcal{L}_{seg}	$\mathcal{L}_{\text{smooth}}$	\mathcal{L}_{VLB}	Breakfast β ($\alpha = 0.3$)				EGTEA
			0.1	0.2	0.3	0.5	All
\times	\times	\checkmark	19.19	18.31	17.70	16.14	74.11
\checkmark	\times	\checkmark	31.66	31.44	32.17	30.91	77.07
\checkmark	\checkmark	\checkmark	32.13	31.83	31.18	30.77	77.33

Table 4. Impact of loss terms. The action segmentation loss \mathcal{L}_{seg} substantially improves the anticipation performance, indicating the importance of recognizing past frames in effectively anticipating future actions. The temporal smoothness loss $\mathcal{L}_{\text{smooth}}$ is optional and slightly improves the results, but not for all settings.

M	Breakfast β ($\alpha = 0.3$)				Steps	Breakfast β ($\alpha = 0.3$)			
	0.1	0.2	0.3	0.5		0.1	0.2	0.3	0.5
6	31.42	31.20	31.38	30.57	25	32.32	31.81	33.05	30.14
8	32.13	31.83	31.18	30.77	50	31.90	31.46	33.14	30.09
10	30.50	30.01	30.43	28.57	100	32.13	31.83	31.18	30.77
12	29.18	28.28	29.39	26.88	200	31.41	31.02	31.38	30.11

(a) # Action queries.

(b) # Inference steps.

Table 5. Ablation study on number of action queries M and inference steps on Breakfast. (a) Using $M = 8$ queries is sufficient. (b) 100 inference steps are sufficient.

tures has been found particularly beneficial for action segmentation [58], we similarly apply a hierarchical attention mask to the attention matrix, forcing the shallow layers to concentrate on neighboring frames, while deeper layers are allowed to attend to global frames. Details can be found in the supplementary material. Additionally, we employ the rotary position embedding technique [52], which allows to endow the transformer with relative positional embeddings without learnable parameters and has been widely adopted in the NLP community. The default setting consistently outperforms the other architectures across all prediction horizons (β), highlighting its efficacy in this context.

Loss Terms. In Table 4, we assess the impact of the action segmentation loss and the temporal smoothness loss on Breakfast and EGTEA Gaze+. As also observed in [18], the action segmentation loss substantially enhances the anticipation performance. This underscores the importance of recognizing past frames to effectively anticipate future actions. While introducing the temporal smoothness loss

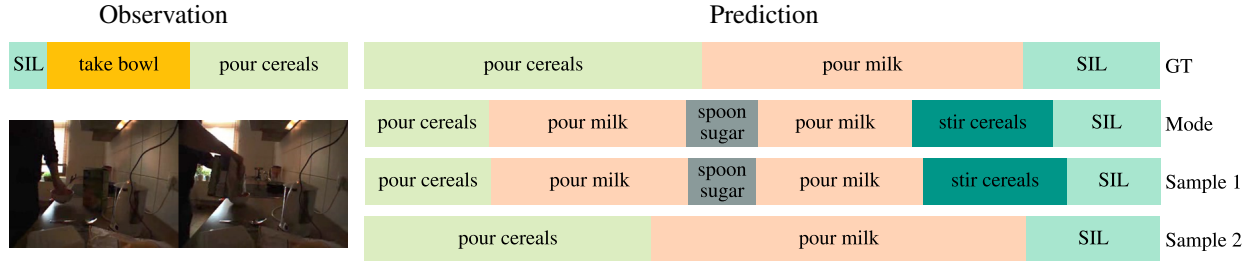


Figure 5. Qualitative results on Breakfast. Observations are shown on the left, while the ground-truth labels and predicted results, *i.e.*, one deterministic prediction and two randomly sampled results, are displayed on the right. We set α as 0.3 and predict all subsequent actions in this experiment. Action labels and durations are decoded as frame-wise action classes. More qualitative results are provided in the supplementary material.

Method	m	$\beta (\alpha = 0.2)$				$\beta (\alpha = 0.3)$			
		0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5
[1]	25	15.69	14.00	13.30	12.95	19.14	17.18	17.38	14.98
Ours	25	24.74	22.89	22.09	22.34	30.92	30.17	28.85	27.46
Ours	5	25.10	22.80	22.00	21.07	30.83	29.64	28.32	25.88
	10	24.94	23.34	22.74	21.78	30.86	29.31	28.66	26.66
	25	24.74	22.89	22.09	22.34	30.92	30.17	28.85	27.46
	50	24.81	23.24	22.26	22.67	30.92	30.37	28.80	27.22
	100	24.77	23.29	22.42	22.18	30.98	30.05	28.67	27.40

Table 6. Diverse anticipation performance on Breakfast averaged over m random samples for each observation. DIFFANT significantly exceeds the probabilistic approach [1].

into the framework improves the performance on EGTEA Gaze+, its effect on Breakfast is not clearly evident.

Number of Action Queries. In Table 5a, we analyze the impact of the number of action queries M . We vary M , starting from 6 and increasing in increments of 2 up to 12. For the Breakfast dataset, 6-8 queries are enough. Using more than 8 queries decreases the accuracy.

Diffusion Inference Steps. We adopt DDIM [47] to skip inference steps in our work. We vary the number of total inference steps and report the results in Table 5b. Using 100 iterations performs for most settings best, but even 25 iterations delivers reasonable results.

5.3. Uncertainty-Aware Anticipation

To evaluate the quality of our method for non-deterministic anticipation, we use the protocols that have been proposed in [1]. The protocols require the generation of m future predictions for the same video observation, calculating either the average accuracy of all generated predictions or the top-1 accuracy. We first report the results for the averaging protocol. For our method, we generate m future predictions by randomly sampling m future action embeddings \mathbf{z}_S from the Gaussian distribution $\mathcal{N}(\mathbf{0}, \mathbf{I})$. We then compute the average performance of these samples across all splits of

Method	m	$\beta (\alpha = 0.2)$				$\beta (\alpha = 0.3)$			
		0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5
[1]	25	28.89	28.43	27.61	28.04	32.38	31.60	32.83	30.79
Ours	25	31.27	29.87	29.40	30.11	37.38	37.01	36.29	34.80
Ours	–	25.33	24.59	24.39	22.74	32.13	31.83	31.18	30.77
	5	27.62	25.32	24.79	24.37	35.12	32.76	31.11	29.27
	10	28.77	27.49	27.36	27.16	35.63	34.16	32.96	31.59
	25	31.27	29.87	29.40	30.11	37.38	37.01	36.29	34.80
	50	31.08	30.54	30.34	31.52	37.35	38.13	37.51	36.22
	100	32.43	32.00	31.87	33.27	38.26	39.10	39.79	38.54

Table 7. Diverse anticipation performance on Breakfast using top-1 accuracy, *i.e.*, accuracy of best match to ground truth of m random samples for each observation. DIFFANT surpasses the probabilistic approach [1]. – denotes the result of our approach in the deterministic setting.

Breakfast. Since the previous probabilistic generative methods [38, 60] rely on ground truth action labels, we only compare our approach with the probabilistic approach [1]. The results for varying m are presented in Table 6, demonstrating that the averaged performance of DIFFANT maintains relatively stable with respect to changes in m . Moreover, our method significantly outperforms the approach [1] if we use the same number of samples $m = 25$.

It is important to recognize that the provided ground truth future in the dataset represents merely one among several feasible futures. This is exemplified in Fig. 5, where the left side depicts the observations, and the right side shows the ground truth labels alongside three predictions from DIFFANT. We label the prediction for zero \mathbf{z}_S as *Mode* and two randomly sampled \mathbf{z}_S as *Sample 1* and *Sample 2*. Notably, while Sample 2 matches with the provided ground truth, the other predictions also appear logical. Consequently, we adopt the top-1 protocol proposed in [1], which focuses on the performance of the prediction that best aligns with the ground truth, acknowledging that the ground truth represents just one of several possible futures. Specifically,

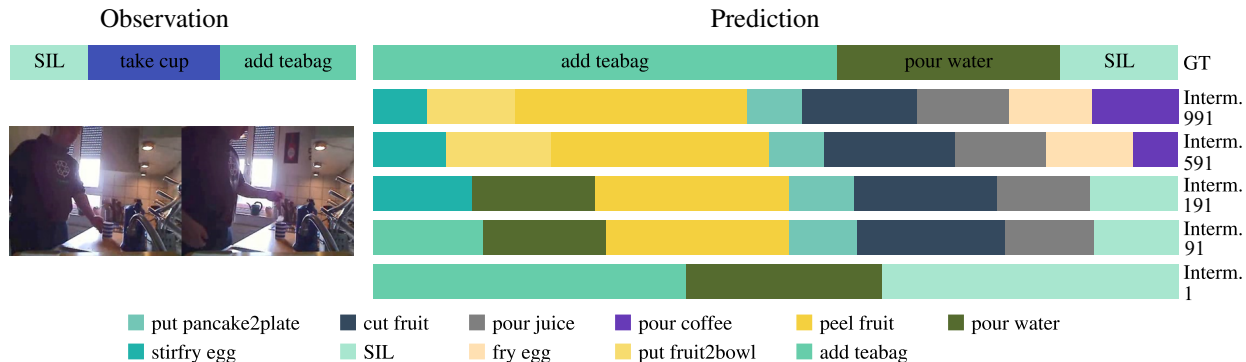


Figure 6. Anticipation results of the intermediate inference steps on Breakfast. Observations are shown on the left, while the ground-truth labels and predicted results are displayed on the right. α is set as 0.3 and all subsequent actions are predicted in this experiment. As the inference step approaches 1, the predictions are gradually refined, and the ground truth actions emerge.

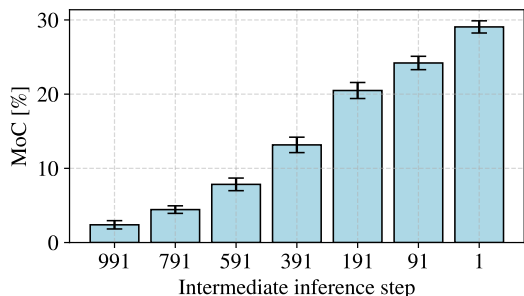


Figure 7. Intermediate anticipation performance on Breakfast. We run our model 25 times for each split and show the overall mean and standard deviation for the intermediate inference steps. α and β are set as 0.3 and 0.5, respectively.

for m generated predictions for the same observation, we select the prediction with the highest number of correctly predicted future frames. The performance for different m values is presented in Table 7. For comparison, we also include results where \mathbf{z}_S is a zero vector (denoted with $-$ for m). The performance of our non-deterministic anticipations consistently surpasses that of our deterministic one and the probabilistic approach [1] with few samples, underscoring the superior quality of DIFFANT. Furthermore, as expected, we observe a consistent increase in performance as the number of samples grows.

5.4. Intermediate Diffusion Anticipation

In this section, we delve into the denoising capability of the proposed method and assess the anticipation performance of the intermediate inference steps. As the total number of diffusion steps is set as $S = 1000$, with 100 steps used in inference, the inference trajectory becomes 991, 981, ..., 1. Fig. 7 provides a quantitative perspective across these varying intermediate inference steps. For this evaluation, we rerun our model 25 times and compute the performance of

each run for the randomly sampled \mathbf{z}_S . We then calculate the mean and standard deviation of these runs for each particular inference step. A discernible trend emerges from the figure: as the inference step approaches 1, there is a marked rise in the performance. This upward trajectory signifies the enhanced accuracy during the reverse diffusion process.

Fig. 6 presents the anticipation results derived at various intermediate inference steps. A noticeable progression is evident in the model’s predictions as the intermediate inference steps approach 1. The predicted actions become sharper and more aligned with the ground truth over time. Beginning from the 991-th step, the model discerns activities related to *egg cooking* and *juice preparation*. As the reverse diffusion progresses, the initially vague future sharpens. Consequently, the model, leveraging the visual observations, gains a deeper understanding and starts predicting actions pertinent to *tea-making* with increased accuracy. Additional qualitative results can be found in the supplementary material.

6. Conclusion

In this work, we have delved into the inherent challenges of predicting future actions, emphasizing the uncertainties that arise, especially when predictions extend far into the future. The proposed method, DIFFANT, leverages diffusion models to capture different possible future actions in the latent space, conditioned on the observed video. By introducing DIFFANT, we have showcased a novel approach to explicitly account for the non-deterministic nature of action anticipation. We have demonstrated the advantages of our method through extensive experiments on four benchmarks, achieving superior or comparable results to state-of-the-art methods. As the realm of action prediction continues to evolve, we believe that solutions like DIFFANT, which embrace a generative perspective, will play a pivotal role in navigating the complexities and uncertainties of future action anticipation.

References

- [1] Yazan Abu Farha and Juergen Gall. Uncertainty-Aware Anticipation of Activities. In *ICCV Workshop*, 2019. 1, 2, 5, 7, 8
- [2] Yazan Abu Farha, Qiuhong Ke, Bernt Schiele, and Juergen Gall. Long-term anticipation of activities with cycle consistency. In *GCPR*, pages 159–173. Springer, 2020. 2, 5
- [3] Tomer Amit, Eliya Nachmani, Tal Shaharbany, and Lior Wolf. SegDiff: Image segmentation with diffusion probabilistic models. *arXiv preprint arXiv:2112.00390*, 2021. 2
- [4] Dmitry Baranchuk, Ivan Rubachev, Andrey Voynov, Valentin Khruikov, and Artem Babenko. Label-efficient semantic segmentation with diffusion models. *ICLR*, 2021. 2
- [5] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *ECCV*, 2020. 1, 2, 3, 6
- [6] Joao Carreira and Andrew Zisserman. Quo Vadis, Action Recognition? A New Model and the Kinetics Dataset. In *CVPR*, 2017. 5, 6
- [7] Shoufa Chen, Peize Sun, Yibing Song, and Ping Luo. DiffusionDet: Diffusion model for object detection. *arXiv preprint arXiv:2211.09788*, 2022. 2
- [8] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. Scaling egocentric vision: The epic-kitchens dataset. In *ECCV*, pages 720–736, 2018. 2, 5, 6
- [9] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. The epic-kitchens dataset: Collection, challenges and baselines. *TPAMI*, 43(11):4125–4141, 2020. 2
- [10] Prafulla Dhariwal and Alexander Nichol. Diffusion models beat GANs on image synthesis. In *NeurIPS*, 2021. 2
- [11] Guodong Ding, Fadime Sener, and Angela Yao. Temporal action segmentation: An analysis of modern technique. *arXiv preprint arXiv:2210.10352*, 2022. 2
- [12] Yazan Abu Farha and Jurgen Gall. Ms-tcn: Multi-stage temporal convolutional network for action segmentation. In *CVPR*, pages 3575–3584, 2019. 4, 5, 2
- [13] Yazan Abu Farha, Alexander Richard, and Juergen Gall. When will you do what? - Anticipating Temporal Occurrences of Activities. In *CVPR*, 2018. 1, 2, 4, 5
- [14] Antonino Furnari and Giovanni Farinella. What Would You Expect? Anticipating Egocentric Actions With Rolling-Unrolling LSTMs and Modality Attention. In *ICCV*, 2019. 2
- [15] Jiyang Gao, Zhenheng Yang, and Ram Nevatia. RED: Reinforced Encoder-Decoder Networks for Action Anticipation. In *BMVC*, 2017. 2
- [16] Rohit Girdhar and Kristen Grauman. Anticipative Video Transformer. In *ICCV*, 2021. 2
- [17] Rohit Girdhar, Deva Ramanan, Abhinav Gupta, Josef Sivic, and Bryan Russell. Actionvlad: Learning spatio-temporal aggregation for action classification. In *CVPR*, pages 971–980, 2017. 6
- [18] Dayoung Gong, Joonseok Lee, Manjin Kim, Seong Jong Ha, and Minsu Cho. Future Transformer for Long-term Action Anticipation. In *CVPR*, 2022. 1, 2, 3, 4, 5, 6
- [19] Shansan Gong, Mukai Li, Jiangtao Feng, Zhiyong Wu, and LingPeng Kong. Diffuseq: Sequence to sequence text generation with diffusion models. *ICLR*, 2023. 2, 3, 4, 5, 1
- [20] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, et al. Ego4D: Around the World in 3,000 Hours of Egocentric Video. In *CVPR*, 2022. 2
- [21] Shuyang Gu, Dong Chen, Jianmin Bao, Fang Wen, Bo Zhang, Dongdong Chen, Lu Yuan, and Baining Guo. Vector quantized diffusion model for text-to-image synthesis. In *CVPR*, 2022. 2
- [22] Akash Gupta, Jingen Liu, Liefeng Bo, Amit K Roy-Chowdhury, and Tao Mei. A-act: Action anticipation through cycle transformations. *arXiv preprint arXiv:2204.00942*, 2022. 5
- [23] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models. In *NeurIPS*, pages 6840–6851, 2020. 1, 2, 3
- [24] Tobias Höppe, Arash Mehrjou, Stefan Bauer, Didrik Nielsen, and Andrea Dittadi. Diffusion models for video prediction and infilling. *arXiv preprint arXiv:2206.07696*, 2022. 2
- [25] Noureldien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Timeception for complex action recognition. In *CVPR*, pages 254–263, 2019. 6
- [26] Noureldien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Videograph: Recognizing minutes-long human activities in videos. In *ICCV Workshop*, 2019. 6
- [27] Qiuhong Ke, Mario Fritz, and Bernt Schiele. Time-Conditioned Action Anticipation in One Shot. In *CVPR*, 2019. 1, 2, 5
- [28] Gwanghyun Kim and Jong Chul Ye. DiffusionClip: Text-guided image manipulation using diffusion models. In *CVPR*, 2022. 2
- [29] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. In *ICLR*, 2014. 3
- [30] Hilde Kuehne, Ali Arslan, and Thomas Serre. The Language of Actions: Recovering the Syntax and Semantics of Goal-Directed Human Activities. In *CVPR*, 2014. 1, 5
- [31] Max WY Lam, Jun Wang, Dan Su, and Dong Yu. BDDM: Bilateral denoising diffusion models for fast and high-quality speech synthesis. *arXiv preprint arXiv:2203.13508*, 2022. 2
- [32] Yichong Leng, Zehua Chen, Junliang Guo, Haohe Liu, Jiawei Chen, Xu Tan, Danilo Mandic, Lei He, Xiang-Yang Li, Tao Qin, et al. Binauralgrad: A two-stage conditional diffusion probabilistic model for binaural audio synthesis. *arXiv preprint arXiv:2205.14807*, 2022. 2
- [33] Xiang Li, John Thickstun, Ishaan Gulrajani, Percy S Liang, and Tatsunori B Hashimoto. Diffusion-lm improves controllable text generation. *NeurIPS*, 35:4328–4343, 2022. 2, 3, 4, 1
- [34] Yin Li, Miao Liu, and James M Rehg. In the eye of beholder: Joint learning of gaze and actions in first person video. In *ECCV*, pages 619–635, 2018. 2, 5, 6

- [35] Daochang Liu, Qiyue Li, Anhdung Dinh, Tingting Jiang, Mubarak Shah, and Chang Xu. Diffusion action segmentation. *arXiv preprint arXiv:2303.17959*, 2023. 2, 3
- [36] Calvin Luo. Understanding diffusion models: A unified perspective. *arXiv preprint arXiv:2208.11970*, 2022. 3
- [37] Esteve Valls Mascaró, Hyemin Ahn, and Dongheui Lee. Intention-conditioned long-term human egocentric action anticipation. In *WACV*, pages 6048–6057, 2023. 2
- [38] Nazanin Mehrasa, Akash Abdu Jyothi, Thibaut Durand, Jiawei He, Leonid Sigal, and Greg Mori. A Variational Auto-Encoder Model for Stochastic Point Processes. In *CVPR*, 2019. 1, 2, 4, 7
- [39] Tushar Nagarajan, Yanghao Li, Christoph Feichtenhofer, and Kristen Grauman. Ego-topo: Environment affordances from egocentric video. In *CVPR*, pages 163–172, 2020. 5, 6
- [40] Megha Nawhal, Akash Abdu Jyothi, and Greg Mori. Rethinking learning approaches for long-term action anticipation. In *ECCV*, pages 558–576. Springer, 2022. 2, 3, 6
- [41] Alexander Quinn Nichol and Prafulla Dhariwal. Improved denoising diffusion probabilistic models. In *International Conference on Machine Learning, ICML*, 2021. 5, 1
- [42] AJ Piergiovanni, Anelia Angelova, Alexander Toshev, and Michael S Ryoo. Adversarial generative grammars for human activity prediction. In *ECCV*, pages 507–523. Springer, 2020. 2
- [43] Alexander Richard, Hilde Kuehne, and Juergen Gall. Weakly supervised action learning with rnn based fine-to-coarse modeling. In *CVPR*, pages 754–763, 2017. 6
- [44] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 2
- [45] Fadime Sener, Dipika Singhania, and Angela Yao. Temporal Aggregate Representations for Long-Range Video Understanding. In *ECCV*, 2020. 1, 2, 5
- [46] Jascha Sohl-Dickstein, Eric Weiss, Niru Maheswaranathan, and Surya Ganguli. Deep unsupervised learning using nonequilibrium thermodynamics. In *ICML*, 2015. 2
- [47] Jiaming Song, Chenlin Meng, and Stefano Ermon. Denoising diffusion implicit models. In *ICLR*, 2021. 1, 2, 3, 4, 7
- [48] Yang Song and Stefano Ermon. Generative modeling by estimating gradients of the data distribution. *NeurIPS*, 2019. 2
- [49] Yang Song and Stefano Ermon. Improved techniques for training score-based generative models. *NeurIPS*, 2020.
- [50] Yang Song, Jascha Sohl-Dickstein, Diederik P Kingma, Abhishek Kumar, Stefano Ermon, and Ben Poole. Score-based generative modeling through stochastic differential equations. *arXiv preprint arXiv:2011.13456*, 2020. 2
- [51] Sebastian Stein and Stephen J McKenna. Combining embedded accelerometers with computer vision for recognizing food preparation activities. In *UbiComp*, pages 729–738, 2013. 5
- [52] Jianlin Su, Yu Lu, Shengfeng Pan, Ahmed Muratdha, Bo Wen, and Yunfeng Liu. Roformer: Enhanced transformer with rotary position embedding. *arXiv preprint arXiv:2104.09864*, 2021. 6, 2, 3
- [53] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In *NeurIPS*, 2017. 1, 2
- [54] Vikram Voleti, Alexia Jolicoeur-Martineau, and Christopher Pal. Masked conditional video diffusion for prediction, generation, and interpolation. *arXiv preprint arXiv:2205.09853*, 2022. 2
- [55] Carl Vondrick, Hamed Pirsiavash, and Antonio Torralba. Anticipating visual representations from unlabeled video. In *CVPR*, pages 98–106, 2016. 2
- [56] Yunke Wang, Xiyu Wang, Anh-Dung Dinh, Bo Du, and Chang Xu. Learning to schedule in diffusion probabilistic models. In *KDD*, 2023. 2
- [57] Ruihan Yang, Prakhar Srivastava, and Stephan Mandt. Diffusion probabilistic modeling for video generation. *arXiv preprint arXiv:2203.09481*, 2022. 2
- [58] Fangqiu Yi, Hongyu Wen, and Tingting Jiang. Asformer: Transformer for action segmentation. In *BMVC*, 2021. 6, 1, 2
- [59] Peiyu Yu, Sirui Xie, Xiaojian Ma, Baoxiong Jia, Bo Pang, Ruigi Gao, Yixin Zhu, Song-Chun Zhu, and Ying Nian Wu. Latent diffusion energy-based model for interpretable text modeling. *arXiv preprint arXiv:2206.05895*, 2022. 2
- [60] He Zhao and Richard P Wildes. On diverse asynchronous activity anticipation. In *ECCV*, pages 781–799. Springer, 2020. 1, 2, 7
- [61] Zeyun Zhong, Manuel Martin, Michael Voit, Juergen Gall, and Jürgen Beyerer. A survey on deep learning techniques for action anticipation. *arXiv preprint arXiv:2309.17257*, 2023. 2, 6
- [62] Zeyun Zhong, David Schneider, Michael Voit, Rainer Stiefelwagen, and Jürgen Beyerer. Anticipative feature fusion transformer for multi-modal action anticipation. In *WACV*, 2023. 2
- [63] Xizhou Zhu, Weijie Su, Lewei Lu, Bin Li, Xiaogang Wang, and Jifeng Dai. Deformable detr: Deformable transformers for end-to-end object detection. *arXiv preprint arXiv:2010.04159*, 2020. 3

DIFFANT: Diffusion Models for Action Anticipation

Supplementary Material

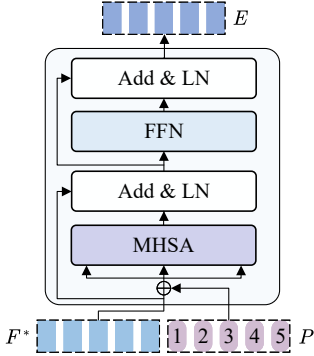


Figure 8. Details of an encoder layer.

7. Additional Methodology Details

Encoder Architecture. Our encoder follows the transformer encoder used in DETR [5] and FUTR [18], as depicted in Fig. 8. The original video features $F \in \mathbb{R}^{L \times K}$ are first passed through a linear layer, which adjusts them to a suitable hidden dimension D , yielding transformed features $F^* \in \mathbb{R}^{L \times D}$. Subsequently, these transformed features are combined with sinusoidal positional encodings $P \in \mathbb{R}^{L \times D}$, facilitating the encoder in producing refined representations $E \in \mathbb{R}^{L \times D}$ via self-attention [53].

Training Objective. Our training objective consists of a past observation recognition loss \mathcal{L}_{seg} , a temporal smoothness loss $\mathcal{L}_{\text{smooth}}$, and a variation lower bound loss \mathcal{L}_{VLB} :

$$\mathcal{L} = \mathcal{L}_{\text{seg}} + \mathcal{L}_{\text{smooth}} + \mathcal{L}_{\text{VLB}}.$$

\mathcal{L}_{seg} is a cross-entropy loss and $\mathcal{L}_{\text{smooth}}$ is computed as the mean squared error of the log-likelihoods between adjacent video frames to promote the local similarity along the temporal dimension. We compute the variational lower bound (\mathcal{L}_{VLB}) following the original diffusion process:

$$\mathcal{L}_{\text{VLB}} = \mathbb{E}_{q(\mathbf{z}_{1:S}|\mathbf{z}_0)} \left[\underbrace{\log \frac{q(\mathbf{z}_S|\mathbf{z}_0)}{p_\theta(\mathbf{z}_S)}}_{\mathcal{L}_S^{\text{emb}}} + \sum_{s=2}^S \underbrace{\log \frac{q(\mathbf{z}_{s-1}|\mathbf{z}_0, \mathbf{z}_s)}{p_\theta(\mathbf{z}_{s-1}|\mathbf{z}_s, E)}}_{\mathcal{L}_{s-1}^{\text{emb}}} \right. \\ \left. + \underbrace{\log \frac{q_\phi(\mathbf{z}_0|\mathbf{a})}{p_\theta(\mathbf{z}_0|\mathbf{z}_1, E)}}_{\mathcal{L}_0^{\text{emb}}} - \underbrace{\log p_\theta(\mathbf{a}|\mathbf{z}_0, E)}_{\mathcal{L}^{\text{pred}}} \right]. \quad (3)$$

Following [19, 33], we further simplify the training objective as follows:

$$\mathcal{L}_{\text{VLB}} = \sum_{s=1}^S \mathcal{L}_{s-1}^{\text{emb}} - \log p_\theta(\mathbf{a}|\mathbf{z}_0, E), \quad (4)$$

$$\text{where } \mathcal{L}_{s-1}^{\text{emb}} = \begin{cases} \|\mathbf{z}_0 - f_\theta(\mathbf{z}_s, s, E)\|^2 & \text{if } 2 \leq s \leq S \\ \|\text{EMB}(\mathbf{a}) - f_\theta(\mathbf{z}_1, 1, E)\|^2 & \text{if } s = 1, \end{cases}$$

and $-\log p_\theta(\mathbf{a}|\mathbf{z}_0, E) = -\log p_\theta(\mathbf{a}^c|\mathbf{z}_0, E) - \log p_\theta(\mathbf{a}^t|\mathbf{z}_0, E)$, under the assumption that action classes and durations are conditionally independent for mathematical convenience, as described in Section 3.3.

8. Additional Implementation Details

Generation of the Local Attention Mask. Drawing from the methodology outlined in [58], we conduct experiments where a hierarchical attention mask is applied to the attention matrix. This approach enables shallow layers to concentrate on neighboring frames, while deeper layers are allowed to attend to global frames. In each of the four encoder layers, we create an attention mask whose window size exponentially increases with the depth of the layer, with specific window sizes set at $\{9, 33, 129, 513\}$. This exponential growth is motivated by the intuition that lower layers capture finer, local details, while higher layers integrate more global, contextual information. During the self-attention computation, this mask is employed to assign minimal attention scores to positions outside the defined window, effectively minimizing their influence.

Training. For training, we tailor the batch size and learning rate to each dataset. Specifically, 50Salads utilizes a batch size of 8 and a learning rate of 1e-3. Breakfast is set with a batch size of 64 and a learning rate of 5e-4. Both EpicKitchens and EGTEA Gaze+ share a batch size of 32, but while EpicKitchens has a learning rate of 2.5e-4, EGTEA Gaze+ uses 5e-4. All experiments are done with a machine equipped with 4 RTX A6000 GPUs.

To sample diffusion step s during training, we employ importance sampling [41] in our experiments, following [19]. In contrast to uniform sampling, the importance-weighted sampling algorithm allocates more steps to diffusion steps with a larger loss value \mathcal{L}_s , and fewer to others:

$$\mathcal{L} = \mathbb{E}_{s \sim p_s} \left[\frac{\mathcal{L}_s}{p_s} \right], \quad p_s \propto \sqrt{\mathbb{E}[\mathcal{L}_s^2]}, \quad \sum_{s=1}^S p_s = 1. \quad (5)$$

9. Additional Analysis

Diffusion Step Sampling Method. In Table 8, we compare importance sampling with uniform sampling on the Breakfast dataset. The advantages of importance sampling over uniform sampling are not definitively clear. While it demonstrates superior performance in scenarios with

Samp.	β ($\alpha = 0.2$)				β ($\alpha = 0.3$)			
	0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5
Uni.	23.67	23.54	23.67	22.52	32.14	32.27	32.82	30.83
Imp.	25.33	24.59	24.39	22.74	32.13	31.83	31.18	30.77

Table 8. Comparison between uniform sampling and importance sampling on Breakfast using MoC.

Encoder	F1{10/25/50}	Edit	Acc.	MoC
DETR [5]	24.84/21.79/16.56	35.19	54.33	31.48
DETR [5] (local)	29.53/25.70/19.14	36.31	52.50	30.28
RoFormer [52]	7.77/5.57/3.26	14.26	37.88	21.96

Table 9. Observation segmentation performance on Breakfast with $\alpha = 0.3$. For direct comparison, we also include the anticipation performance (marked as gray), averaged over all prediction horizons (β).

shorter observation periods, its efficacy diminishes in cases involving longer observation durations.

Observation Segmentation Performance. Our model, enhanced with a classification head attached to the encoder, also exhibits the capability to segment past observations. In line with previous work in action segmentation [12, 35, 58], we report frame-wise accuracy (Acc), edit score (Edit), and F1 scores at overlap thresholds 10%, 25%, 50% (F1@10, 25, 50) in Table 9. Accuracy assesses the results at the frame level, while the edit score and F1 scores measure the performance at the segment level. We note that comparing the segmentation performance of our model, optimized for action anticipation, to state-of-the-art action segmentation methods is not straightforward due to our model’s focus on partial video inputs. For direct comparison, we also include the anticipation performance averaged across all prediction horizons. As indicated in Table 9, variants of DETR [5] significantly surpass RoFormer [52] across all metrics, highlighting the significance of precise past observation recognition in action anticipation. Although local attention contributes positively to observation recognition, its benefit for action anticipation is not conclusively established.

Computational Cost. Table 10 presents a comprehensive computation analysis of our model on Breakfast, capturing the mean over classes (MoC) across various prediction horizons. The model comprises 80 million parameters and utilizes 28.1 GB of GPU memory during training for a batch size of 64. We compute the floating-point operations per second (FLOPs) of DIFFANT using THOP¹. The FLOPs and inference time are calculated when inferring a video sequence of 1000 frames on a single RTX A6000 GPU, both

¹<https://github.com/Lyken17/pytorch-OpCounter>

# Steps	MoC	# params	GPU Mem.	FLOPs	Inference time
10	27.41	80.0M	28.1G	5.94G	39.04ms
25	27.32	80.0M	28.1G	10.09G	94.31ms
50	27.83	80.0M	28.1G	17.02G	185.41ms
100	27.87	80.0M	28.1G	30.87G	370.78ms

Table 10. The mean over classes (MoC) averaged across all observation ratios ($\alpha = 0.2$ and 0.3) and all prediction horizons on Breakfast, the number of parameters, the FLOPs at inference for a video of 1000 frames, the GPU memory cost during training for a batch size of 64, and the inference time on an A6000 GPU of DIFFANT.

of which scale with the number of diffusion inference steps. While our model achieves optimal performance with 100 steps, maintaining an acceptable inference time for long-term action anticipation, it also produces satisfactory outcomes using just 10 steps, with a notably quick inference time of 39.04 milliseconds.

10. Additional Ablation Results

In Table 11, we provide the comprehensive results for the ablation study (Section 5.2) with two observation ratios $\alpha \in \{0.2, 0.3\}$ for Breakfast and all sets for EGTEA Gaze+. The experimental trends generally show consistency across both observation ratios when examining variations in the encoder architecture and inference steps. However, there are minor differences in the effects of varying the loss terms and the number of action queries. Specifically, the temporal smoothness loss positively impacts performance across all sets on EGTEA Gaze+, but its advantage for the Breakfast dataset is less pronounced. In fact, this loss term slightly deteriorates performance in scenarios with shorter observations. Regarding the number of action queries, using six queries performs slightly better for shorter observations, whereas eight queries performs slightly better for longer observations.

11. Additional Qualitative Results

Additional qualitative results are showcased in Fig. 9, where the observations are displayed on the left, and the right side features the ground truth labels alongside three predictions from our model. The prediction corresponding to future action embeddings \mathbf{z}_S initialized with zero vectors is denoted as *Mode*. In contrast, predictions derived from randomly sampled \mathbf{z}_S are identified as *Sample 1* and *Sample 2*. Moreover, Fig. 10 presents a collection of 20 random sample results. Given an observed video with the actions *taking knife* and *cutting bun*, our model successfully predicts actions associated with *preparing juice* and *making sandwich*, demonstrating its proficiency in generating varied yet pertinent future actions. In instances where the observation

Encoder	Breakfast β ($\alpha = 0.2$)				Breakfast β ($\alpha = 0.3$)			
	0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5
DETR Encoder [5]	25.33	24.59	24.39	22.74	32.13	31.83	31.18	30.77
DETR Encoder [5] (local)	23.39	24.59	24.28	23.48	31.01	30.58	30.38	29.16
RoFormer [52]	18.38	18.86	19.35	19.43	22.26	21.94	22.42	21.20

(a) Encoder architecture.

Loss			Breakfast β ($\alpha = 0.2$)				Breakfast β ($\alpha = 0.3$)				EGTEA		
\mathcal{L}_{seg}	$\mathcal{L}_{\text{smooth}}$	\mathcal{L}_{VLB}	0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5	All	Freq	Rare
\times	\times	\checkmark	13.31	12.71	12.95	12.71	19.19	18.31	17.70	16.14	74.11	80.59	57.26
\checkmark	\times	\checkmark	25.35	25.81	24.83	23.64	31.66	31.44	32.17	30.91	77.07	83.27	60.95
\checkmark	\checkmark	\checkmark	25.33	24.59	24.39	22.74	32.13	31.83	31.18	30.77	77.33	83.47	61.35

(b) Loss terms.

# Queries	Breakfast β ($\alpha = 0.2$)				Breakfast β ($\alpha = 0.3$)			
	0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5
4	23.78	23.88	24.02	23.43	31.55	31.24	30.57	32.78
6	25.35	25.05	25.44	25.33	31.42	31.20	31.38	30.57
8	25.33	24.59	24.39	22.74	32.13	31.83	31.18	30.77
10	24.53	24.85	24.32	22.22	30.50	30.01	30.43	28.57
12	26.27	24.83	24.09	23.49	29.18	28.28	29.39	26.88

(c) Number of action queries.

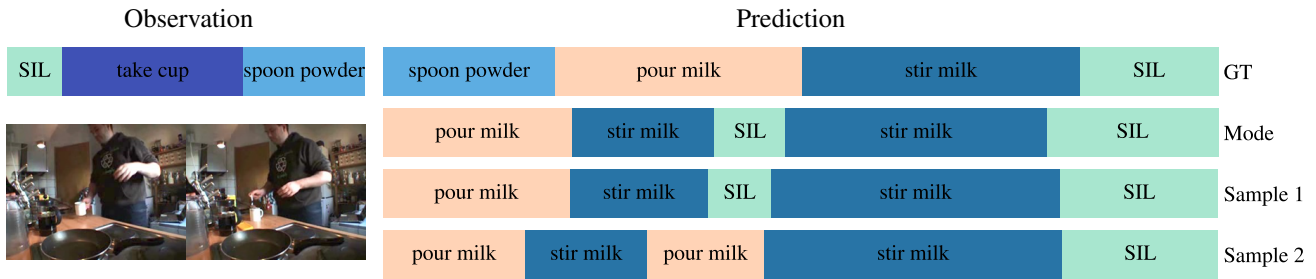
# Steps	Breakfast β ($\alpha = 0.2$)				Breakfast β ($\alpha = 0.3$)				EGTEA		
	0.1	0.2	0.3	0.5	0.1	0.2	0.3	0.5	All	Freq	Rare
5	25.29	24.15	24.09	22.74	32.05	31.64	31.35	29.58	75.18	80.95	60.15
10	24.72	23.52	23.33	22.50	32.41	31.45	31.52	29.85	76.75	83.56	59.06
25	23.76	23.23	22.71	21.55	32.32	31.81	33.05	30.14	76.69	82.58	60.78
50	24.29	23.97	23.94	23.88	31.90	31.46	33.14	30.09	76.33	82.69	59.79
100	25.33	24.59	24.39	22.74	32.13	31.83	31.18	30.77	77.33	83.47	61.35
200	23.95	23.80	23.49	22.28	31.41	31.02	31.38	30.11	76.47	82.78	60.04

(d) Number of inference steps.

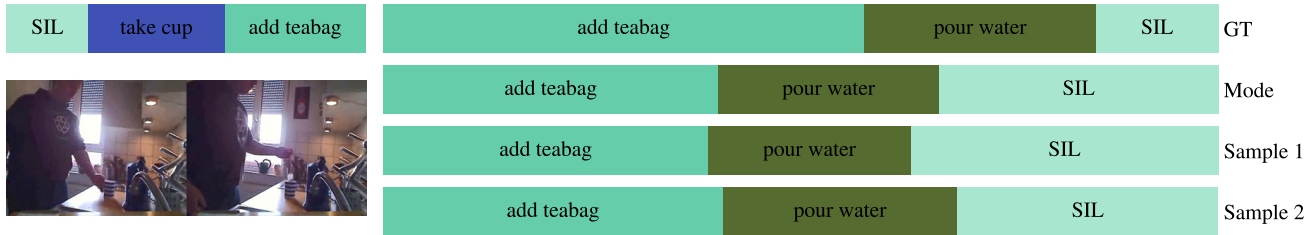
Table 11. Comprehensive results of the ablation study (Section 5.2).

offers only limited insight into the current activity, as illustrated in Fig. 11, our model exhibits a broader range of potential futures. These include actions such as *preparing cereals*, *cooking egg*, *preparing fruit*, and *making sandwich*, alongside the actual ground truth activity of *preparing milk*.

Additional intermediate diffusion results are presented in Fig. 12. Consistent with the observations detailed in Section 5.4, the clarity and alignment of the predicted actions with the ground truth progressively improve over time. This improvement is evident as the reverse diffusion process advances, transforming initially indistinct future predictions into more defined and accurate outcomes.



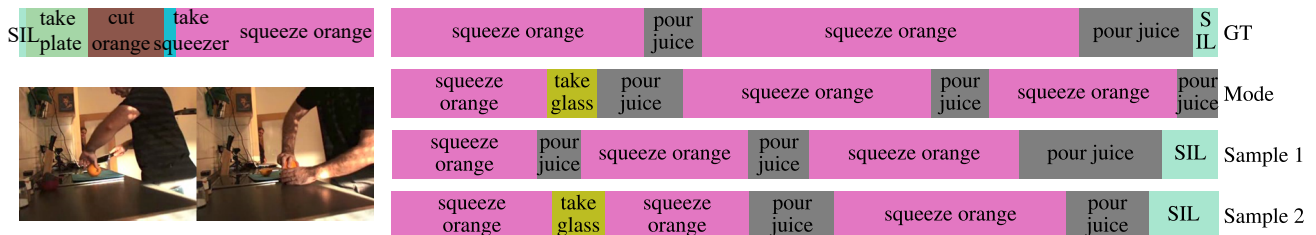
(a) Activity: Milk



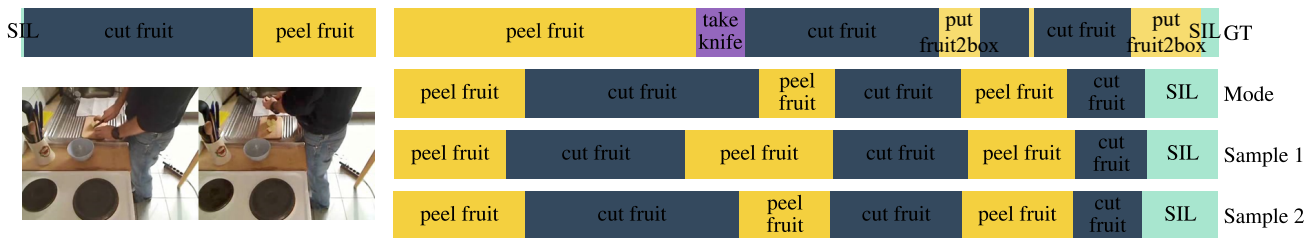
(b) Activity: Tea



(c) Activity: Cereals



(d) Activity: Juice



(e) Activity: Salat

Figure 9. Additional qualitative results on Breakfast. Observations are shown on the left, while the ground-truth labels and predicted results, *i.e.*, one deterministic prediction and two randomly sampled results, are displayed on the right. We set α as 0.3 and predict all subsequent actions in this experiment. Action labels and durations are decoded as frame-wise action classes.



Figure 10. Uncertainty-aware anticipations on Breakfast of the activity *make sandwich*. Observation rate α is set as 0.3.

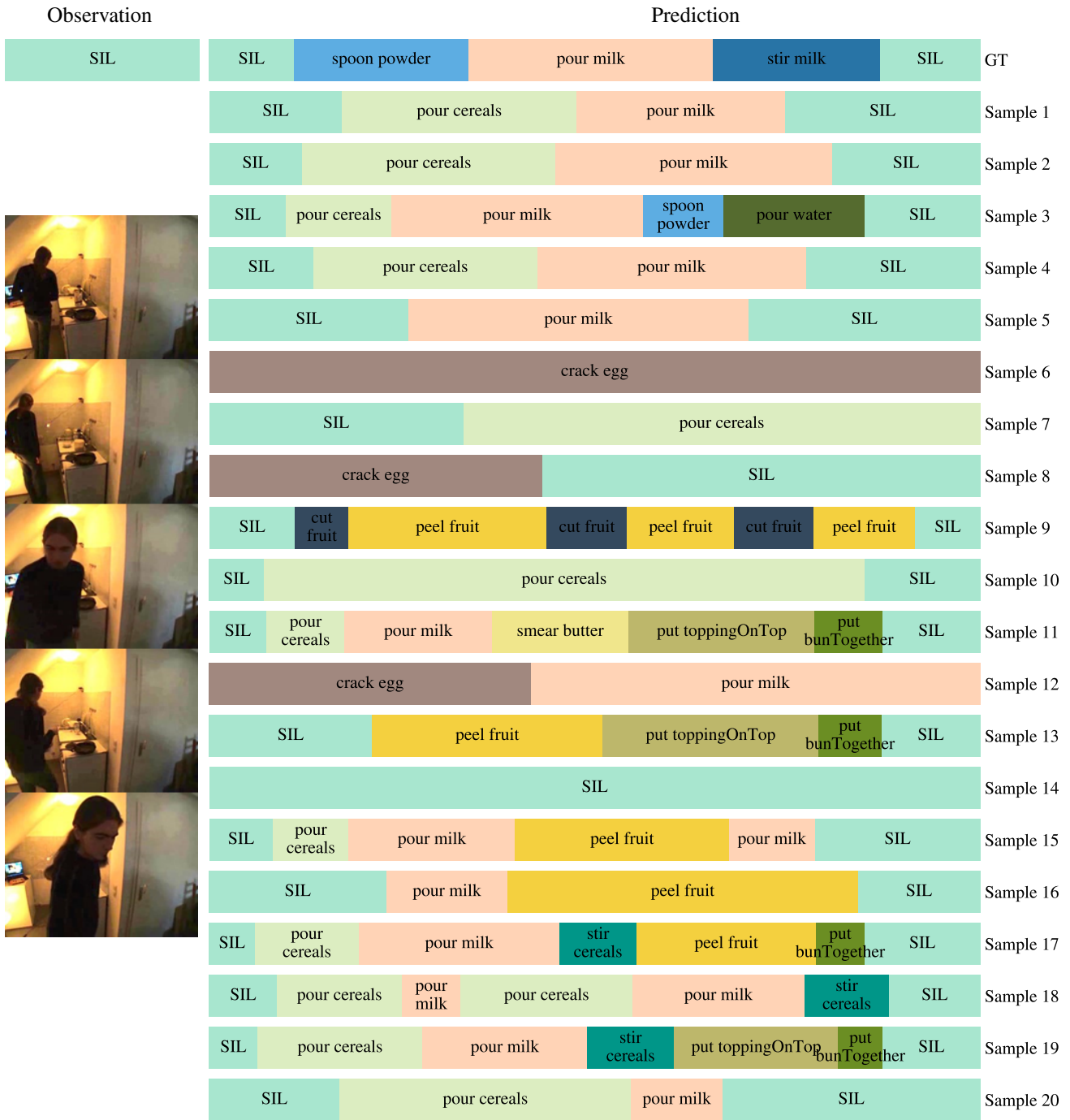


Figure 11. Uncertainty-aware anticipations on Breakfast of the activity *prepare milk*. Observation rate α is set as 0.2.

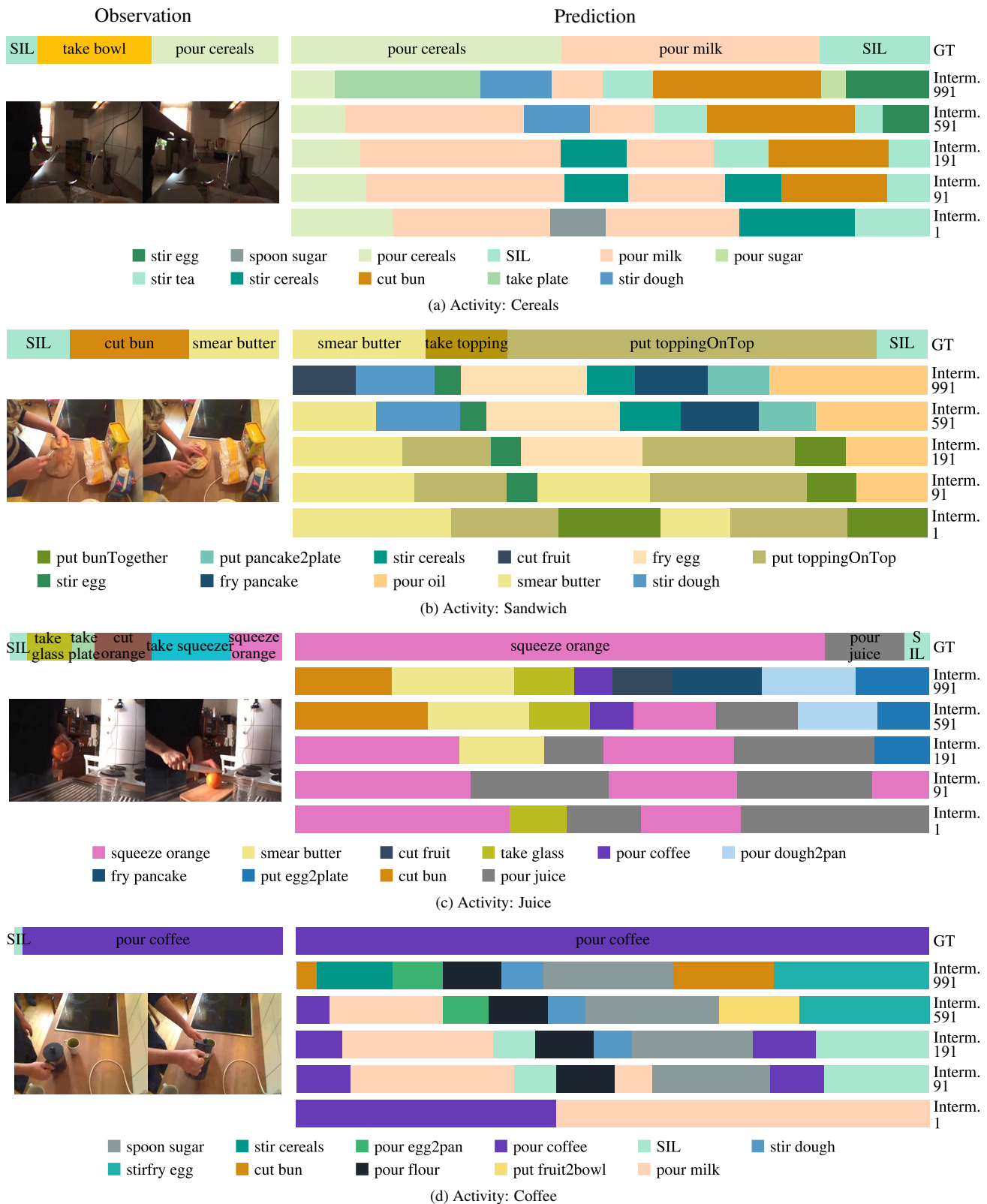


Figure 12. Anticipation results of the intermediate inference steps on Breakfast. Observations are shown on the left, while the ground-truth labels and predicted results are displayed on the right. α is set as 0.3 and all subsequent actions are predicted in this experiment. As the inference step approaches 1, the predictions are gradually refined, and the ground truth actions emerge.