AttTrack: Online Deep Attention Transfer for Multi-object Tracking

Keivan Nalaie McMaster University nalaiek@mcmaster.ca

Abstract

Multi-object tracking (MOT) is a vital component of intelligent video analytics applications such as surveillance and autonomous driving. The time and storage complexity required to execute deep learning models for visual object tracking hinder their adoption on embedded devices with limited compute power. In this paper, we aim to accelerate *MOT by transferring the knowledge from high-level features* of a complex network (teacher) to a lightweight network (student) at both training and inference time. The proposed AttTrack framework has three key components: 1) crossmodel feature learning to align intermediate representations from the teacher and student models, 2) interleaving the execution of the two models at inference time, and 3) incorporating the updated predictions from the teacher model as prior knowledge to assist the student model. Experiments on pedestrian tracking tasks are conducted on the MOT17 and MOT15 datasets using two different object detection backbones YOLOv5 and DLA34 show that AttTrack can significantly improve student model tracking performance while sacrificing only minor degradation of tracking speed.

1. Introduction

Multi-object tracking (MOT) is a fundamental problem in computer vision that aims to identify trajectories in video frames. MOT is a key building block of many applications such as human-computer interaction, surveillance, and autonomous driving tasks. Existing methods handle a MOT task by treating it as two sub-problems: *object detection* and *object association*. Object detection identifies bounding boxes of objects of interest in each frame whereas object association links detected objects to form trajectories over time.

As neural network models for MOT become more complex, improved accuracy usually comes at the cost of longer Rong Zheng McMaster University rzheng@mcmaster.ca

inference time. To accelerate the execution of deep models, techniques such as model quantization [27, 17, 2] and pruning [20, 25, 41] have been widely utilized to reduce computations and redundant connections. Computation acceleration from model quantization is generally hardware dependent [12]. Extensive parameter tuning is required for network pruning to work without significant loss of accuracy [5]. Recently, a few works utilized temporal[6, 24, 22] and spatial [35] correlations by finding configurations (e.g., frame rate, frame resolution) that achieve a good trade-off between computation complexity and model performance. Unfortunately, the optimal configuration is not only input sensitive but also dependent on run time environments such as the available CPU or memory resources.

Another line of work to reduce model complexity is through knowledge distillation (KD) [38]. KD uses soft labels generated by a large model (teacher) to train a small neural network with fewer parameters (student). The soft labels provide useful information that allows the student model to learn the behavior of the teacher to improve generalization. However, the lack of distilled knowledge for the student model during inference may hinder it from correctly detecting harder instances (e.g., crowded scenes with smaller objects).

In this work, we propose an *attention transfer approach for object tracking* aiming to exploit the knowledge of a teacher model at both training and inference stages. The proposed online deep attention transfer network (AttTrack), is inspired by the idea of attention transfer first proposed in [38]. Unlike existing tracking methods, our detection model receives additional information in the form of previously detected objects from the teacher model. We only run the teacher model every few frames (called *key frame*) during inference to improve the representation capability of the student model and help it to discover likely object positions in the remaining frames. The student model can gain information about obstructed or barely visible objects by leveraging extracted information from outputs of the teacher model on prior frames at no extra cost. The student model is trained to fuse extracted features of input frames and the detection estimation from the teacher model.

Extensive experiments show that AttTrack improves the tracking performance of small models with marginal increases at interference time. Choosing the intervals to run the teaching model, results in different trade-offs between performance and efficiency. In summary, the main contributions of this paper are as follows:

- We conduct an empirical study to investigate the impacts of model size on tracking performance and speed.
- We propose an end-to-end trainable AttTrack framework to transfer the knowledge of a complex teacher model to a lightweight student model at both training and inference time.
- Extensive testing on the MOT17 and MOT15 datasets demonstrates the effectiveness of AttTrack. For example, our approach can improve the MOTA score on YOLOv5 and DLA34 architectures by up to 12% with comparable computation time when compared to existing attention-based methods.

The remainder of the paper is structured as follows. We first describe related work in Section 2 and study the impact of model size in tracking performance and speed in Section 3. In Section 4 details of AttTrack is presented. We present experimental results in Section 5, followed by a conclusion in Section 6.

2. Related Work

2.1. Knowledge distillation

Knowledge Distillation (KD) was first proposed by Hinton et al. [14], which aims to train a student (a smaller and faster) model by transferring knowledge from a teacher (a bigger and slower) model. This knowledge in the form of softened outputs of the teacher model is more informative than one-hot vectors in training data. Subsequent studies improved upon [14] and devised various ways to ease the training of small models with few trainable parameters. Romero et al. [31] proposed FitNet, which uses the intermediate representation learned by a teacher model to change the structure of a small model from being wide and shallow to narrow and deep. Knowledge transmission was considered as a distribution matching problem in [16]. Despite success of KD in classification problems, the needs for bounding box regression and heatmaps estimation in object detection introduce additional obstacles. To extend KD to object detection Chen et al. [3] included a feature imitation loss into the detection loss to use the intermediate features of the teacher as hints for the student model. In [21], the authors devised MIMIC, an extended KD for detector models by employing a fully convolutional feature mimic architecture to transfer knowledge for each pixel individually. In order to avoid teacher supervision for background regions, Mehta *et al.* [23] introduced *objectness scaled distillation* for one-shot object detectors. Similarly, our method uses the attention mechanism of the teacher to train the student model on softer labels. However, the key distinction between our AttTrack and the aforementioned approaches is that during inference time, the student model uses real teacher outputs, to calibrate tracking outputs and achieve better accuracy.

2.2. Attention mechanisms in object tracking

There is a long line of studies that combine the concept of attention with machine learning. Human attention mechanism theories [30] inspired early efforts on attentionbased learning such as [19, 8]. Attention has been used in a wide range of machine learning tasks including deep learning-based video object tracking. Fiaz et al. [10] proposed a channel attention method that gives higher weights to channels that help with target classification and localization. Huang et al. [15] proposed an attentional online update paradigm for siamese visual tracking to improve the performance of a tracker by utilizing knowledge extracted from prior tracking tasks. In [32] residual attention modules are introduced in similarity tracking at multiple levels of feature representation, resulting in improved discrimination quality for similarity searching. Zhang et al. [40] created an attention retrieval network that uses learning masks to conduct soft spatial constrains on features from a tracking backbone network, mitigating the impact of background clutter.

2.3. Trainable attention mechanism for object detection

Researchers have explored attention mechanisms in object detection to enhance feature representation. The encoder and decoder stages of the object detecting system presented in [7] use a dynamic attention approach. The attentions are determined by size, feature dimension, and spatial features using a convolution-based encoder. In [36] the authors proposed a feature pyramid network to object detection in remote sensing images, adapting two types of attention mechanisms: a) global spatial attention that extracts spatial location-related features to improve the positioning ability of the object detector, and b) pixel feature attention that expands the size of receptive fields that makes the model learn more image details. Reverse attention was explored in [4] to assist top-down side-output residual learning in order to acquire more accurate residual features and handle missing object areas and details. In [33], Wang et al. applied a pyramid attention module in their deep saliency model to give more weight to salient regions while extracting multi-scale characteristics from input images. In contrast to the previous researches, to achieve domain sensitivity in object detection, Wang *et al.* [34] utilized a domain attention module for universal object detection.

AttTrack is orthogonal to the previous attempts and can be used in conjunction with other attention-based methods. It takes advantage of the knowledge of more complex teacher models at both training and inference time.

3. Preliminary Study of Model Sizes on Tracking Performance

To motivate our approach we first conduct empirical evaluations on the MOT17 dataset to assess the performance of tracking models with different model sizes.

3.1. Computation time break-down

Object detection, feature extraction, and object association are the three components of a conventional deep neural network (DNN)-based tracker. In [26], it has been reported that object detection is the most time-consuming part of the tracking process. Therefore, this study focuses on reducing the computing cost of object detection while maintaining the overall tracking performance.

3.2. Experimental Setting

Object detection. There are two main categories of approaches in DNN-based object detection. *Two-stage approaches* first extract regions of interest (RoIs) and then classify and regress the RoIs. R-CNN [13] and Faster-RCNN [29] are two widely used object detection models in this category. In the second category, *one-stage approaches*, directly identify and regress objects of interest. For example, YOLO [28] divides an input image into $S \times S$ grids and performs region classification and regression.

In this work, we choose FairMOT [39], a state-of-art (SOTA) one-stage object detection model for three reasons. First, one-stage object detectors tend to be faster than two-stage object detectors. Second, with the YOLOv5 [18] backbone, FairMOT results in a good trade-off between speed and computation complexity. Third, for each identified object, FairMOT computes re-ID features, which can be utilized in object association and tracking.

Model size. We evaluate three models following the YOLOv5 architecture but with different sizes: YOLOv5, YOLOv5-mid (a model with half of the channels in each layer of the base model), and YOLOv5-small (a model with a quarter of the number of channels in each layer of the base model). All three models are pre-trained on the COCO dataset (for object detection task) and then trained on the MOT-17 dataset (for multi-object tracking task).

Performance metric. MOTA and IDF1(F1) are commonly used in MOT to assess tracking performance. An

MOTA score is calculated as follows:

$$MOTA = 1 - \frac{\sum_{t} FP_t + FN_t + IDSW_t}{\sum_{t} GT_t}, \quad (1)$$

where t denotes the frame index, FP_t , FN_t , and $IDSW_t$ denote the number of false positive, false negative, and ID switched objects, respectively, and GT denotes the number of ground truth bounding boxes.

Experiments are conducted on an NVIDIA GTX 3080 graphical card with 8 GB GDDR6, running a docker on Ubuntu 20.04. The system is built using Pytorch v1.8 and CUDA v11.3.

3.3. Results and Observations

The tracking performance of the three models on the validation dataset is shown in Table 1. It is clear that lowering model size leads to a decrease in tracking accuracy but accelerated inference (measured in frame per second (FPS)). YOLOv5-mid, for example, is faster than the base model at the cost of a 5.5% drop in MOTA score. Similarly, YOLOv5-small suffers around a 23.5% drop in MOTA but is 1.35 times faster than the base model.

Figure 1 illustrates the tracking performance for the full model and the small model. Compared to the full model, the small model fails to detect some objects, especially those that are far away, partially occluded, or small sizes. Therefore, our main goal is to train an efficient small neural network using knowledge from a big model to attain comparable performance. Unlike existing work on attention-based approaches, knowledge transfer is performed both during the training and inference stages.

Table 1: Impact of model size on tracking performance

Model	IDF1(%)↑	MOTA(%)↑	FPS↑	Parameter size
YOLOv5	65.90	62.40	43.93	5.01 M
YOLOv5-mid	63.20	56.90	46.16	1.38 M
YOLOv5-small	44.70	38.90	59.32	0.31 M

4. The AttTrack Framework

Figure 2 shows the schematics of the proposed online attention transfer approach. AttTrack employs a teacher model to accurately detect objects from every K frames at the inference time, and a student model combines this knowledge in its tracking model in the interim K-1 frames as depicted in Figure 2.

We formulate the video based object detection and tracking problem as follows: given a set of N input frames $X = \{x_1, x_2, ..., x_N\}$ where $x_n \in \mathcal{R}^{3 \times H \times W}$, the objective is to first obtain set of M_i bounding boxes $B_i = \{b_{i,1}, b_{i,2}, ..., b_{i,M_i}\}$, where $b_{i,j} = \{rect_{i,j}, \phi_{i,j}\}$, $rect_{i,j}$ denotes the 4-dim vector (center coordinates, height and



(a) Extracted trajectories for MOT17-04 using YOLOv5-small (b) Extracted trajectories for MOT17-04 using YOLOv5 model

Figure 1: Demonstration of tracking results from a small and a large model. The YOLOv5-small model performs less accurately than the YOLOv5 model due to partially occluded or small-size objects.



(a) Teacher/student detection performance comparison: a powerful teacher model vs a small student model



(b) Inference stage: attention transfer for two cycles

Figure 2: Schematic illustration of attention transfer. The teacher model is used every K = 4 frames.

width) associated with the *j*th bounding box and $\phi_{i,j}$ represents the extracted visual features of the bounding box in frame *i*. Consequently, object tracking aims to construct the set of trajectories $T_t = \{\mathcal{B}_t, id_t\}$, where \mathcal{B}_t is the set of detected bounding boxes in trajectory *t*, id_t denotes the trajectory ID.

4.1. Online Attention Transfer

Modern object detectors such as FairMOT output heatmaps in addition to bounding boxes, where the value of each pixel in the heatmap is its likelihood of being an object center. Let the heatmap output by the teacher for keyframe k be h_k^t :

$$h_k^t = \mathcal{H}_t(x_k) \tag{2}$$

where \mathcal{H} represents the function associated with the heatmap head of the teacher model. Then, we denote the student model output y_{k+i}^s at frame k + i as below:

$$y_{k+i}^{s} = g(f(x_{k+i}), \Phi(h_{k}^{t}, i))$$
(3)

where *i* stands for number of frames after the keyframe k and $\Phi(h_k^t, i)$ extrapolates the the heatmap of teacher in frame k to get its heatmap in frame k + i, f is the generated features by backbone of the student model and g is the fusion function to be explained in Section 4.3.

Figure 3 depicts system architecture of AttTrack. A non-key frame is processed by the student model, to generate intermediate features. The student model then incorporates updated attention features based on the heatmap of teacher on the most recent keyframe. With the fused features, bounding box regression and re-Id networks are applied to generate the bounding box and re-ID features of each object. During tracking, a trajectory is constructed from detected objects that are similar in appearance-based re-ID features and have a large intersection of union (IoU). Specifically, object association is done in two steps: first, visual features are used to match a trajectory and a detected object. Second, if a match is found, the IoU measure is applied to determine whether a true match is obtained. If the object is matched to a trajectory, the trajectory is extended; otherwise, a new trajectory is initiated. Cosine similarity is used in computing the similarity of visual-based features. A Kalman filter [39] is then applied to update the position state of each trajectory in the current frame.

Since the teacher model is applied for frames between two keyframes, the heatmaps (attention) of teacher are outdated for any frame in-between. To extrapolate the heatmap of teacher for these frames, we devise an attention update approach next.



Figure 3: System architecture: AttTrack applies a teacher model every K frames, and computes updated states $\Phi(h_k^t, i)$ (attention) for intermediate frames based on the teacher's predicted heatmap h_k^t for frame k. The updated state is then fused with the prediction on frame k + i by a student model using a fusion network, resulting in object bounding boxes and re-identification features for frame k + i.

4.2. Attention Update

The knowledge computed in earlier frames by the teacher is beneficial to the student model. However, due to the presence of moving objects, such information becomes more outdated as the time elapses between the current frame and the most recent key frame. Therefore, updates need to be made from the teacher heatmap (Figure 4). Consider B_k^t the set of objects detected by the teacher model in frame k. We first estimate the velocity of each identified object based on bounding box locations from previous frames. The velocity is then used to predict the subsequent locations of the corresponding object in frame k + i, i = 1, 2, ..., K - 1 using a simple linear kinematic equation. The heatmaps are updated accordingly.

The updated heatmaps are most beneficial when the student model fails to detect an object due to poor visibility. However, when object movements are irregular, a new object enters the scene or an object exits the scene, the information of teacher can still be stale. Therefore, the updated heatmaps should be combined with the prediction from the student model for the current frame.

4.3. Network Design

For the teacher model, we can utilize any existing object detection backbone such as DLA34 [37]. The student model, like the teacher model, creates bounding boxes and re-ID features for observed objects. For faster computations, the student model employs fewer parameters than the teacher model in its network backbone. The student model receives attention features and input image as inputs and fuses the attention features and its own calculated features as:

$$g(f(x_{k+i}), \Phi(h_k^t, i)) = (f(x_{k+i}), \Phi(h_k^t, i))$$
(4)

where fusion function g appends the extrapolated features from teacher with new features calculated by the student



(a) States of outputs of teacher are updated using linear kinematic equation



(b) Inference stage: attention transfer

Figure 4: Attention state update. We choose K = 4 as number of frames between every keyframes. At frame F_k and F_{k+4} teacher model is performed whereas between these two keyframes the student model is applied.

backbone. The fused features are fed into heatmap and re-ID branches as defined in [39].

The student model is trained using the following loss function:

$$\mathcal{L}_{student} = \frac{1}{2} \left(\frac{1}{e^{\omega_1}} (\mathcal{L}_{heatmap} + \mathcal{L}_{box}) + \frac{1}{e^{\omega_2}} \mathcal{L}_{identity} + \omega_1 + \omega_2 \right),$$
(5)

which consists of learnable task based parameters ω_1 and ω_2 , heatmap loss, box-size loss, and re-identification loss defined in [39].

4.4. Cross-model Feature Learning

Switching from one model (for example, teacher) to another (for example, student) can result in re-ID features that follow different distributions. Running AttTrack on a video clip necessitates multiple transitions between the teacher



Figure 5: Student network architecture. The student model has fewer parameters and takes two inputs: the input-frame at time t and estimated heatmaps of teacher up to time t.



Figure 6: Cross model feature learning. In the object association step, pivot features for EFM and re-ID features for IFM are used. The re-ID features of the student model in IFM are aligned with the teacher model.

and student models. When re-ID features mismatch between teacher and predictions of student, identity fragmentation occurs, leading to reduced tracking accuracy. To mitigate the domain gap between the re-ID features produced by the student and the teacher models, we propose two crossmodel feature learning approaches.

Explicit Feature Mapping (EFM): We use pivot features to induce correlation between the computed re-ID features by the teacher and the student models. This is done by applying a single linear layer to map the re-ID features to the number of identities (encoded as a one-hot vector) in the training set. Both student and teacher models are subjected to the linear layer. By mapping each identity to those learned pivots in the training time, this approach lowers the distance across two model domains in the inference time.

Implicit Feature Mapping(IFM): the former method requires an additional compute unit in both the student and teacher models, resulting in increased total computation costs. In the second approach, we perform feature mapping implicitly for the student model by mincing the re-ID feature layer in the teacher model. During training, the extracted features from the teacher model are used as an additional supervision signal and the loss function is updated as:

$$\mathcal{L}_W = \mathcal{L}_{student}(W) + \mathcal{L}_{id-att}(W).$$
(6)

where $\mathcal{L}_{student}$ is the loss function in Eq.5 and \mathcal{L}_{id-att} is L2 loss of re-ID features.

5. Performance Evaluation

To qualify the performance of AttTrack, we conduct experiments on a pedestrian tracking task. The MOTA score is used to evaluate tracking accuracy, while FPS is used to quantify tracking speed.

5.1. Datasets

We evaluate AttTrack on pedestrian tracking tasks on two MOT datasets. We use 11 training video clips in the publicly accessible dataset MOT17. We also provide the results of AttTrack on MOT15 [1], which in addition to low-resolution street view videos it includes videos from PETS [9] and KITTI [11] datasets. Since the ground truth of test sequences is not made public, each video is split into two halves with the first half in the training set and the second half in the test set. We also use the COCO dataset to pre-train the models. To evaluate the tracking performance of our models, we utilize the official evaluation method from MOT Challenge [1].

5.2. Implementation Details

Experiments are conducted using the same hardware and software setup as in Section 3.2. We use the Adam optimizer to train our model across 35 epochs, with a starting learning rate of 1e-5 that lowers every 25 epochs. A batch size of 12 is used. Rotation and scaling are applied to augment the training set. The input frame size is 1088×608 pixels.

To evaluate the performance of AttTrack, we consider two backbone models for the student model. The first is DLA34-small, which offers a good trade-off in tracking performance and speed, and is based on DLA34 used in [39]. It has in total 16.55M parameters. The second one, YOLOv5small, has the same architecture as YOLOv5 [18] but with only one-fourth of the parameters, allowing for fast computations and acceptable performance.

5.3. Baseline Methods

To evaluate AttTrack, we consider the following baselines: *Teacher-only* and *Student-only*: object detection in the tracking pipeline only uses teacher and student models, respectively. *Naive-Mix*: which alternates between a teacher and a student model every K frames and merge the tracking outputs of the two with no further information sharing in training or inference; AttTrack w/o attention update (*AttTrack-no-update*), which transfers attention from Teacher to Student at inference but does not update the attention (i.e., $\Phi(h_k^t, i) = h_k^t$); and *Layerwise*, which is similar to Naive-Mix but allows layerwise attention transfer from the teacher to the student model in training[38].

5.4. Online-based Attention Transfer

Table 2 summarizes the performance of teacher models with full-fledged DLA34 and YOLOv5 backbones. The results for the student models with and without AttTrack, AttTrack w/o attention update are given in Table 3. In Table 3, with or without AttTrack, the teacher model is executed every 6 frames. The difference between the two lies in whether attention transfer is performed or not. Similar to the results in Section 3, smaller models have fast inference time but suffer from lower accuracy. DLA34-small without AttTrack, for example, is $1.52 \times$ faster at the price of 6.9 percent MOTA degradation. AttTrack improves the tracking performance of the student model by 1.6% and 5% for DLA and YOLOv5, respectively. This shows the effectiveness of attention transfer from the teacher model. Because AttTrack invokes the teacher model every 6 frames, the running time of AttTrack is longer than those without. In Table 3, we further compare AttTrack and AttTrack-noupdate. As expected, AttTrack-no-update is faster due to less computation but has slightly degraded performance. The relative small gap between the two can attributed to small changes in the scenes when K = 6.

To better understand the impact of K on AttTrack, Table 4 lists the results of different student models under various Ks. As expected a smaller K means more frequent execution of the teacher model, resulting in slower processing time and more accurate tracking outputs, and vice versa. YOLOv5-small runs faster than its DLA34 counterpart but with lower accuracy.

 Table 2: Performance of Teacher-only and Student-only baselines on the MOT17 dataset

Baseline	Model	MOTA (%)	FPS
Taaahar only	DLA34	68.30	20.78
reacher-only	YOLOv5	62.40	40.46
Student only	DLA34	61.40	37.69
Student-Only	Model MOTA (%) DLA34 68.30 YOLOv5 62.40 DLA34 61.40 YOLOv5 38.90	38.90	59.32

Table 3: AttTrack model experiments with K = 6 on the MOT17 dataset

Model	AttTrack		AttTrack-no-update		Naive-mix	
	MOTA	FPS	MOTA	FPS	MOTA	FPS
DLA34	63.00	30.80	62.90	31.30	61.40	31.65
YOLOv5	43.60	50.90	43.40	52.24	38.60	53.13

Table 4: YOLOv5 and DLA34 models with IFM on the MOT17 dataset

K	YOLOv	5-small	DLA34-small		
	MOTA	FPS	MOTA	FPS	
2	48.50	45.09	64.50	26.24	
4	43.90	49.96	63.20	29.28	
6	43.60	50.91	63.00	30.80	

5.5. Alternative teacher

We conduct further investigations to see whether the representational power of teachers can affect the performance of the student model. Specifically, we compare the use of DLA34 in the teacher model and transfer the knowledge to YOLOv5-small student model. The heatmap computed by a DLA34 teacher can still be useful to the YOLOv5-small student model, and the re-Id features can be aligned using the mechanism in Section 4.4.

The DLA34 teacher provides better tracking performance than the YOLOv5-based equivalent, as shown in Table 5, although it runs slower than the YOLOv5 teacher. The gap in tracking performance between YOLOv5 and DLA teacher-based models reduces as K increases as the impact of YOLOv5-small becomes more dominant. Overall, the results in Table 5 show that tracking performance and processing time can be considerably impacted by the choice of the teacher architecture.

 Table 5: Compression of Different Teacher Models using EFM on the MOT17 dataset.

Κ	YOLOv5	$5 \rightarrow$ YOLOv5-small	$DLA34 \rightarrow YOLOv5$ -small		
	MOTA	FPS	MOTA	FPS	
2	50.40	44.69	52.00	28.72	
4	46.90	47.60	47.70	36.41	
6	45.90	48.83	46.20	40.32	

Table 6: Importance of cross-model feature learning on the MOT17 dataset. EFM: employing an additional convolution layer to translate characteristics from the teacher and student models to the common features space. IFM: student model mimics re-ID feature generated by the student model.

Model		EE	M	IE	M	No Eoo	Loorning
WIGUEI	K		IVI	<u>111</u>	.v1	NO Fea.	Leanning
		MOTA	FPS	MOTA	FPS	MOTA	FPS
	2	65.30	25.40	64.50	26.24	64.30	26.30
DLA34	4	64.10	28.36	63.20	29.28	63.20	29.59
	6	63.80	29.99	63.10	30.80	63.20	31.01
	2	50.40	44.69	48.50	45.09	47.50	45.57
YOLOv5	4	46.90	47.60	43.90	49.96	43.20	50.22
	6	45.90	48.83	43.60	50.91	42.90	51.03

5.6. Cross-model feature learning

The usefulness of the learned features for transferring knowledge between teacher and student models is evaluated in Table 6. In the experiments, EFM is done by applying a single linear layer on the generated re-ID features. As can be observed, the inclusion of this extra layer reduces the visual feature distance between the teacher and the student, and produces more precise tracking output. The use of EFM for K = 4 has the greatest influence on the YOLOv5 student model, accounting for 3.7% of more accurate tracking performance. As we can see, cross-model feature learning is beneficial and has more impacts on YOLOv5 than it does on DLA34-small. Furthermore, EFM yields better tracking performance than IFM but incurs higher computation costs. On DLA34-small, for instance, utilizing EFM with K = 2 achieves a 65.30% MOTA score and 25.40 frame rate, whereas the use of IFM results in a 0.9% lower MOTA score and a 0.84 faster FPS.

5.7. Comparison with layer-wise attention transfer

In this set of experiments, we compare AttTrack with the layer-wise attention transfer proposed in [38] (baseline Layerwise). The main difference between our approach and [38] is that the layer-wise solution transfers attention knowledge to the student model during the training time only, and the student model performs tracking entirely on its own, while our AttTrack leverages teacher knowledge in both training and inference phases. We implement a layerwise attention approach for the MOT task since [38] is originally built for object classification tasks. The results are shown in Figure 7 for 11 different Ks between two and twelve. For fair comparison, in the baseline layer-wise attention transfer, we also invoke the teacher model every Kframes though there is no knowledge transfer between the teacher and the student models at inference time. 2nd order polynomial fitting functions for the AttTrack and baseline results are also displayed in the figures. When comparing AttTrack to layer-wise attention, we find that AttTrack exceeds the baseline significantly with comparable processing time on tracking accuracy. The difference with YOLOv5 is more pronounced. For example, AttTrack is 4 percent better with only 2 percent lower FPS when K is between two and four. The gap in computation time between AttTrack and baseline drops for the DLA34-based tracker. This is because the overhead of attention transfer in AttTrack is shadowed by the high compute cost of the DLA34 backbone.

5.8. Experiments on MOT15

To verify the generalizability of AttTrack to other datasets, we further conduct experiments on MOT15. The performance of Teacher-only and Student-only is given in Table 7, and the comparison between AttTrack and Layerwise for different K's is given in Table 8. Similar to the trends with MOT17, we observe that AttTrack outperforms Layerwise and Student-only in MOTA, and is considerably



(a) Tracking performance comparison(b) Tracking performance comparison for DLA34 architecture for YOLOv5 architecture

Figure 7: Results of attention transfer for AttTrack and Layerwise under 11 different $Ks \in [2, 12]$ on the MOT17 dataset

faster than Teacher-only.

 Table 7: Performance of Teacher-only and Student-only baselines on the MOT15 dataset

	Baseline	Model	MOTA (%)	FPS
		DLA34	68.80	21.70
	Teacher-only	YOLOv5	61.00	54.41
	Student only	DLA34	66.90	41.47
	Student-only	YOLOv5	52.70	58.80

Table 8: AttTrack and Layerwise on the MOT15 dataset

Model	v	K AttTrack-EFM		Layerwise			
	L L	MOTA	FPS	MOTA	FPS		
-	2	67.60	27.99	65.90	28.64		
DLA34	4	67.40	32.32	65.70	32.55		
	6	67.30	35.19	65.60	35.90		
	2	56.50	55.17	55.00	56.63		
YOLOv5	4	54.10	56.13	52.60	57.93		
	6	52.90	56.46	52.30	58.39		

6. Conclusion

AttTrack is a teacher-student attention transfer approach for accelerating multi-object tracking tasks. It transfers knowledge from a complex teacher to a lightweight student model in both the training and inference stages. AttTrack is model agnostic and can be used in conjunction with other techniques to accelerate neural network inference. Because AttTrack adopts cross-model feature learning, it is capable to transfer knowledge from any teacher to any student network with different network architectures (e.g. YOLOv5 or DLA34). When compared with traditional attentionbased methods, our work improves tracking accuracy with marginal degradation in inference time. As part of future work, we are interested in investigating attention mechanisms with adaptive window sizes.

References

- [1] MOT Challenge Website. https://motchallenge.net.
- [2] Zhaowei Cai, Xiaodong He, Jian Sun, and Nuno Vasconcelos. Deep Learning With Low Precision by Half-wave Gaussian Quantization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5918–5926, 2017.
- [3] Guobin Chen, Wongun Choi, Xiang Yu, Tony Han, and Manmohan Chandraker. Learning efficient object detection models with knowledge distillation. *Advances in neural information processing systems*, 30, 2017.
- [4] Shuhan Chen, Xiuli Tan, Ben Wang, and Xuelong Hu. Reverse attention for salient object detection. In *Proceedings of the European conference on computer vision (ECCV)*, pages 234–250, 2018.
- [5] Yu Cheng, Duo Wang, Pan Zhou, and Tao Zhang. Model Compression and Acceleration for Deep Neural Networks: The Principles, Progress, and challenges. *IEEE Signal Processing Magazine*, 35(1):126–136, 2018.
- [6] Ting-Wu Chin, Ruizhou Ding, and Diana Marculescu. Adascale: Towards Real-time Video Object Detection Using Adaptive Scaling. arXiv preprint arXiv:1902.02910, 2019.
- [7] Xiyang Dai, Yinpeng Chen, Jianwei Yang, Pengchuan Zhang, Lu Yuan, and Lei Zhang. Dynamic detr: End-toend object detection with dynamic attention. In *Proceedings* of the IEEE/CVF International Conference on Computer Vision, pages 2988–2997, 2021.
- [8] Misha Denil, Loris Bazzani, Hugo Larochelle, and Nando de Freitas. Learning where to attend with deep architectures for image tracking. *Neural computation*, 24(8):2151–2184, 2012.
- [9] Anna Ellis and James Ferryman. PETS2010 and PETS2009 evaluation of results using individual ground truthed single views. In 2010 7th IEEE international conference on advanced video and signal based surveillance, pages 135–142. IEEE, 2010.
- [10] Mustansar Fiaz, Arif Mahmood, Ki Yeol Baek, Sehar Shahzad Farooq, and Soon Ki Jung. Improving object tracking by added noise and channel attention. *Sensors*, 20(13):3780, 2020.
- [11] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? the kitti vision benchmark suite. In 2012 IEEE conference on computer vision and pattern recognition, pages 3354–3361. IEEE, 2012.
- [12] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. A survey of Quantization Methods for Efficient neural Network Inference. arXiv preprint arXiv:2103.13630, 2021.
- [13] Ross Girshick, Jeff Donahue, Trevor Darrell, and Jitendra Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 580–587, 2014.
- [14] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling The Knowledge In A Neural Network. *arXiv preprint arXiv:1503.02531*, 2015.

- [15] Bo Huang, Tingfa Xu, Ziyi Shen, Shenwang Jiang, Bingqing Zhao, and Ziyang Bian. Siamatl: online update of siamese tracking network via attentional transfer learning. *IEEE Transactions on Cybernetics*, 2021.
- [16] Zehao Huang and Naiyan Wang. Like what you like: Knowledge distill via neuron selectivity transfer. arXiv preprint arXiv:1707.01219, 2017.
- [17] Itay Hubara, Matthieu Courbariaux, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Quantized Neural Networks: Training Neural Networks With Low Precision Weights and Activations. *The Journal of Machine Learning Research*, 18(1):6869–6898, 2017.
- [18] Glenn Jocher, Alex Stoken, Jirka Borovec, Ayush Chaurasia, L Changyu, AV Laughing, A Hogan, J Hajek, L Diaconu, YK Marc, et al. ultralytics/yolov5: v5. 0-yolov5-p6 1280 models aws supervise. ly and youtube integrations. *Zenodo*, 11, 2021.
- [19] Hugo Larochelle and Geoffrey E Hinton. Learning to combine foveal glimpses with a third-order boltzmann machine. Advances in neural information processing systems, 23, 2010.
- [20] Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning Filters for Efficient Convnets. arXiv preprint arXiv:1608.08710, 2016.
- [21] Quanquan Li, Shengying Jin, and Junjie Yan. Mimicking very efficient network for object detection. In *Proceedings* of the ieee conference on computer vision and pattern recognition, pages 6356–6364, 2017.
- [22] Miaomiao Liu, Xianzhong Ding, and Wan Du. Continuous, Real-Time Object Detection on Mobile Devices without Offloading. In 2020 IEEE 40th International Conference on Distributed Computing Systems (ICDCS), pages 976–986. IEEE, 2020.
- [23] Rakesh Mehta and Cemalettin Ozturk. Object detection at 200 frames per second. In *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, pages 0–0, 2018.
- [24] Yue Meng, Chung-Ching Lin, Rameswar Panda, Prasanna Sattigeri, Leonid Karlinsky, Aude Oliva, Kate Saenko, and Rogerio Feris. Ar-net: Adaptive Frame Resolution For Efficient Action Recognition. In *European Conference on Computer Vision*, pages 86–104. Springer, 2020.
- [25] Pavlo Molchanov, Stephen Tyree, Tero Karras, Timo Aila, and Jan Kautz. Pruning Convolutional Neural Networks for Resource Efficient Inference. arXiv preprint arXiv:1611.06440, 2016.
- [26] Keivan Nalaie, Renjie Xu, and Rong Zheng. DeepScale: Online Frame Size Adaptation for Multi-object Tracking on Smart Cameras and Edge Servers.
- [27] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet Classification Using Binary Convolutional Neural Networks. In *European conference on computer vision*, pages 525–542. Springer, 2016.
- [28] Joseph Redmon, Santosh Divvala, Ross Girshick, and Ali Farhadi. You Only Look Once: Unified, Real-time Object Detection. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 779–788, 2016.

- [29] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: Towards Real-time Object Detection with Region Proposal Networks. *IEEE transactions on pattern* analysis and machine intelligence, 39(6):1137–1149, 2016.
- [30] Ronald A Rensink. The dynamic representation of scenes. *Visual cognition*, 7(1-3):17–42, 2000.
- [31] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. arXiv preprint arXiv:1412.6550, 2014.
- [32] Kokul Thanikasalam, Clinton Fookes, Sridha Sridharan, Amirthalingam Ramanan, and Amalka Pinidiyaarachchi. Target-specific siamese attention network for real-time object tracking. *IEEE Transactions on Information Forensics* and Security, 15:1276–1289, 2019.
- [33] Wenguan Wang, Shuyang Zhao, Jianbing Shen, Steven CH Hoi, and Ali Borji. Salient object detection with pyramid attention and salient edges. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1448–1457, 2019.
- [34] Xudong Wang, Zhaowei Cai, Dashan Gao, and Nuno Vasconcelos. Towards universal object detection by domain attention. In *Proceedings of the IEEE/CVF Conference* on Computer Vision and Pattern Recognition, pages 7289– 7298, 2019.
- [35] Mengwei Xu, Mengze Zhu, Yunxin Liu, Felix Xiaozhu Lin, and Xuanzhe Liu. DeepCache: Principled Cache for Mobile Deep Vision. In Proceedings of the 24th Annual International Conference on Mobile Computing and Networking, pages 129–144, 2018.
- [36] Xiang Ying, Qiang Wang, Xuewei Li, Mei Yu, Han Jiang, Jie Gao, Zhiqiang Liu, and Ruiguo Yu. Multi-attention object detection model in remote sensing images based on multiscale. *IEEE Access*, 7:94508–94519, 2019.
- [37] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep Layer Aggregation. In Proceedings of the IEEE conference on computer vision and pattern recognition, pages 2403–2412, 2018.
- [38] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. *arXiv preprint arXiv:1612.03928*, 2016.
- [39] Yifu Zhang, Chunyu Wang, Xinggang Wang, Wenjun Zeng, and Wenyu Liu. FairMOT: On The Fairness Of Detection And Re-Identification In Multiple Object Tracking. arXiv preprint arXiv:2004.01888, 2020.
- [40] Zhipeng Zhang, Yufan Liu, Bing Li, Weiming Hu, and Houwen Peng. Toward accurate pixelwise object tracking via attention retrieval. *IEEE Transactions on Image Processing*, 30:8553–8566, 2021.
- [41] Hattie Zhou, Janice Lan, Rosanne Liu, and Jason Yosinski. Deconstructing Lottery Tickets: Zeros, Signs, and the Supermask. Advances in neural information processing systems, 32, 2019.