

# End-to-End Spatio-Temporal Action Localisation with Video Transformers

Alexey A. Gritsenko Xuehan Xiong Josip Djolonga Mostafa Dehghani  
 Chen Sun Mario Lučić Cordelia Schmid Anurag Arnab  
 Google

{agritsenko, aarnab}@google.com

## Abstract

The most performant spatio-temporal action localisation models use external person proposals and complex external memory banks. We propose a fully end-to-end, purely-transformer based model that directly ingests an input video, and outputs tubelets – a sequence of bounding boxes and the action classes at each frame. Our flexible model can be trained with either sparse bounding-box supervision on individual frames, or full tubelet annotations. And in both cases, it predicts coherent tubelets as the output. Moreover, our end-to-end model requires no additional pre-processing in the form of proposals, or post-processing in terms of non-maximal suppression. We perform extensive ablation experiments, and significantly advance the state-of-the-art on five different spatio-temporal action localisation benchmarks with both sparse keyframes and full tubelet annotations.

## 1. Introduction

Spatio-temporal action localisation is an important problem with applications in advanced video search engines, robotics and security among others. It is typically formulated in one of two ways: Firstly, predicting the bounding boxes and actions performed by an actor at a single keyframe given neighbouring frames as spatio-temporal context [18, 28]. Or alternatively, predicting a sequence of bounding boxes and actions (*i.e.* “tubes”), for each actor at each frame in the video [21, 49].

The most performant models [3, 15, 40, 62], particularly for the first, keyframe-based formulation of the problem, employ a two-stage pipeline inspired by the Fast-RCNN object detector [17]: They first run a separate person detector to obtain proposals. Features from these proposals are then aggregated and classified according to the actions of interest. These models have also been supplemented with memory banks containing long-term contextual information from other frames [40, 53, 61, 62], and/or detections of other potentially relevant objects [2, 53] to capture additional scene context, achieving state-of-the-art results.

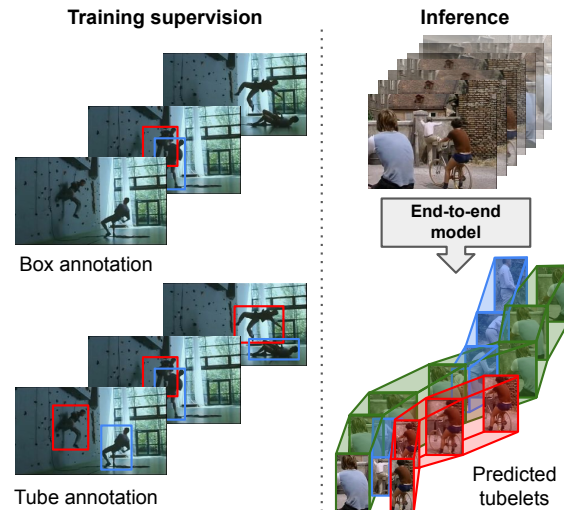


Figure 1. We propose an end-to-end Spatio-Temporal Action Recognition model named STAR. Our model is end-to-end in that it does not require any external region proposals to predict tubelets – sequences of bounding boxes associated with a given person in every frame and their corresponding action classes. Our model can be trained with either sparse box annotations on selected keyframes, or full tubelet supervision.

And whilst proposal-free algorithms, which do not require external person detectors, have been developed for detecting both at the keyframe-level [9, 25, 51] and tubelet-level [22, 66], their performance has typically lagged behind their proposal-based counterparts. Here, we show for the first time that an end-to-end trainable spatio-temporal model outperforms a two-stage approach.

As shown in Fig. 1, we propose our Spatio-Temporal Action Transformer (STAR) that consists of a transformer architecture, and is based on the DETR [6] detection model. Our model is “end-to-end” in that it does not require pre-processing in the form of proposals, nor post-processing in the form of non-maximal suppression (NMS) in contrast to the majority of prior work. The initial stage of the model is a vision encoder. This is followed by a decoder that processes learned latent queries, which represent each actor in the video, into output tubelets – a sequence of bounding boxes and action classes at each time step of

the input video clip. Our model is versatile in that we can train it with either fully-labeled tube annotations, or with sparse keyframe annotations (when only a limited number of keyframes are labelled). In the latter case, our network still predicts tubelets, and learns to associate detections of an actor, from one frame to the next, without explicit supervision. This behaviour is facilitated by our formulation of factorised queries, decoder architecture and tubelet matching in the loss which all contain temporal inductive biases.

We conduct thorough ablation studies of these modelling choices, confirming the benefit of temporal inductive biases in our model design. Informed by these experiments, we achieve state-of-the-art on both keyframe-based action localisation datasets like AVA [18] and AVA-Kinetics [28], and also tubelet-based datasets like UCF101-24 [49], JHMDB [21] and MultiSports [32]. In particular, we achieve a Frame mAP of 45.1 on AVA-Kinetics, outperforming the best previous results achieved by a massive video foundation model [57]. Moreover, our state-of-the-art results are achieved with a single forward-pass through the model, using only a video clip as input, and without any separate external person detectors providing proposals [57, 60, 62], complex memory banks [40, 62, 66], or additional object detectors [2, 53], as used by the prior works. Furthermore, we outperform these complex, prior, state-of-the-art two-stage models whilst also having additional functionality in that our model predicts tubelets, that is, temporally consistent bounding boxes at each frame of the input video clip. This capability is demonstrated by our results on tube-based datasets like UCF101-24 where we surpass the prior work [66] by 13.2 points on Video AP50.

## 2. Related Work

Models for spatio-temporal action localisation have typically built upon advances in object detectors for images. The most performant methods for action localisation [3, 15, 40, 53, 62] are based on “two-stage” detectors like FastRCNN [17]. These models use external, pre-computed person detections, and use them to ROI-pool features which are then classified into action classes. Although these models are cumbersome in that they require an additional model and backbone to first detect people, and therefore additional detection training data as well, they are currently the leading approaches on datasets such as AVA [18]. Such models using external proposals are also particularly suited to datasets such as AVA [18] as each person is exhaustively labelled as performing an action, and therefore there are fewer false-positives from using action-agnostic person detections compared to datasets such as UCF101 [49].

The performance of these two-stage models has further been improved by incorporating more contextual information using feature banks extracted from additional frames in the video [40, 53, 61, 62] or by utilising detections of ad-

ditional objects in the scene [2, 5, 58, 64]. Both of these cases entail significant extra computation and complexity to train additional auxiliary models, and to precompute features from them that are then used during training and inference of the localisation model.

Our proposed method, in contrast, is end-to-end in that it directly produces detections without any additional inputs besides a video clip. Moreover, it outperforms these prior works without resorting to external proposals or memory banks, showing that a transformer backbone is sufficient to capture long-range dependencies in the input video. In addition, unlike previous two-stage methods, our method directly predicts tubelets: a sequence of bounding boxes and actions for each frame of the input video, and can do so even when we do not have full tubelet annotations available.

A number of proposal-free action localisation models have also been developed [9, 16, 22, 25, 51, 66]. These methods are based upon alternative object detection architectures such as SSD [35], CentreNet [67], YOLO [43], DETR [6] and Sparse-RCNN [52]. However, in contrast to our approach, they have been outperformed by their proposal-based counterparts. Moreover, some of these methods [16, 25, 51] also consist of separate network backbones for learning video feature representations and proposals for a keyframe, and are thus effectively two networks trained jointly, and cannot predict tubelets either.

Among prior works that do not use external proposals, and also directly predict tubelets [22, 29, 31, 47, 48], our work is most similar to TubeR [66] given that our model is also based on DETR. Our model, however, is designed with additional temporal inductive biases which improves accuracy (without using external memory banks precomputed offline like [66]). And moreover, unlike TubeR, we also demonstrate how our model can predict tubelets (*i.e.* predictions at every frame of the input video), even when we only have sparse keyframe supervision (*i.e.* ground truth annotation for a limited number of frames) available.

Finally, we note that DETR has also been extended as a proposal-free method to addressing other localisation tasks in video such as instance segmentation [59], temporal localisation [36, 39, 63] and moment retrieval [27].

## 3. Spatio-Temporal Action Transformer

Our proposed model ingests a sequence of video frames, and directly predicts tubelets (a sequence of bounding boxes and action labels). No external person detections [40, 54, 57], or memory banks [62, 66], are needed.

As summarised in Fig. 2, our model consists of a vision encoder (Sec. 3.1), followed by a decoder which processes learned query tokens into output tubelets (Sec. 3.2). We incorporate temporal inductive biases into our decoder to improve accuracy and tubelet prediction with weaker supervision. Our model is inspired by the DETR architecture [6]

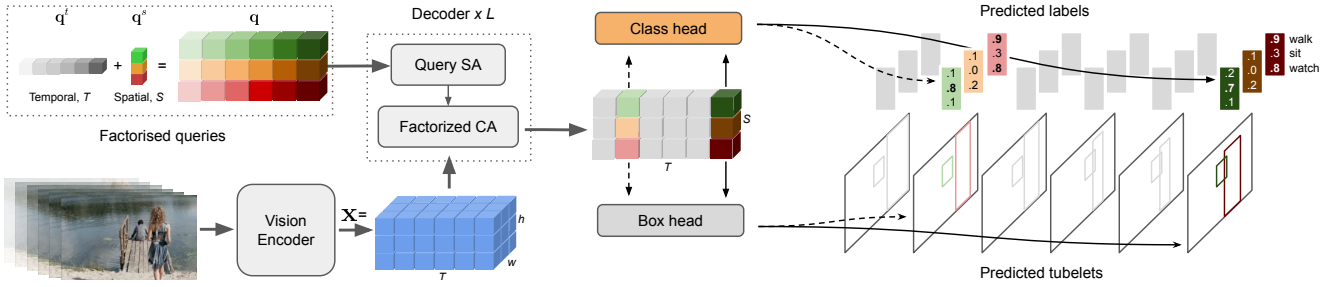


Figure 2. Our model processes a fixed-length video clip, and for each frame, outputs tubelets (*i.e.* linked bounding boxes with associated action class probabilities). It consists of vision encoder which outputs a video representation,  $\mathbf{x} \in \mathbb{R}^{T \times h \times w \times d}$ . The video representation, along with learned queries,  $\mathbf{q}$  (which are factorised into spatial  $\mathbf{q}^s$  and temporal components  $\mathbf{q}^t$ ) are decoded into tubelets by a decoder of  $L$  layers followed by shallow box and class prediction heads.

for object detection in images, and is also trained with a set-based loss and Hungarian matching. We detail our loss, and how we can train with either sparse keyframe supervision or full tubelet supervision, in Sec. 3.3.

### 3.1. Vision Encoder

The vision backbone processes an input video,  $\mathbf{X} \in \mathbb{R}^{T \times H \times W \times 3}$  to produce a feature representation of the input video  $\mathbf{x} \in \mathbb{R}^{t \times h \times w \times d}$ . Here,  $T$ ,  $H$  and  $W$  are the original temporal-, height- and width-dimensions of the input video respectively, whilst  $t$ ,  $h$  and  $w$  are the spatio-temporal dimensions of their feature representation, and  $d$  its latent dimension. When using a transformer backbone, these spatio-temporal dimensions depend on the patch size when tokenising the input, and when using a convolutional backbone, they depend on the overall stride. To retain spatio-temporal information, we remove the spatial- and temporal-aggregation steps at the end of the original backbone. And if the temporal patch size (or stride) is larger than 1, we bilinearly upsample the final feature map along the temporal axis to maintain the original temporal resolution.

### 3.2. Tubelet Decoder

Our decoder processes the visual features,  $\mathbf{x} \in \mathbb{R}^{T \times h \times w \times c}$ , along with learned queries,  $\mathbf{q} \in \mathbb{R}^{T \times S \times d}$ , to output tubelets,  $\mathbf{y} = (\mathbf{b}, \mathbf{a})$  which are a sequence of bounding boxes,  $\mathbf{b} \in \mathbb{R}^{T \times S \times 4}$  and corresponding actions,  $\mathbf{a} \in \mathbb{R}^{T \times S \times C}$ . Here,  $S$  denotes the maximum number of bounding boxes per frame (padded with “background” as necessary) and  $C$  denotes the number of output classes.

The idea of decoding learned queries into output detections using the transformer decoder architecture of Vaswani *et al.* [56] was used in DETR [6]. In summary, the decoder of [6, 56] consists of  $L$  layers, each performing a series of self-attention operations on the queries, and cross-attention between the queries and encoder outputs.

We modify the queries, self-attention and cross-attention operations for our spatio-temporal localisation scenario, as shown in Fig. 2 and 3 to include additional temporal inductive biases, and to improve accuracy as detailed below.

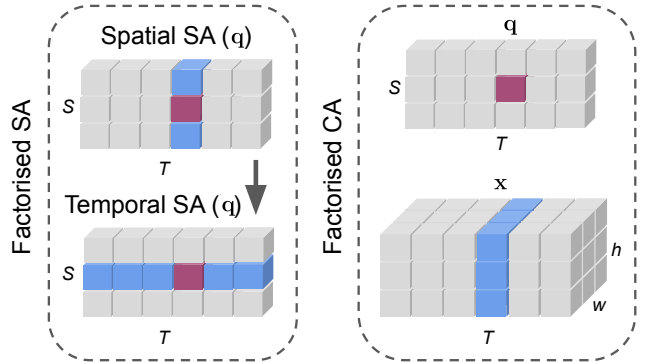


Figure 3. Our decoder layer consists of factorised self-attention (SA) (left) and cross-attention (CA) (right) operations designed to provide a spatio-temporal inductive bias and reduce computation. Both operations restrict attention to the same spatial and temporal slices as the query token, as illustrated by the receptive field (blue) for a given query token (magenta). Factorised SA consists of two operations, whilst in Factorised CA, there is one operation.

**Queries** Queries,  $\mathbf{q}$ , in DETR, are decoded using the encoded visual features,  $\mathbf{x}$ , into bounding box predictions, and are analogous to the “anchors” used in other detection architectures such as Faster-RCNN [44].

The most straightforward way to define queries is to randomly initialise  $\mathbf{q} \in \mathbb{R}^{T \times S \times d}$ , where there are  $S$  bounding boxes at each of the  $T$  input frames in the video clip.

However, we find it is more effective to factorise the queries into separate learned spatial,  $\mathbf{q}^s \in \mathbb{R}^{S \times d}$ , and temporal,  $\mathbf{q}^t \in \mathbb{R}^{T \times d}$  parameters. To obtain the final tubelet queries, we simply repeat the spatial queries across all frames, and add them to their corresponding temporal embedding at each location, as shown in Fig. 2. More concretely  $\mathbf{q}_{ij} = \mathbf{q}_i^t + \mathbf{q}_j^s$  where  $i$  and  $j$  denote the temporal and spatial indices respectively.

The factorised query representation means that the same spatial embedding is used across all frames. Intuitively, this encourages the  $i^{th}$  spatial query embedding,  $\mathbf{q}_i^s$ , to bind to the same location across different frames of the video, and since objects typically have small displacements from frame to frame, may help to associate bounding boxes within a tubelet together. We verify this intuition empirically in the experimental section.

**Decoder layer** The decoder layer in the original transformer [56] consists of self-attention on the queries,  $\mathbf{q}$ , followed by cross-attention between the queries and the outputs of the encoder,  $\mathbf{x}$ , and then a multilayer perceptron (MLP) layer [19, 56]:

$$\mathbf{u}^\ell = \text{MHSA}(\mathbf{q}^\ell) + \mathbf{q}^\ell, \quad (1)$$

$$\mathbf{v}^\ell = \text{CA}(\mathbf{u}^\ell, \mathbf{x}) + \mathbf{u}^\ell, \quad (2)$$

$$\mathbf{z}^\ell = \text{MLP}(\mathbf{v}^\ell) + \mathbf{v}^\ell, \quad (3)$$

where  $\mathbf{z}^\ell$  is the output of the  $\ell^{\text{th}}$  decoder layer,  $\mathbf{u}$  and  $\mathbf{v}$  are intermediate variables, MHSA denotes multi-headed self-attention and CA cross-attention. Note that the inputs to the MLP, self- and cross-attention operations are layer-normalised [4], which we omit here for clarity.

In our model, we factorise the self- and cross-attention layers across space and time respectively as shown in Fig. 3, to introduce a temporal locality inductive bias, and also to increase model efficiency. Concretely, when applying MHSA, we first compute the queries, keys and values, over which we attend twice: first independently at each time step with each frame, and then, independently along the time axis at each spatial location. Similarly, we modify the cross-attention operation so that only tubelet queries and backbone features from the same time index attend to each other.

**Localisation and classification heads** We obtain the final predictions of the network,  $\mathbf{y} = (\mathbf{b}, \mathbf{a})$ , by applying a small feed-forward network to the outputs to the decoder,  $\mathbf{z}$ , following DETR [6]. The sequence of bounding boxes,  $\mathbf{b}$ , is obtained with a 3-layer MLP, and is parameterised by the box center, width and height for each frame in the tubelet. A single-layer linear projection is used to obtain class logits,  $\mathbf{a}$ . As we predict a fixed number of  $S$  bounding boxes per frame, and  $S$  is more than the maximum number of ground truth instances in the frame, we also include an additional class label,  $\emptyset$ , which represents the “background” class which tubelets with no action class can be assigned to.

### 3.3. Training objective

Our model predicts bounding boxes and action classes at each frame of the input video. Many datasets, however, such as AVA [18], are only sparsely annotated at selected keyframes of the video. In order to leverage the available annotations, we compute our training loss, Eq. 4, only at the annotated frames of the video, after having matched the predictions to the ground truth:

$$\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} \mathcal{L}_{\text{frame}}(\mathbf{y}^t, \hat{\mathbf{y}}^t), \quad (4)$$

where  $\mathcal{T}$  is the set of labelled frames;  $\mathbf{y}$  and  $\hat{\mathbf{y}}$  denote the ground truth and predicted tubelets after matching. Following DETR [6], our training loss at each frame,  $\mathcal{L}_{\text{frame}}$ , is a

sum of an  $L_1$  regression loss on bounding boxes, the generalised IoU loss [45] on bounding boxes, and a cross-entropy loss on action labels:

$$\begin{aligned} \mathcal{L}_{\text{frame}}(\mathbf{b}^t, \hat{\mathbf{b}}^t, \mathbf{a}^t, \hat{\mathbf{a}}^t) = & \sum_i \mathcal{L}_{\text{box}}(\mathbf{b}_i^t, \hat{\mathbf{b}}_i^t) + \mathcal{L}_{\text{iou}}(\mathbf{b}_i^t, \hat{\mathbf{b}}_i^t) \\ & + \mathcal{L}_{\text{class}}(\mathbf{a}_i^t, \hat{\mathbf{a}}_i^t). \end{aligned} \quad (5)$$

**Matching** Set-based detection models such as DETR can make predictions in any order, which is why the predictions need to be matched to the ground truth before computing the training loss.

The first form of matching that we consider is to independently perform bipartite matching at each frame to align the model’s predictions to the ground truth (or the  $\emptyset$  background class) before computing the loss. In this case, we use the Hungarian algorithm [26] to obtain  $T$  permutations of  $S$  elements,  $\hat{\pi}^t \in \Pi^t$ , at each frame, where the permutation at the  $t^{\text{th}}$  frame minimises the per-frame loss,

$$\hat{\pi}^t = \arg \min_{\pi \in \Pi^t} \mathcal{L}_{\text{frame}}(\mathbf{y}^t, \hat{\mathbf{y}}_{\pi(i)}^t). \quad (6)$$

An alternative is to perform *tubelet matching*, where all queries with the same spatial index,  $\mathbf{q}^s$ , must match to the same ground truth annotation across all frames of the input video. Here the permutation is obtained over  $S$  elements as

$$\hat{\pi} = \arg \min_{\pi \in \Pi} \frac{1}{|\mathcal{T}|} \sum_{t \in \mathcal{T}} \mathcal{L}_{\text{frame}}(\mathbf{y}^t, \hat{\mathbf{y}}_{\pi(i)}^t). \quad (7)$$

Intuitively, tubelet matching provides stronger supervision when we have full tubelet annotations available. Note that regardless of the type of matching that we perform, the loss computation and the overall model architecture remains the same. Note that we do not weight terms in Eq. 5, for both matching and loss calculation, for simplicity, and to avoid having additional hyperparameters, as also done by [38].

### 3.4. Discussion

As our approach is based on DETR, it does not require external proposals nor non-maximal suppression for post-processing. The idea of using DETR for action localisation has also been explored by TubeR [66] and WOO [9]. There are, however, a number of key differences: WOO does not detect tubelets at all, but only actions at the center keyframe. We also factorise our queries in the spatial and temporal dimensions (Sec. 3.2) to provide inductive biases urging spatio-temporal association. Moreover, we predict action classes separately for each time step in the tubelet, meaning that each of our queries binds to an actor in the video. TubeR, in contrast, parameterises queries such that they are each associated with separate actions (features are average-pooled over the tubelet, and then linearly classified into a single action class). This choice also means that TubeR requires an additional “action switch” head to predict when

tubelets start and end, which we do not require as different time steps in a tubelet can have different action classes in our model. Furthermore, we show experimentally (Tab. 1) that TubeR’s parameterisation obtains lower accuracy. We also consider two types of matching in the loss computation (Sec. 3.3) unlike TubeR, with “tubelet matching” designed for predicting more temporally consistent tubelets. And in contrast to TubeR, we experimentally show how our decoder design allows our model to accurately predict tubelets even with weak, keyframe supervision.

Finally, TubeR requires additional complexity in the form of a “short-term context module” [66] and the external memory bank of [61] which is computed offline using a separate model to achieve strong results. As we show experimentally in the next section, we outperform TubeR without any additional modules, meaning that our model does indeed produce tubelets in an end-to-end manner.

## 4. Experimental Evaluation

### 4.1. Experimental set-up

**Datasets** We evaluate on four spatio-temporal action localisation benchmarks. AVA and AVA-Kinetics contain sparse annotations at each keyframe, whereas UCF101-24 and JHMDB51-21 contain full tubelet annotations.

AVA [18] consists of 430, 15-minute video clips from movies. Keyframes are annotated at every second in the video, with about 210 000 labelled frames in the training set, and 57 000 in the validation set. There are 80 atomic actions labelled for every actor in the clip, of which 60 are used for evaluation [18]. Following standard practice, we report the Frame Average Precision (fAP) at an IoU threshold of 0.5 using the latest v2.2 annotations [18].

AVA-Kinetics [28] is a superset of AVA, and adds detection annotations following the AVA protocol, to a subset of Kinetics 700 [7] videos. Only a single keyframe in a 10-second Kinetics clip is labelled. In total, about 140 000 labelled keyframes are added to the training set, and 32 000 to the validation sets of AVA. Once again, we follow standard practice in reporting the Frame AP at a 0.5 IoU threshold.

UCF101-24 [49] is a subset of UCF101, and annotates 24 action classes with full spatio-temporal tubes in 3 207 untrimmed videos. Note that actions are not labelled exhaustively as in AVA, and there may be people present in the video who are not performing any labelled action. Following standard practice, we use the corrected annotations of [47]. We report both the Frame AP, which evaluates the predictions at each frame independently, and also the Video AP. The Video AP uses a 3D, spatio-temporal IoU to match predictions to targets. And since UCF101-24 videos are up to 900 frames long (median length of 164 frames), and our network typically processes  $T = 32$  frames at a time, we link together tubelet predictions from our network into full-video-tubes using the same causal linking algo-

Table 1. Comparison of detection architectures on AVA controlling for the same resolution (160p) and training settings. Binding each query to a person, rather than to an action (as done in TubeR [66]) yields solid improvements. We report the mean AP for both ViViT-B and CSN-152 backbones.

	ViViT-B	CSN-152
Query binds to action	23.6	25.7
Ours, query binds to person	<b>26.7</b>	<b>27.8</b>

rithm as [22, 31] for fair comparison.

JHMDB51-21 [21] also contains full tube annotations in 928 trimmed videos. However, as the videos are shorter and at most 40 frames, we can process the entire clip with our network, and do not need to perform any linking.

MultiSports [32] is a recent dataset with full tube annotations for 66 fine-grained action classes in 2 129 untrimmed videos, which are up to 3 676 frames long. It is challenging, as it contains fast motion with multiple actors perform different concurrent actions, and there is a large variance in action durations. Evaluation is the same as for UCF101-24.

**Implementation details** For our vision encoder backbone, we consider both transformer-based (ViViT Factorised Encoder [1]), and convolutional (CSN [55]) backbones. For ViViT, we use the “Base” and “Large” model sizes [11, 12], which are typically first pretrained on image datasets like ImageNet-21K [10] and then finetuned on video datasets like Kinetics [23]. We also use CSN-152 pretrained on Instagram65M [37] and then Kinetics following [66]. Our models process  $T = 32$  frames unless otherwise specified, with  $S = 64$  spatial queries per frame and latent decoder dimensionality of  $d = 2048$ . Exhaustive implementation details and training hyperparameters are included in the supplement. We will also release all code and models upon acceptance.

### 4.2. Ablation studies

We analyse the design choices in our model by conducting experiments on both AVA (with sparse per-frame supervision) and on UCF101-24 (where we can evaluate the quality of our predicted tubelets). Unless otherwise stated, our backbone is ViViT-Base pretrained on Kinetics 400, and the frame resolution is 160 pixels (160p) on the smaller side.

**Comparison of detection architectures** Tab. 1 compares our model, where each query represents a person, and all of their actions (Sec. 3.2) to the approach of TubeR [66] (Sec. 3.4), where there is a separate query for each action being performed. We observe that our parameterisation has a substantial impact, with our method outperforming binding to actions by 3.1 points with a ViViT backbone, and 2.1 points with a CSN backbone on the AVA dataset, therefore motivating the design of our decoder. the supplement shows that this trend is consistent on UCF101-24 and JHMDB too.

Another architectural baseline that we can compare to is that of a two-stage Fast-RCNN model using external person

Table 2. Comparison of independent and factorised queries on the AVA, UCF101-24 and JHMDB51-21 datasets. Factorised queries are particularly beneficial for predicting tubelets, as shown by the VideoAP on UCF101-24 and JHMDB51-21 which has full tube annotations. Both models use tubelet matching in the loss.

Query	AVA		UCF101-24			JHMDB51-21		
	fAP	fAP	vAP20	vAP50	vAP50:95	fAP	vAP20	vAP50
Independent	25.2	85.6	86.3	59.5	28.9	85.0	88.5	85.2
Factorised	<b>26.3</b>	<b>86.5</b>	<b>87.4</b>	<b>63.4</b>	<b>29.8</b>	<b>86.9</b>	<b>89.5</b>	<b>88.2</b>

Table 3. Comparison of independent and tubelet matching for computing the loss on AVA, UCF101-24 and JHMDB51-21. Tubelet matching helps for tube-level evaluation metrics like the Video AP (vAP) on UCF101-24 and JHMDB51-21. Note that tubelet matching is actually still possible on AVA as the annotations are at 1fps with actor identities.

Query	AVA		UCF101-24			JHMDB51-21		
	fAP	fAP	vAP20	vAP50	vAP50:95	fAP	vAP20	vAP50
Per-frame matching	<b>26.7</b>	<b>88.2</b>	85.7	63.5	29.4	86.0	88.3	86.0
Tubelet matching	26.3	86.5	<b>87.4</b>	63.4	<b>29.8</b>	<b>86.9</b>	<b>89.5</b>	<b>88.2</b>

detections from [61], as used by [3, 13, 15, 62]. This baseline using the same ViViT-B backbone achieved a mean AP of 25.2, which is still 1.5 points below our model, emphasising the promise of our end-to-end approach. Note that the proposals of [61] obtain an AP50 of 93.9 for person detection on the AVA validation set. They were obtained by first pretraining a Faster-RCNN [44] detector on COCO keypoints, and then finetuning on the person boxes from the AVA training set, using a resolution of 1333 on the longer side. Our model is end-to-end, and does not require any external proposals generated by a separate model at all.

**Comparison to TubeR** The second row of Tab. 1 using a CSN-152 backbone corresponds to our reimplement of TubeR. By keeping all other training hyperparameters constant, we observe that our query binding provides an improvement of 2.1 mAP points in a fair comparison. Note that we could not use the public TubeR code [65], as it does not reproduce the paper’s results: A higher resolution 256p model achieved only 20 mAP when trained with the public code, whilst it is reported to achieve 31.1. Exhaustive details on our attempts to reproduce TubeR with the authors’ public code is in the supplement.

**Query parameterisation** Tab. 2 compares our independent and factorised query methods (Sec. 3.2) on AVA and UCF101-24. We observe that factorised queries consistently provide improvements on both the Frame AP and the Video AP across both datasets. As hypothesised in Sec. 3.2, we believe that this is due to the inductive bias present in this parameterisation. Note that we can only measure the Video AP on UCF101-24 as it has tubes labelled. We also show in the supplement that these observations are consistent on the JHMDB dataset too.

**Matching for loss calculation** As described in Sec. 3.3, when matching the predictions to the ground truth for loss

Table 4. Our model can predict tubelets even when the ground truth annotations are sparse. We show this by subsampling training annotations from the UCF101-24 dataset. Our model sees minimal performance deterioration even when using only 1/24 or 4% of the annotated frames.

Sampling	Labelled frames	fAP	vAP20	vAP50	vAP50:95
All frames	458 814	86.5	87.4	63.4	29.8
Every 12	39 237	85.2	87.2	63.0	29.3
Every 24	20 243	84.9	86.8	63.2	28.1
One per video	2 284	70.2	77.1	48.5	20.4

computation, we can either independently match the outputs at each frame to the ground truths at each frame, or, we can match entire predicted tubelets to the ground truth tubelets. Tab. 3 shows that tubelet matching does indeed improve the quality of the predicted tubelets, as shown by the Video AP on UCF101-24. However, this comes at the cost of the quality of per-frame predictions, (*i.e.* Frame AP). This suggests that tubelet matching improves the association of bounding boxes predicted at different frames (hence higher Video AP), but may also impair the quality of the bounding boxes predicted at each frame (Frame AP). Note that it is technically possible for us to also perform tubelet matching on AVA, since AVA is annotated at 1fps with actor identities, and our model is input 32 frames at 12.5fps (therefore 2.56 seconds of temporal context) meaning that we have sparse tubelets with 2 or 3 annotated frames.

As tubelet matching helps with the overall Video AP, we use it for subsequent experiments on UCF101-24 and JHMDB51-21. For AVA, we use per-frame matching as the standard evaluation metric is the Frame AP, and annotations are sparse at 1fps.

**Weakly-supervised tubelet detection** Our model can predict tubelets even when the ground truth annotations are sparse and only labelled at certain frames (such as the AVA dataset). We quantitatively measure this ability of our model on UCF101-24 which has full tube annotations. We do so by subsampling labels from the training set, and evaluating the full tubes on the validation set.

As shown in Tab. 4, we still obtain meaningful tube predictions, with a Video AP20 of 77.1, when using only a single frame of annotation from each UCF video clip. When retaining 1 frame of supervision for every 24 labelled frames (which is roughly 1fps and corresponds to the AVA dataset’s annotations), we observe minimal deterioration with respect to the fully supervised model (all Video AP metrics are within 0.7 points). Retaining 1 frame of annotation for every 12 consecutive labelled frames also performs similarly to using all frames in the video clip. These results suggest that due to the redundancy in the data (motion between frames is often limited), and the inductive bias of our model, we do not require each frame in the tube to be labelled in order to predict accurate tubelets.

**Decoder design** Tabs. 5 and 6 analyse the effect of the decoder depth and the type of attention in the decoder (de-

Table 5. Effect of decoder depth on performance on the AVA dataset. Performance saturates at  $L = 6$  layers.

Layers ( $L$ )	0	1	3	6	9
mAP $\uparrow$	23.4	24.6	26.2	26.5	<b>26.7</b>

Table 6. Effect of the type of attention used in the decoder on AVA. Factorised attention is both more accurate and efficient (almost half of the GFLOPs per decoder layer).

Decoder attention	mAP	GFLOPs
Full	26.4	10.5
Factorised	<b>26.7</b>	<b>5.3</b>

Table 7. Increasing the image resolution on the AVA dataset leads to consistent accuracy improvements, primarily on small objects. APs, APm and API denote the AP at 0.5 IoU threshold on small, medium and large boxes respectively following the COCO protocol [33]. AVA videos have a median aspect ratio of 16:10, and we pad the larger side when the aspect ratio is different.

Resolution	mAP	APs	APm	API
$140 \times 224$	25.4	7.2	11.2	27.8
$160 \times 256$	26.7	11.5	12.5	28.7
$220 \times 352$	28.8	12.0	15.1	30.7
$260 \times 416$	29.4	13.3	15.8	31.0
$320 \times 512$	30.0	17.5	16.0	32.0

Table 8. Comparison of pretraining for our models with ViViT-B and ViViT-L backbones on AVA using a resolution of  $160 \times 256$ . Larger models benefit more from additional initial pretraining.

Pretrain	STAR/B	STAR/L
IN21K [10] $\rightarrow$ K400 [23]	26.7	27.0
IN21K [10] $\rightarrow$ K700 [7]	27.3	27.6
CLIP [42] $\rightarrow$ K700 [7]	30.3	<b>36.2</b>

scribed in Sec. 3.2). As seen in Tab. 5, detection accuracy on AVA increases with the number of decoder layers, plateauing at around 6 layers. It is possible to use no decoder layers too: In this case, instead of learning queries  $\mathbf{q}$  (Sec. 3.2), we simply interpret the outputs of the vision encoder (Sec. 3.1),  $\mathbf{x}$ , as our queries and apply the localisation and classification heads directly upon them. Using decoder layers, however, can provide a performance increase of up to 3.3 mAP points (14% relative), emphasising their utility.

Tab. 6 shows that factorised attention in the decoder is more accurate than standard, “full” attention between all queries and visual features. Moreover, it is more efficient too, using almost half of the GFLOPs at each decoder layer.

**Effect of resolution and pretraining** Scaling up the image resolution is critical to achieving high performance for object detection in images [20, 46]. However, we are not aware of previous works studying this for video action localisation. Tab. 7 shows that we do indeed observe substantial improvements from higher resolution, improving by up to 4.6 points on AVA. As expected, higher resolutions help more for detection at small sizes, where we follow the COCO [33] convention of object sizes. Note that AVA videos have a median aspect ratio of 16:10, and we pad the larger side for videos with different aspect ratios.

Similarly, Tab. 8 shows the effect of different pretraining

datasets. Video vision transformers are typically pretrained on an image dataset (like ImageNet-21K [10]), before being finetuned on a video dataset, such as Kinetics [23]. We find that the initial image checkpoint plays an important role, with CLIP [42] pretraining significantly outperforming supervised pretraining on ImageNet-21K [12, 50]. This improvement grows further when using a “Large” backbone. Our observations are consistent with prior works which have shown that CLIP-pretraining outperforms ImageNet-21K and even JFT pretraining [24, 30, 34, 41].

**Qualitative examples** We include example result videos of our proposed model in the supplementary.

### 4.3. Comparison to state-of-the-art

We compare our model to the state-of-the-art on datasets with both sparsely annotated keyframes (AVA and AVA-K), and full tubes (UCF101-24, JHMDB and MultiSports).

**AVA and AVA-Kinetics** Tab. 9 compares to prior work on AVA and AVA-Kinetics. The best previous methods relied on external proposals [3, 54, 60] and external memory banks [40, 62] which we outperform. There are fewer end-to-end approaches, and we outperform these by an even larger margin. Note that though TubeR [66] is a proposal-free approach, their best results are actually obtained with the external memory of [61]. Consequently, we have reported the end-to-end, and external-memory versions of TubeR (“TubeR + LTC”) separately in Tab. 9. Furthermore, as detailed in the supplement, the public TubeR training code produces significantly lower performance (20.0 mAP). Hence, to compare to it, we report results with our reimplementation in addition to the paper’s results. Observe that we outperform TubeR using the same CSN-152 backbone, and then improve further using larger transformer backbones.

We achieve greater relative improvements on AVA-Kinetics, showing that our end-to-end approach can leverage larger datasets more effectively. To our knowledge, we surpass the best previous results on AVA-Kinetics, achieving a Frame AP of 45.1. Notably, we outperform InternVideo [60] and VideoMAE-v2 [57], which are two recent video foundation models using more powerful backbones and larger, private, web-scale video datasets. Note that InternVideo consists of two different encoders, one of which is also initialised from CLIP. And like [60], we achieve our best AVA results by training a model on AVA-Kinetics, and then evaluating it only on the AVA validation set. Finally, we do not use any test-time augmentation, unlike previous works that ensemble results over multiple resolutions and/or left/right flips as shown by the “Views” column.

**UCF101-24** Tab. 10 shows that we outperform prior work on UCF101-24, both in terms of frame-level (Frame AP), and tube-level metrics (Video AP). We achieve state-of-the-art results with a CSN-152 backbone, and improve further

Table 9. Comparison to the state-of-the-art (reported with mean Average Precision; mAP  $\uparrow$ ) on AVA v2.2 [18] and AVA-Kinetics (AVA-K) [28]. Methods using external proposals (i.e. not end-to-end) are also trained on additional object detection and human pose data. Unless otherwise stated, separate models are trained for AVA and AVA-Kinetics. \* denotes the model was trained on AVA-Kinetics and evaluated on AVA. “Res.” denotes the frame resolution of the shorter side. Web-scale foundational models are denoted in grey.  $\dagger$  We also report our reimplemented results for TubeR, as the authors’ public code does not reproduce their reported results, as detailed in the supplement.

	Pretraining	Views	AVA	AVA-K	Res.	Backbone	End-to-end
MViT-B [13]	K400	1	27.3	–	–	MViT	$\times$
Unified [2]	K400	6	27.7	–	320	SlowFast	$\times$
AIA [53]	K700	18	32.3	–	320	SlowFast	$\times$
ACAR [40]	K700	6	33.3	36.4	320	SlowFast	$\times$
TubeR [66] with LTC [61]	IG65M [37]→K400, COCO	2	33.6	–	256	CSN-152	$\times$
MeMViT [62]	K700	–	34.4	–	312	MViT v2	$\times$
Co-finetuning [3]	IN21K→K700, MiT, SSv2	1	32.8	33.1	320	ViViT/L	$\times$
	JFT,WTS→K700, MiT, SSv2	1	36.1	36.2	320	ViViT/L	$\times$
VideoMAE [54]	SSL K700 → Sup. K700.	–	39.3	–	–	ViViT/L	$\times$
InternVideo* [60]	7 different datasets	–	41.0	42.5	–	Uniformer v2	$\times$
VideoMAE v2 [57]	6 different datasets	–	42.6	43.9	–	ViViT/g	$\times$
Action Transformer [28]	K400	1	–	23.0	400	I3D	$\checkmark$
WOO [9]	K600	1	28.3	–	320	SlowFast	$\checkmark$
TubeR [66]	IG65M [37]→K400, COCO	1	31.1	–	256	CSN-152	$\checkmark$
TubeR reimplemented $\dagger$	IG65M [37]→K400	1	29.5	33.6	256	CSN-152	$\checkmark$
STAR/CSN-152 (ours)	IG65M→K400	1	31.4	35.8	256	CSN-152	$\checkmark$
STAR/B (ours)	IN21K→K400	1	30.0	36.6	320	ViViT/B	$\checkmark$
	CLIP→K700	1	33.9	39.1	320	ViViT/B	$\checkmark$
STAR/L (ours)	CLIP→K700	1	39.2	44.5	320	ViViT/L	$\checkmark$
STAR/L (ours)*	CLIP→K700	1	<b>42.5</b>	<b>45.1</b>	420	ViViT/L	$\checkmark$

Table 10. State-of-the-art comparison on datasets with tubelet annotations, UCF101, JHMDB51 and MultiSports.

	Pretraining	UCF101-24				JHMDB51-21			MultiSports			Backbone
		fAP	vAP20	vAP50	vAP50:95	fAP	vAP20	vAP50	fAP	vAP20	vAP50	
ACT [22]	IN1K	67.1	77.2	51.4	25.0	65.7	74.2	73.7	–	–	–	VGG
MOC [31]	IN1K → COCO	78.0	82.8	53.8	28.3	70.8	77.3	77.2	25.2	12.9	0.6	DLA34
Unified [2]	K600	79.3	–	–	–	–	–	–	–	–	–	SlowFast
WOO [9]	K600	–	–	–	–	80.5	–	–	–	–	–	SlowFast
TubeR [66]	IG65M→K400	83.2	83.3	58.4	28.9	–	87.4	82.3	–	–	–	CSN-152
HIT [14]	K700	84.8	88.8	74.3	–	83.8	89.7	88.1	33.3	27.8	8.8	SlowFast
FBIL [8]	K700	86.0	86.9	69.1	–	–	–	–	40.8	30.3	9.9	SlowFast
STAR/CSN-152	IG65M→K400	86.7	87.0	65.4	30.6	<b>93.5</b>	<b>96.3</b>	<b>95.4</b>	45.4	50.1	18.6	CSN-152
STAR/B	IN21K→K400	87.3	88.2	68.6	31.7	86.9	89.5	88.2	46.3	51.6	25.6	ViViT/B
STAR/L	CLIP→K700	<b>90.3</b>	<b>89.8</b>	<b>73.4</b>	<b>35.8</b>	92.1	93.1	92.6	<b>59.3</b>	<b>62.0</b>	<b>36.2</b>	ViViT/L

by scaling up to ViViT-Large, consistent with our results on AVA (Tab. 9). Moreover, note how we substantially outperform TubeR [66] using the same CSN-152 backbone. Our margin of improvement over TubeR grows from vAP20 (+3.7) to vAP50 (+7.0) with the same backbone, showing that our tubelets are more precise, and is in line with our visual observations in the supplement. To our knowledge, we outperform the best previous reported Video AP50 result by 13.2 points. Note that as UCF videos are up to 900 frames, and as our network processes  $T = 32$  frames, we follow prior works and link together tubelets using the same causal algorithm as [2, 22, 31, 47] for fair comparison.

**JHMDB51-21** Tab. 10 shows that we surpass the state-of-the-art on JHMDB. Once again, we significantly outperform TubeR [66], by 13.1 vAP50, with the same CSN-152 backbone. The CSN-152 backbone outperforms ViViT in this case, possibly because this is the smallest dataset and larger backbones can overfit more easily. The videos in this dataset are trimmed (meaning that labelled actions are being performed on each frame), and also shorter. Hence the

Video AP is not as strict as it is on UCF. As videos are at most 40 frames, we set  $T = 40$  in our model so that we process the entire clip at once without needing to link tubelets.

**MultiSports** Finally, Tab. 10 shows that we improve on prior state-of-the-art on MultiSports substantially on both Frame AP and Video AP. Although this is a recent benchmark with few methods to compare against, our improvement is significant due to the challenging nature of the dataset: it contains fast motion, many simultaneous actors and large variation in action duration.

## 5. Conclusion and Future Work

We have presented STAR, an end-to-end spatio-temporal action localisation model that can output tubelets, when either sparse keyframe, or full tubelet annotation is available. Our approach achieves state-of-the-art results on four action localisation datasets for both frame-level and tubelet-level predictions (in particular, we obtain 45.1% mAP on the challenging AVA-Kinetics dataset), outperforming complex methods that use external proposals and memory banks. Future work is to extend our method to open vocabularies.



## References

- [1] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. ViViT: A Video Vision Transformer. In *ICCV*, 2021. 5
- [2] Anurag Arnab, Chen Sun, and Cordelia Schmid. Unified graph structured models for video understanding. In *ICCV*, 2021. 1, 2, 8
- [3] Anurag Arnab, Xuehan Xiong, Alexey Gritsenko, Rob Romijnders, Josip Djolonga, Mostafa Dehghani, Chen Sun, Mario Lučić, and Cordelia Schmid. Beyond transfer learning: Co-finetuning for action localisation. In *arXiv preprint arXiv:2207.03807*, 2022. 1, 2, 6, 7, 8
- [4] Jimmy Lei Ba, Jamie Ryan Kiros, and Geoffrey E Hinton. Layer normalization. In *arXiv preprint arXiv:1607.06450*, 2016. 4
- [5] Fabien Baradel, Natalia Neverova, Christian Wolf, Julien Mille, and Greg Mori. Object level visual reasoning in videos. In *ECCV*, 2018. 2
- [6] 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, 4
- [7] Joao Carreira, Eric Noland, Chloe Hillier, and Andrew Zisserman. A short note on the kinetics-700 human action dataset. In *arXiv preprint arXiv:1907.06987*, 2019. 5, 7
- [8] Keke Chen, Xiangbo Shu, Guo-Sen Xie, Rui Yan, and Jinhui Tang. Foreground/background-masked interaction learning for spatio-temporal action detection. In *ACM MM*, 2023. 8
- [9] Shoufa Chen, Peize Sun, Enze Xie, Chongjian Ge, Jiannan Wu, Lan Ma, Jiajun Shen, and Ping Luo. Watch only once: An end-to-end video action detection framework. In *ICCV*, 2021. 1, 2, 4, 8
- [10] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *CVPR*, 2009. 5, 7
- [11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*, 2019. 5
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 5, 7
- [13] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, 2021. 6, 8
- [14] Gueter Josmy Faure, Min-Hung Chen, and Shang-Hong Lai. Holistic interaction transformer network for action detection. In *WACV*, 2023. 8
- [15] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *ICCV*, 2019. 1, 2, 6
- [16] Rohit Girdhar, Joao Carreira, Carl Doersch, and Andrew Zisserman. Video action transformer network. In *CVPR*, 2019. 2
- [17] Ross Girshick. Fast r-cnn. In *ICCV*, 2015. 1, 2
- [18] Chunhui Gu, Chen Sun, David A Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, et al. Ava: A video dataset of spatio-temporally localized atomic visual actions. In *CVPR*, 2018. 1, 2, 4, 5, 8
- [19] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). In *arXiv preprint arXiv:1606.08415*, 2016. 4
- [20] Jonathan Huang, Vivek Rathod, Chen Sun, Menglong Zhu, Anoop Korattikara, Alireza Fathi, Ian Fischer, Zbigniew Wojna, Yang Song, Sergio Guadarrama, et al. Speed/accuracy trade-offs for modern convolutional object detectors. In *CVPR*, 2017. 7
- [21] H. Jhuang, J. Gall, S. Zuffi, C. Schmid, and M. J. Black. Towards understanding action recognition. In *ICCV*, 2013. 1, 2, 5
- [22] Vicky Kalogeiton, Philippe Weinzaepfel, Vittorio Ferrari, and Cordelia Schmid. Action tubelet detector for spatio-temporal action localization. In *ICCV*, 2017. 1, 2, 5, 8
- [23] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. In *arXiv preprint arXiv:1705.06950*, 2017. 5, 7
- [24] Skanda Koppula, Yazhe Li, Evan Shelhamer, Andrew Jaegle, Nikhil Parthasarathy, Relja Arandjelovic, João Carreira, and Olivier Hénaff. Where should i spend my flops? efficiency evaluations of visual pre-training methods. In *arXiv preprint arXiv:2209.15589*, 2022. 7
- [25] Okan Köpüklü, Xiangyu Wei, and Gerhard Rigoll. You only watch once: A unified cnn architecture for real-time spatiotemporal action localization. In *arXiv preprint arXiv:1911.06644*, 2019. 1, 2
- [26] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 4
- [27] Jie Lei, Tamara L Berg, and Mohit Bansal. Detecting moments and highlights in videos via natural language queries. In *NeurIPS*, 2021. 2
- [28] Ang Li, Meghana Thotakuri, David A Ross, João Carreira, Alexander Vostrikov, and Andrew Zisserman. The avakinetis localized human actions video dataset. In *arXiv preprint arXiv:2005.00214*, 2020. 1, 2, 5, 8
- [29] Dong Li, Zhaofan Qiu, Qi Dai, Ting Yao, and Tao Mei. Recurrent tubelet proposal and recognition networks for action detection. In *ECCV*, 2018. 2
- [30] Kunchang Li, Yali Wang, Yinan He, Yizhuo Li, Yi Wang, Limin Wang, and Yu Qiao. Uniformerv2: Spatiotemporal learning by arming image vits with video uniformer. In *arXiv preprint arXiv:2211.09552*, 2022. 7
- [31] Yixuan Li, Zixu Wang, Limin Wang, and Gangshan Wu. Actions as moving points. In *ECCV*, 2020. 2, 5, 8
- [32] Yixuan Li, Lei Chen, Runyu He, Zhenzhi Wang, Gangshan Wu, and Limin Wang. Multisports: A multi-person video dataset of spatio-temporally localized sports actions. In *ICCV*, 2021. 2, 5

- [33] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft coco: Common objects in context. In *ECCV*, 2014. 7
- [34] Ziyi Lin, Shijie Geng, Renrui Zhang, Peng Gao, Gerard de Melo, Xiaogang Wang, Jifeng Dai, Yu Qiao, and Hongsheng Li. Frozen clip models are efficient video learners. In *ECCV*, 2022. 7
- [35] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. Ssd: Single shot multibox detector. In *ECCV*, 2016. 2
- [36] Xiaolong Liu, Qimeng Wang, Yao Hu, Xu Tang, Shiwei Zhang, Song Bai, and Xiang Bai. End-to-end temporal action detection with transformer. *IEEE Transactions on Image Processing*, 31:5427–5441, 2022. 2
- [37] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens Van Der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, 2018. 5, 8
- [38] Matthias Minderer, Alexey Gritsenko, Austin Stone, Maxim Neumann, Dirk Weissenborn, Alexey Dosovitskiy, Aravindh Mahendran, Anurag Arnab, Mostafa Dehghani, Zhuoran Shen, et al. Simple Open-Vocabulary Object Detection with Vision Transformers. In *ECCV*, 2022. 4
- [39] Megha Nawhal and Greg Mori. Activity graph transformer for temporal action localization. In *arXiv preprint arXiv:2101.08540*, 2021. 2
- [40] Junting Pan, Siyu Chen, Mike Zheng Shou, Yu Liu, Jing Shao, and Hongsheng Li. Actor-context-actor relation network for spatio-temporal action localization. In *CVPR*, 2021. 1, 2, 7, 8
- [41] Junting Pan, Ziyi Lin, Xiatian Zhu, Jing Shao, and Hongsheng Li. St-adapter: Parameter-efficient image-to-video transfer learning. In *NeurIPS*, 2022. 7
- [42] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *ICML*, 2021. 7
- [43] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You only look once: Unified, real-time object detection. In *CVPR*, 2016. 2
- [44] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster r-cnn: Towards real-time object detection with region proposal networks. In *NeurIPS*, 2015. 3, 6
- [45] Hamid Rezaatoughi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union: A metric and a loss for bounding box regression. In *CVPR*, 2019. 4
- [46] Bharat Singh and Larry S Davis. An analysis of scale invariance in object detection snip. In *CVPR*, 2018. 7
- [47] Gurkirt Singh, Suman Saha, Michael Sapienza, Philip HS Torr, and Fabio Cuzzolin. Online real-time multiple spatiotemporal action localisation and prediction. In *ICCV*, 2017. 2, 5, 8
- [48] Lin Song, Shiwei Zhang, Gang Yu, and Hongbin Sun. Tacnet: Transition-aware context network for spatio-temporal action detection. In *CVPR*, 2019. 2
- [49] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. In *arXiv preprint arXiv:1212.0402*, 2012. 1, 2, 5
- [50] Andreas Steiner, Alexander Kolesnikov, Xiaohua Zhai, Ross Wightman, Jakob Uszkoreit, and Lucas Beyer. How to train your vit? data, augmentation, and regularization in vision transformers. *TMLR*, 2022. 7
- [51] Chen Sun, Abhinav Shrivastava, Carl Vondrick, Kevin Murphy, Rahul Sukthankar, and Cordelia Schmid. Actor-centric relation network. In *ECCV*, 2018. 1, 2
- [52] Peize Sun, Rufeng Zhang, Yi Jiang, Tao Kong, Chenfeng Xu, Wei Zhan, Masayoshi Tomizuka, Lei Li, Zehuan Yuan, Changhu Wang, et al. Sparse r-cnn: End-to-end object detection with learnable proposals. In *CVPR*, 2021. 2
- [53] Jiajun Tang, Jin Xia, Xinzhi Mu, Bo Pang, and Cewu Lu. Asynchronous interaction aggregation for action detection. In *ECCV*, 2020. 1, 2, 8
- [54] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In *NeurIPS*, 2022. 2, 7, 8
- [55] Du Tran, Heng Wang, Lorenzo Torresani, and Matt Feiszli. Video classification with channel-separated convolutional networks. In *ICCV*, 2019. 5
- [56] 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. 3, 4
- [57] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. Videomae v2: Scaling video masked autoencoders with dual masking. In *arXiv preprint arXiv:2303.16727*, 2023. 2, 7, 8
- [58] Xiaolong Wang and Abhinav Gupta. Videos as space-time region graphs. In *ECCV*, 2018. 2
- [59] Yuqing Wang, Zhaoliang Xu, Xinlong Wang, Chunhua Shen, Baoshan Cheng, Hao Shen, and Huaxia Xia. End-to-end video instance segmentation with transformers. In *CVPR*, 2021. 2
- [60] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, et al. Internvideo: General video foundation models via generative and discriminative learning. In *arXiv preprint arXiv:2212.03191*, 2022. 2, 7, 8
- [61] Chao-Yuan Wu, Christoph Feichtenhofer, Haoqi Fan, Kaiming He, Philipp Krahenbuhl, and Ross Girshick. Long-term feature banks for detailed video understanding. In *CVPR*, 2019. 1, 2, 5, 6, 7, 8
- [62] Chao-Yuan Wu, Yanghao Li, Karttikeya Mangalam, Haoqi Fan, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. MeMVit: Memory-augmented multiscale vision transformer for efficient long-term video recognition. In *CVPR*, 2022. 1, 2, 6, 7, 8
- [63] Chuhan Zhang, Ankush Gupta, and Andrew Zisserman. Temporal query networks for fine-grained video understanding. In *CVPR*, 2021. 2
- [64] Yubo Zhang, Pavel Tokmakov, Martial Hebert, and Cordelia Schmid. A structured model for action detection. In *CVPR*, 2019. 2

- [65] Jiaojiao Zhao, Yanyi Zhang, Xinyu Li, Hao Chen, Bing Shuai, Mingze Xu, Chunhui Liu, Kaustav Kundu, Yuanjun Xiong, Davide Modolo, et al. Tuber public code. <https://github.com/amazon-science/tubelet-transformer>, 2022. 6
- [66] Jiaojiao Zhao, Yanyi Zhang, Xinyu Li, Hao Chen, Bing Shuai, Mingze Xu, Chunhui Liu, Kaustav Kundu, Yuanjun Xiong, Davide Modolo, et al. TubeR: Tubelet Transformer for Video Action Detection. In *CVPR*, 2022. 1, 2, 4, 5, 7, 8
- [67] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. In *arXiv preprint arXiv:1904.07850*, 2019. 2