# PromptonomyViT: Multi-Task Prompt Learning Improves Video Transformers using Synthetic Scene Data

Roei Herzig[* 1,3,4], Ofir Abramovich[*2], Elad Ben Avraham[1],
Assaf Arbelle[4], Leonid Karlinsky[5], Ariel Shamir[2], Trevor Darrell[3], Amir Globerson[1]

[1]Tel-Aviv University, [2]Reichman University, [3]UC Berkeley, [4]IBM Research, [5]MIT-IBM Watson AI Lab
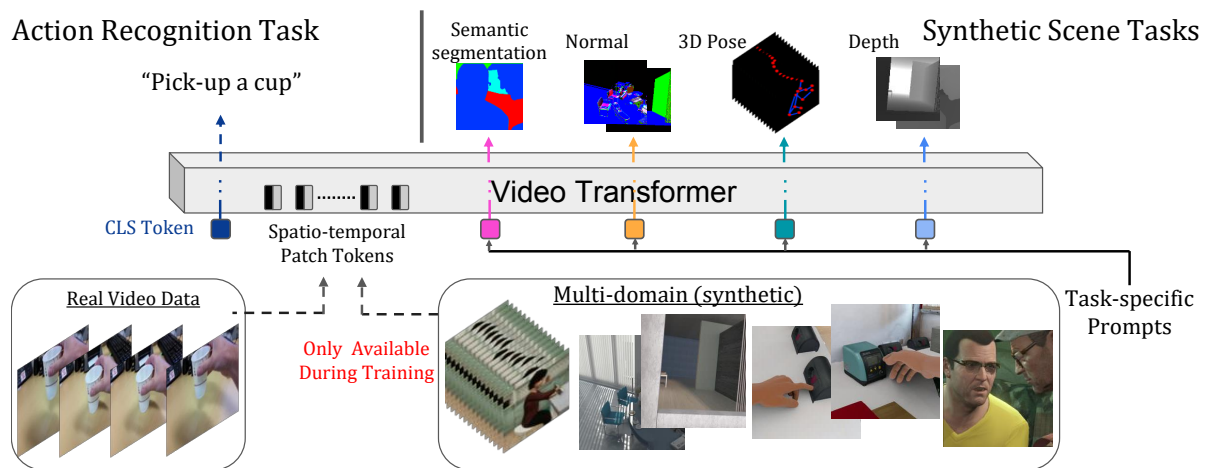
Figure 1. Our PromptonomyViT (PViT) adds a set of multiple prompts to a video transformer to capture inter-task structure and solve a downstream task. We consider the setting where automatically generated synthetic scene data for scene-level tasks (e.g., depth, semantic segmentation) is used for improving an action recognition model on real data. Our PViT model utilizes a multi-task prompt learning approach for video transformers, where a shared transformer backbone is enhanced with task-specific prompts (colored squares). The task prompts predict the synthetic labels for each task, and a CLS token (blue square) is used to predict the action recognition label. The use of task-specific prompts allows the model to benefit from task-related information.

## Abstract

*Action recognition models have achieved impressive results by incorporating scene-level annotations, such as objects, their relations, 3D structure, and more. However, obtaining annotations of scene structure for videos requires a significant amount of effort to gather and annotate, making these methods expensive to train. In contrast, synthetic datasets generated by graphics engines provide powerful alternatives for generating scene-level annotations across multiple tasks. In this work, we propose an approach to leverage synthetic scene data for improving video understanding. We present a multi-task prompt learning approach for video transformers, where a shared video transformer backbone is enhanced by a small set of specialized parameters for each task. Specifically, we add a set of "task prompts", each corresponding to a different task, and let each prompt predict task-related annotations. This design allows the model to capture information shared among synthetic scene tasks as well as information shared between synthetic scene tasks and a real video downstream task throughout the entire network. We refer to this approach as "Promptonomy", since the prompts model task-related structure. We propose the PromptonomyViT model (PViT), a video transformer that incorporates various types of scene-level information from synthetic data using the "Promptonomy" approach. PViT shows strong performance improvements on multiple video understanding tasks and datasets. Project page: https://ofir1080.github.io/PromptonomyViT*

---

*Equal contribution. The order of authors is determined by a coin flip.

# 1. Introduction

Video understanding is a key challenge for machine vision and artificial intelligence. It is intuitively clear that video models should benefit from incorporating spatio-temporal scene-level information including objects, their relations, sizes of instances, 3D structure of a scene, its layout, depth and more. Indeed, several recent studies have explored the use of scene-level information for a variety of video tasks, such as action recognition [21, 26, 40, 65], action detection [50, 112], 3D understanding [1, 18, 76], and structured representations for videos [3, 33, 34, 40, 41, 46, 97, 102]. However, collecting and annotating real large-scale video datasets [36, 52] requires an extensive amount of effort and a large budget. This is especially true for complex labels such as 3D structure and segmentation maps.

In the absence of real-world data, synthetic datasets generated by graphics engines [29, 84] provide a powerful alternative for automatically generating scene-level annotations. Graphics engines can be used to generate a large amount of various types of labeled examples of scene-level information. However, learning from synthetic data requires models that can capture those aspects of the synthetic data that are relevant for downstream tasks, and overcome domain gap issues. An additional challenge is how to benefit from multiple types of scene labels (e.g., depth, normal, segmentation maps, 3D joints positions, and more). In this work, we propose a novel approach that can utilize synthetic data of various sources with multiple types of scene annotations to enhance video understanding models.

Our approach employs Vision Transformers (ViT) [25], which have recently emerged as the leading model for many vision applications [2, 13, 27], including for video understanding [5, 40, 60, 103]. Our key insight is that ViT can be naturally extended to multiple synthetic sources through the use of prompt learning. The key idea of prompt learning methods is to augment the transformer input with a set of additional learnable parameters. The notion of prompt learning has been used successfully in NLP [56], and more recently in machine vision [113, 114]. Inspired by this, we present a prompt learning approach for video transformers, where a shared backbone is enhanced by a small set of specialized parameters for each task. More specifically, we add a set of "task prompts", each dedicated to a unique task. With this design, it is possible to capture information shared among synthetic tasks as well as information shared between synthetic tasks and a real video downstream task, even without applying any domain gap techniques.[1]

The "task prompts" construction can be viewed as implementing "streams of information", each stream representing a task. This facilitates incorporating information from other tasks into the downstream task, starting from early layers and propagating into the spatio-temporal representations throughout the network. We refer to our prompt-per-task approach as "Promptonomy" since the prompts are intended to manage multiple tasks and capture inter-task structure, and name our model PromptonomyViT (PViT).[2] See Figure 1 for an overview.

Recently, the general idea of prompt tuning has been adapted to vision models by VPT [48], suggesting better efficiency of large vision models. Our model differs from recent "prompt tuning" approaches in that we refine a full transformer model rather than optimize a limited set of prompt tokens. As a result, information is propagated from the "task tokens" to all other tokens, enabling interaction across the entire network between the synthetic tasks and the real video downstream task. Furthermore, our multi-task prompts are supervised by auxiliary tasks, and not the primary action recognition task.

To summarize, our main contributions are as follows: (i) we propose a new method for exploiting synthetically generated labels for several tasks to improve video understanding models; (ii) we propose the concept of special "multi-task prompts" to capture task-related information through task supervision, while also interacting with prompts of other tasks and the downstream video task; (iii) we demonstrate improved performance on five tasks and five datasets on video understanding benchmarks: compositional and few-shot action recognition on SomethingElse, spatio-temporal action detection on AVA, standard action recognition on Something-Something V2, Diving48, and PNR Temporal Localization task on Ego4D, highlighting the effectiveness of the proposed approach.

# 2. Related Work

**Prompt Tuning**. Natural language prompting is a method of reformatting NLP tasks as natural language responses to natural language input. Recently, the concept of prompt tuning for efficient fine-tuning of language models was introduced by [56]. Several recent works [4, 82, 95], have explored prompt tuning in the context of multi-task learning in natural language processing. ATTEMPT [4] suggested a soft prompt tuning approach for parameter efficient multi-task knowledge sharing, UNIFIED PROMPT [82] suggested to use multi-task text prompting for zero-shot tasks, and the authors in [95] suggested the soft prompt tuning method for efficient fine-tuning. Additional recent works [48, 99, 100] suggested exploring the usage of prompt tuning in vision transformers. Specifically, VPT [48] uses prompt tuning to efficiently fine tune vision transformers, while others [99, 100] use prompts for continual learning. As opposed to these works, our focus is on the addition of multiple prompts that incorporate various types of scene-level information learned from synthetic data, which will lead to better video understanding. Last, we note that, since our focus is not on efficiency, the entire model is fine-tuned without freezing any parameters.

**Learning from Synthetic Data**. In the field of computer vision, synthetic data has been widely used as an alternative

---

[1] Such techniques may improve performance further, but are orthogonal to our approach.

[2] The name also refers to the classic work on Taskonomy [110], which studied the structure and management of multiple tasks in images.

to real-world training data to solve various problems [24, 29, 70, 75, 81, 84, 93]. Many works attempted to generate synthetic data that mimics real data for image classification [29, 70], semantic segmentation [81, 98], action recognition [24, 93], object detection [74, 75], representation learning [71, 101], and more [84, 104]. Instead, our approach focuses on learning multiple tasks simultaneously from several synthetic domains and then transferring knowledge into the real world task by developing a multi-task prompting model and training scheme.

**Multi-task Learning from Synthetic Data**. The multi-task setting refers to the ability to learn multiple tasks simultaneously in which all model parameters are subject to a shared influence [12, 15, 22, 31, 32, 53, 63, 64, 87, 88, 108, 115]. Many recent works employ multi-task learning in CNNs [66, 92] and Transformers [11, 43, 107] to exploit the potential advantages of fast training, stronger results, and fewer parameters. MTFormer [107] is a transformer-based architecture, where multiple tasks share the same transformer encoder and decoder but has multiple modules layered on top of that for each task. MulT [11] is a transformer-based encoder-decoder model with shared attention to learn task inter-dependencies, and UniT [43] jointly learns multiple tasks across different domains, from object detection to vision-and-language reasoning and natural language understanding. In contrast to these works, our work is a form of prompt-driven auxiliary task learning which uses synthetic scene-level annotations to train video transformers for improving action recognition on real video data.

**Scene Understanding Models**. Recently, scene understanding models that use scene-level annotations have been successfully applied to a wide range of computer vision applications: panoptic segmentation [20, 77], video relation understanding [61, 83, 89], vision and language [19, 57, 58, 90], relational reasoning [8, 9, 42, 45, 54, 78, 106, 109], human-object interactions [30, 51, 105], action recognition [3, 33, 34, 40, 41, 46, 72, 86, 97, 102, 111], and even image & video generation [7, 39, 49]. In our work, we demonstrate how video transformers can utilize shared representations from a variety of multiple different synthetic tasks to perform video downstream tasks.

**Video Transformers**. Vision Transformers [25, 91] recently proposed a new approach to image recognition by discarding the convolutional inductive bias entirely and instead employing self-attention operations. With the advent of ViT, and the fact that attention-based architectures are a natural choice for modeling long-range contextual relationships in video, a number of video transformer models, including TimeSformer [10], ViViT [2], Mformer (MF) [73], ORViT [40], MViT [27], MViTv2 [60] and Video Swin [67], form the latest era in action recognition. We choose to work with MViTv2, although our method can be used on top of any of these. Our work exploits the seamless ability of the transformer architecture to process multiple domains and to integrate the underlying structure among tasks for several downstream video-related tasks.

## 3. The PViT Model

Our PViT approach utilizes synthetic data of various domains with multiple types of scene annotations to enhance video understanding models. We consider the setting in which the main goal is to learn downstream video-understanding tasks, such as action recognition or action detection, while leveraging multiple synthetic scene-annotated datasets. The key idea of our work is that multi-task prompt learning can be used to incorporate synthetic scene tasks into the video model. This is achieved by adding a set of *task prompts*, each corresponding to a different task, and letting each prompt predict task-related annotations. Importantly, all prompts are part of the computation for any video, regardless of the underlying task, and thus enables sharing information among auxiliary tasks.

We begin by describing the video transformer architecture and the training setup (Section 3.1). We then introduce our Multi-task Prompts (Section 3.2) and the Training losses (Section 3.3). Our method is illustrated in Figure 2.

### 3.1. Preliminaries

**Video Transformer Architecture**. A typical Video Transformer model takes as input a video $X \in \mathbb{R}^{T \times 3 \times H \times W}$, extracts $N$ non-overlapping per-frame patches $x_i \in \mathbb{R}^{3 \times h \times w}$ and projects them into a lower-dimension $d$ (e.g., see [25]). Denote the transformer patches by $Ex_i$, which we refer to as "patch tokens". Then, spatio-temporal position embeddings $PE \in \mathbb{R}^{N \times d}$ are added for providing location and time location information, resulting in a new embedding: $z_i = Ex_i + PE_i$. This forms the sequence of input tokens to the video transformer:

$$z = [z_{CLS}, z_1, z_2, \cdots, z_N] \quad (1)$$

where $z_{CLS}$ is a CLS token used for the downstream task. Next, a transformer is comprised of a stack the Multi-headed Attention (MHSA) blocks, which apply the self-attention operation over all patch tokens $z$ (including the CLS token $z_{CLS}$) followed by a Feed-Forward Network (FFN), a layer normalization (LayerNorm [6]) step and a non-linear operation with residual connections [38].

**Training Setup for Various Domains**. In our approach, we aim to process batches of videos from various domains for $n$ different tasks. A key desideratum in this context is to be able to input both videos of synthetic scene data across various domains for multiple tasks, as well as videos from the real domain into the same model. In contrast to standard training, where each sample contains a full set of annotations (e.g., depth, normal, etc.), in our case, only partial annotations are included. This is explained in greater detail in Section 3.3.

### 3.2. Multi-task Prompts

As mentioned earlier, our key observation is that multi-task prompt learning can be used to incorporate synthetic scene tasks into the video model. Towards this end, we add a set of "task
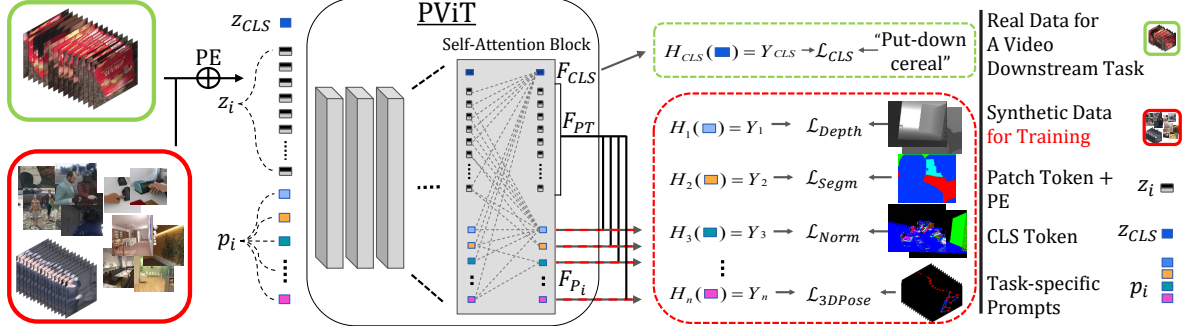
Figure 2. **PViT architecture**. We extend a transformer with a set of "task prompts", $p_i$, that are designed to capture information regarding each task, as well as capture the inter-task structure. The prompts are supervised by synthetic scene auxiliary tasks (depth, segmentation, normal, and 3D pose) available only during training, in order to enhance performance on a video task (predicting "put-down cereal"). Each task prompt in the attention block interacts with the patch tokens and CLS token, as well as other task prompts within the block.

prompts" designed to capture information regarding each task, as well as capture the inter-task structure. Specifically, we define a fixed number of $n$ learned vectors $p_1, p_2, \cdots, p_n \in \mathbb{R}^{1 \times d}$ for tasks $T_1, \cdots, T_n$. We refer to these vectors as the learned task prompts.

Let $P = \{p_1, p_2, \cdots, p_n\}$ be the set of task prompts. These prompts are concatenated to the patch tokens to obtain the following set of inputs to the transformer:

$$z = [z_{CLS}, z_1, z_2, ..., z_N, p_1, p_2, ..., p_n] \tag{2}$$

The transformer processes the input $z$, resulting in a new representation for each token $z$ (i.e., the CLS token, the patch tokens, and the task-prompts). We denote $F_{CLS}(z)$ as the representation of the CLS token, and let $F_{P_i}(z)$ denote the representation of the $i^{th}$ task-prompt. We also use $F_{PT}(z)$ as the final representation of all the patch tokens.

Next, these final output tokens are used for predicting labels. For the action recognition task, we simply predict using $F_{CLS}(z)$ and a prediction head $\hat{Y}_{CLS} = H_{CLS}(F_{CLS}(z))$. For the synthetic tasks, the task $i$ has a prediction head $\hat{Y}_i = H_i(F_{P_i}(z), F_{PT}(z))$ that is used for predicting labels corresponding to this task. It uses the patch tokens only for cases where a dense prediction is required (e.g., segmentation maps, normal and depth estimation). The task heads $H_i$ for localization tasks (e.g., boxes and 3D poses), are a simple FC layer, while for dense prediction tasks, we upsample patch token outputs from several layers and concatenate them with the corresponding task token to predict the task output map. Figure 3 also visualizes the "task prompts" learned by our model. For more info about the prediction heads see Section C in Supplementary.

### 3.3. Training and Inference

Our training data consists of labeled examples from $n$ synthetic tasks, as well as the downstream task of action recognition. As mentioned above, we have $n+1$ predictions heads corresponding to those. During training, for each training video we add a loss corresponding to the labels provided for that video. For example, if the synthetic video $X$ contains

labels for task 2 (e.g., depth) and task 5 (e.g., normal), we take the output of prediction heads $F_2$ and $F_5$ and compare them to the ground-truth labels for these two tasks. We formally describe the task-specific losses below. We use $\hat{Y}$ to refer to predicted labels, and $Y$ for ground-truth labels.

**Losses.** For Depth Estimation, we first downsample the ground-truth depth map $Y_{Depth}$ to a fixed scale of $\tilde{h} \times \tilde{w}$ map. Next, we predict a fixed scale map $\hat{Y}_{depth}$, and clip large values to focus on relatively closer objects. Finally, we use the MSE loss for computing the per-pixel depth error:

$$\mathcal{L}_{Depth} = \frac{1}{\tilde{h} \times \tilde{w}} \cdot \text{MSE}\left(\hat{Y}_{Depth}, Y_{Depth}\right) \tag{3}$$

For Normal Estimation, we predict the normal map $\hat{Y}_{Normal} \in \mathbb{R}^{h \times w \times 3}$ for every axis in world coordinates. We again down-sample the ground truth map $Y_{Normal}$ and compute the MSE loss:

$$\mathcal{L}_{Normal} = \frac{1}{\tilde{h} \times \tilde{w}} \cdot \text{MSE}\left(\hat{Y}_{Normal}, Y_{Normal}\right) \tag{4}$$

For Semantic Segmentation, we use per-pixel multi-label classification to compute a map for different semantic instances in the scene. We downsample the ground-truth map $Y_{Segm}$ and compute pixel-level cross-entropy loss followed by a Softmax function:

$$\mathcal{L}_{Segm} = \frac{1}{\tilde{h} \times \tilde{w}} \cdot \text{CE}\left(\hat{Y}_{Segm}, Y_{Segm}\right) \tag{5}$$

For 3D Pose estimation, we predict a tensor $\hat{Y}_{Pose3d} \in \mathbb{R}^{1 \times 75}$ corresponding to a $25 \times 3$ of 3D joins in KinectV2 format [16]. Each training sample consists of a single individual. We define the loss for 3D Pose Estimation to be:

$$\mathcal{L}_{3DPose} = \frac{1}{75} \cdot \text{MSE}\left(\hat{Y}_{3DPose}, Y_{3DPose}\right) \tag{6}$$

For Bounding Box Prediction, we set a fixed number of $O$ objects per training sample and use the $L1$ loss function
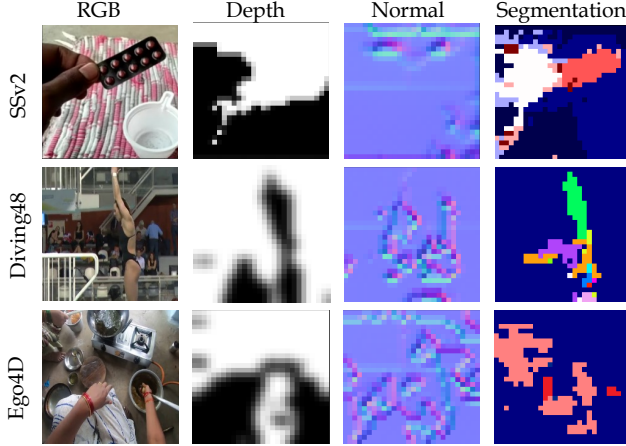
Figure 3. **"Task Prompts" Visualization**. Visualization of the output of the "task prompts" prediction heads on frames from the SSv2, Diving48, and Ego4D datasets. The model was trained with Something-Else as the action recognition dataset. Shown are prediction head outputs (i.e., $H_i$) for depth, normal, and semantic segmentation. It can be seen that the task prompts produce meaningful maps, despite not receiving such labels for real videos.

to compute boxes predictions $\hat{Y}_{Boxes} \in \mathbb{R}^{O \times 4}$ and the corresponding ground-truth coordinates $Y_{Boxes}$:

$$\mathcal{L}_{Boxes} = \text{L}_1\left(\hat{Y}_{Boxes}, Y_{Boxes}\right) + \text{GIoU}\left(\hat{Y}_{Boxes}, Y_{Boxes}\right) \quad (7)$$

where the GIoU is used as in [79].

Last, for the video downstream task (denoted as $DT$), on which we evaluate our model, we consider the standard cross-entropy loss between the predicted logits $\hat{Y}_{CLS}$ and the true video labels $Y_{CLS}$ as follows:

$$\mathcal{L}_{DT} = \text{CE}\left(\hat{Y}_{CLS}, Y_{CLS}\right) \quad (8)$$

The total loss is the sum of all of the losses described above. We note that only losses for which the samples have ground truth are added since the ground truth changes across instances, as our training does not use explicit correspondences between different input modalities. Each of the task terms in the loss is multiplied by a hyper-parameter ($\lambda$), and these were chosen such that all loss components have the same scale (see Supplementary). The total loss is the weighted combination of all terms:

$$\mathcal{L}_{Total} = \lambda_{DT}\mathcal{L}_{DT} + \lambda_{Depth}\mathcal{L}_{Depth} + \lambda_{Normal}\mathcal{L}_{Normal} \\ + \lambda_{Segm}\mathcal{L}_{Segm} + \lambda_{3DPose}\mathcal{L}_{Pose3d} + \lambda_{Boxes}\mathcal{L}_{Boxes} \quad (9)$$

For simplicity, we omit the temporal dimension when predicting the losses above per frame.

**Inference**. For inference, PViT receives input from the real videos without requiring any additional synthetic data.

Finally, our method can be applied on top of a variety of video transformers (MViT [27], TimeSformer [10], Mformer [73]). For our experiments, we use the MViTv2 [60] model because it performs well empirically.

## 4. Experiments and Results

We begin by describing the datasets (Section 4.1), implementation details (Section 4.2), and baselines (Section 4.3). Next, we evaluate our approach on several benchmarks and tasks. Specifically, we consider the following tasks: Compositional Action Recognition (Section 4.4), Object State Change Classification & Localization (Section 4.5), Action Recognition (Section 4.6), and Spatio-Temporal Action Detection (Section 4.7).

### 4.1. Datasets

We first describe the datasets used for the downstream video tasks, followed by the datasets used as auxiliary synthetic datasets including their annotations. We use the following video datasets: **(1) Something-Something v2 (SSv2)** [35] is a dataset containing 174 action categories of common human-object interactions. **(2) SomethingElse [69]** which exploits the compositional structure of SSv2, where a combination of a verb and a noun defines an action. We follow the official compositional split from [69], which assumes the set of noun-verb pairs available for training is disjoint from the set given at test time. **(3) Ego4D** [36] is a new large-scale dataset of more than 3,670 hours of video data, capturing the daily-life scenarios of more than 900 unique individuals from nine different countries around the world. **(4) Diving48** [59] contains 48 fine-grained categories of diving activities. **(5) Atomic Visual Actions (AVA)** [37] is a benchmark for human action detection, we report Mean Average Precision (mAP) on AVA-V2.2. For "auxiliary" synthetic datasets, we use **(1) SURREACT** [94], a novel synthetic data generation method based on real human motion from real datasets. The method renders 3D SMPL [68] sequences with randomized cloth textures, lighting, and body shapes from 3D skeleton joints extracted by Kinect V2 [52] from the two following datasets: **(i) NTU RGB+D** [85] is a large-scale multi-view video dataset of RGB-D human actions with 56,880 samples collected from 40 subjects, including depth maps and 3D skeleton joints. **(ii) UESTC RGB-D** [47] is also a multi-view action dataset that with 40 categories of aerobic exercise along with depth maps and 3D skeleton joints. **(2) HyperSim** [80] is a photorealistic synthetic dataset for holistic indoor scene understanding. This dataset contains 77,400 HD images of 461 indoor scenes as well as ground truth depth and normal values for each pixel. **(3) Procedural Human Action Videos (PHAV)** [23] is a human action video dataset which relies on procedural generation and other computer graphics techniques of modern game engines. There are 39,982 actions in 35 categories, annotated with optical flow, segmentation, and depth maps. **(4) KIST SynADL** [44] generated by the ElderSim engine, is a large-scale synthetic dataset of elders' activities. There are 462K RGB videos representing 55 action classes, along with 2D, 3D skeleton joints positions used as ground truth. **(5) EHOI** [55] consists of 20K synthetic image dataset of first-person view, annotated with segmentation masks, and hand-object interaction boxes of 19 categories.

| Model | Compositional | | Base | | Few-Shot | |
|---|---|---|---|---|---|---|
| | Top-1 | Top-5 | Top-1 | Top-5 | 5-Shot | 10-Shot |
| I3D [14] | 42.8 | 71.3 | 73.6 | 92.2 | 21.8 | 26.7 |
| SlowFast [28] | 45.2 | 73.4 | 76.1 | 93.4 | 22.4 | 29.2 |
| TimeSformer [10] | 44.2 | 76.8 | 79.5 | 95.6 | 24.6 | 33.8 |
| Mformer [73] | 60.2 | 85.8 | 82.8 | 96.2 | 28.9 | 33.8 |
| MViTv2 [60] | 63.3 | 87.5 | 83.7 | 96.8 | 32.7 | 40.2 |
| MViTv2 MT | 63.0 | 87.6 | 79.8 | 95.8 | 32.7 | 40.6 |
| MViTv2 VPT | 53.0 | 81.8 | 76.8 | 94.8 | 31.8 | 39.0 |
| **PViT (Ours)** | **65.5** | **89.0** | **85.0** | **97.4** | **34.3** | **41.3** |
| | (+2.2) | (+2.5) | (+1.3) | (+0.6) | (+1.6) | (+1.1) |

Table 1. **Compositional and Few-Shot Action Recognition** on the SomethingElse dataset.

## 4.2. Implementation Details

PViT is implemented in PyTorch, and the code will be released upon acceptance and is included in the supplementary. Our training recipes and code are based on the MViTv2-S, $16 \times 4$ model, and were taken from `https://github.com/facebookresearch/mvit`. We pretrain the PromptonomyViT model on the K400 [52] video dataset. Then, we fine-tune on the downstream video task (detailed in Section 4.1) with the synthetic datasets and the PromptonomyViT loss. In the training batch, there are 64 videos with the number of synthetic videos being at most $\times 3$ the number of real videos. For more implementation details, see Section C in Supplementary.

## 4.3. Baselines

In our experiments, we compare PViT to several models reported in previous work for the corresponding datasets. These include the following methods: BMN [62], *I3D* [14], *SlowFast* [28], as well as the state-of-the-art transformers – *SViT* [5], *TimeSformer* [10], *ViViT* [2], and *MViTv2* [60].

Additionally, we explore two alternative ViT-based baselines. First, we consider a model we call *MViTv2 multi-task (MViTv2 MT)*, and is perhaps the simplest application of ViT to our task. It augments the MViTv2 model with multiple prediction heads (one per synthetic task) operating on the CLS token, but *does not use additional task prompts*. The prediction heads have the same architecture as $H_i$ used in PViT. We also consider a model we refer to as *MViTv2 VPT*, which is an implementation of the VPT [48] approach for action recognition. This is a simple prompt-based approach utilizes the additional task prompts included in PViT but does not use additional synthetic data and *keeps the backbone frozen*. The advantage of MViTv2 VPT is training efficiency, as fewer parameters are used in training. Considering VPT trains only a few parameters, we assume that the parameters are insufficient to account for differences between pretraining (K400) and target datasets (AVA, SSv2). This is in marked difference to the case of images where VPT worked (their pretraining and target benchmarks are similar in distribution and tasks). Nevertheless, we still find it important to evaluate their method in our setting.

| Model | Temporal Localization Error | PNR Classification Top-1 |
|---|---|---|
| Bi-LSTM | 0.790 | 65.3 |
| BMN [62] | 0.780 | - |
| I3D ResNet-50 [14] | 0.739 | 68.7 |
| MViTv2 [60] | 0.702 | 71.6 |
| MViTv2 MT | 0.640 | 73.6 |
| MViTv2 VPT | 0.791 | 64.2 |
| **PViT (Ours)** | **0.637** (-0.065) | **74.8** (+3.2) |

Table 2. **PNR Temporal Localization** results on Ego4D.

## 4.4. Compositional & Few-Shot Action Recognition

In several video datasets, an action is defined as the combination of a verb and a noun. Hence, one of the challenges is to identify combinations of words that were not seen during training. This "compositional" setting was explored in the "SomethingElse" dataset [69], where verb-noun combinations in the test data do not occur in the training data. We also evaluate the few-shot compositional action recognition task in [69] (See Section C.2 in supplementary).

Table 1 reports the results for these two tasks. PViT outperforms MViTv2 baseline for both the *Compositional* and *Few-shot* tasks by 2.2% for the compositional task, and by 1.6%, 1.1% for 5 and 10-shot tasks. Furthermore, PViT outperforms MViTv2 MT, suggesting that the design of our task prompts approach is beneficial for learning from synthetic data. It can also be seen that MViT VPT performance is adversely affected, as suggested above, resulting in 53%.

## 4.5. Object State Change Tasks

Human activity relies heavily on hands and objects. Two tasks studying hand-object interaction have recently been introduced to the Ego4D [36] dataset. The first is temporal localization, which involves finding key frames that indicate a change in object state within a video clip. The second is the classification of object state changes, which indicates whether an object state has changed or not.

Table 2 reports results on the above two tasks in Ego4D. We observe that PViT performs better than MViTv2 by 3.2%/-0.065 on the classification/localization tasks. As in Section 4.4, it can be seen that PViT consistently outperforms MViTv2 MT and MViTv2 VPT baselines. Overall, these results indicate that PViT successfully leverages scene data, even for another downstream video task.

## 4.6. Action Recognition

Tables 3a and 3b report results for the standard action recognition task on the SSv2 and Diving48 datasets. It can be seen that in Diving48, our method improves over the MViTv2 baseline by 6.0%, outperforming the other methods. We hypothesize that this relatively high gain is due to (i) the large availability of synthetic pose annotations (which is likely to help in human actions in the Diving dataset; See Figure 4d). (ii) Since Diving is a small dataset, the introduction of additional synthetic

| (a) Something–Something V2 | | | | (b) Diving48 | | | | | (c) AVA-V2.2 | | |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Model | Pretrain | Top-1 | Top-5 | Model | Pretrain | Frames | Top-1 | | Model | Pretrain | mAP |
| SlowFast [28], R101 | K400 | 63.1 | 87.6 | SlowFast [28], R101 | K400 | 16 | 77.6 | | SlowFast [28], R50 | K400 | 22.7 |
| ViViT-L [2] | IN+K400 | 65.4 | 89.8 | TimeSformer [10] | IN | 16 | 74.9 | | SlowFast [28], R101 | K400 | 23.8 |
| MViTv1 [27] | K400 | 64.7 | 89.2 | MViTv2 [60] | K400 | 16 | 73.1 | | MViTv1 [27] | K400 | 25.5 |
| MViTv2 [60] | K400 | 68.2 | 91.4 | SViT [5] | K400 | 16 | 79.8 | | MViTv2 [60] | K400 | 26.8 |
| MViTv2 MT | K400 | 68.4 | 91.4 | MViTv2 MT | K400 | 16 | 82.2 | | MViTv2 MT | K400 | 26.3 |
| MViTv2 VPT | K400 | 61.5 | 87.5 | MViTv2 VPT | K400 | 16 | 69.8 | | MViTv2 VPT | K400 | 19.0 |
| **PViT (Ours)** | K400 | **69.4** (+1.2) | **91.6** (+0.2) | **PViT (Ours)** | K400 | 16 | **85.8** (+6.0) | | **PViT (Ours)** | K400 | **28.4** (+1.6) |

Table 3. **Results on SSv2, Diving48, and AVA-V2.2 datasets.** We report (a) Top-1 and top-5 accuracy on SSv2. (b) Top-1 on Diving48. (c) mAP metric on AVA. IN refers to ImageNet-21K. For additional comparisons, see Section A.1 in supplementary.

supervision results in a larger effect. Finally, PViT achieves a 1.2%, improvement in SSv2, indicating that PViT can improve on large datasets (180K videos). Last, PViT consistently outperforms MViTv2 MT and MViTv2 VPT baselines, as above.

## 4.7. Spatio-temporal Action Detection

Gu et al. [37] describes the action detection task on AVA as a two-stage prediction procedure. As a first step, boxes are detected using an off-the-shelf person detector, followed by a prediction of the action of each detected box. For fair comparisons, the person boxes are kept identical across approaches, and the final result is measured by the Mean Average Precision (MAP) metric.

Table 3c reports results for spatio-temporal action detection on the AVA dataset. We observe that PViT improves the MViTv2 baseline by 1.6%, thereby demonstrating the ability to leverage "task prompts" to detect and localize human actions. In addition, PViT consistently outperforms MViTv2 MT and MViTv2 VPT baselines, as above.

## 4.8. Ablations

We perform a comprehensive ablation study on the "SomethingElse" [69] dataset to measure the contribution of the different PViT components (See Table 4). For more ablations, see Section A in supplementary.

**The Role of Prompts and Tuning**. PViT contains two main concepts: (i) the addition of multiple task-specific prompts dedicated to unique tasks. (ii) training these prompt representations to predict task-related labels from synthetic data. We present results for different combinations of these two factors in Table 4a. First, to demonstrate the importance of having multiple prompts, one per task, we suggest the *MViTv2 one-prompt (OP)* variant. This variant is similar to PViT but uses *a single prompt* instead of $n$ prompts for $n$ auxiliary tasks. Since the number of prompts decreases, we compensate by increasing the dimension size. As shown in Table 4a, PViT outperforms the OP variant, suggesting that multiple prompts are important for integrating information across tasks.

Next, we consider the MViTv2 *neutral-prompts (NP)* variant, which is simply MViTv2 with additional prompts but without

additional synthetic supervision (similar to the MViTv2 VPT, but with an unfrozen backbone). The purpose of this variant is to examine whether the model performance is due to the increased model capacity. This result (63.4) is similar to the baseline without synthetic data (MViTv2, 63.3), suggesting that the gain of PViT is due to the use of synthetic data. Last, the *PViT VPT* variant is a simply PViT with a frozen backbone. MViTv2 VPT differs from this variant since here, synthetic data is used for training. The result (53.9) emphasizes the importance of fine-tuning the backbone even when using synthetic data.

**Model Capacity and Efficiency Analysis**. To determine whether the performance improvement is a result of increasing parameter size, Table 4a compares the number of parameters, FLOPS, and inference runtime between the methods. The main difference between the models is due to the additional task prompts and the task heads since the latter contains the most overhead (only during training). In our setting, task prompts only add 20K parameters, while the task heads add 6.8M parameters. However, during test time, the heads are not used, and thus the parameter sizes are almost equal to the baseline (i.e., 38.2M), resulting in similar inference runtime and FLOPS as the baseline.

**Effect of Synthetic Data Size**. Here, we examine the impact of the synthetic data portion on performance. In Figure 4b, we plot the performance of PViT as a function of the synthetic data portion when the largest value is obtained using all synthetic data. The positive slope suggests that adding synthetic data consistently improves results, which is an advantage since synthetic data is abundant. We note that the synthetic data we used is the size of the real data.

**Contribution from Auxiliary Tasks**. To investigate the impact of each auxiliary task on performance, we examined in Table 4c how the auxiliary tasks contribute to performance individually, as well as the most effective combinations of auxiliary tasks. As can be seen, we find that performing PViT on auxiliary tasks individually does improve performance (see also *Dataset Task Agreement* below). However, using all tasks (last line) improves more than any individual task, and is also close to the optimal combination. This reinforces our strategy of simply training on all tasks.

**Dataset-Task Agreement**. We next aim to explore how dif-

(a) **The Role of Prompts and Tuning**

| Model | Top-1 | Top-5 | Synthetic Data | Train/Test Params ($\times10^6$) | FLOPS ($\times10^6$) | Runtime (ms) |
|---|---|---|---|---|---|---|
| MViTv2 [60] | 63.3 | 87.5 | ✗ | 38.2/38.2 | 70.6 | 132.2 |
| MViTv2 MT | 62.7 | 87.6 | ✓ | 45.0/38.2 | 89.3 | 131.4 |
| MViTv2 OP | 63.5 | 88.0 | ✓ | 45.0/38.2 | 89.5 | 137.7 |
| MViTv2 NP | 63.4 | 87.8 | ✗ | 38.2/38.2 | 79.9 | 154.5 |
| MViTv2 VPT | 53.0 | 81.8 | ✗ | 0.13/38.2 | 82.3 | 154.5 |
| PViT VPT | 53.9 | 82.4 | ✓ | 7.2/38.2 | 93.9 | 143.7 |
| PViT | **65.5** | **89.0** | ✓ | 45.0/38.2 | 93.9 | 142.8 |

(b) **Effect of Synthetic Data Size**



(c) **Auxiliary Tasks Contribution**

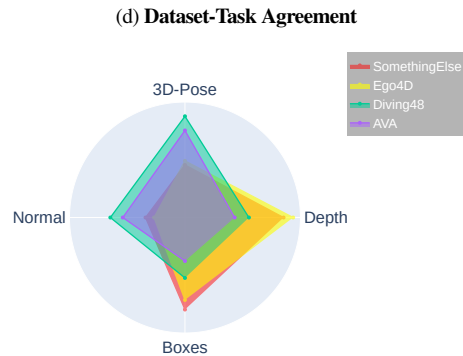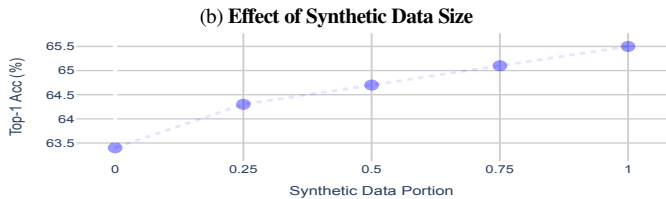| Datasets | Depth | Segm. | Normal | 3D Poses | 2D Boxes | Top-1 | Top-5 |
|---|---|---|---|---|---|---|---|
| - | ✗ | ✗ | ✗ | ✗ | ✗ | 63.3 | 87.5 |
| PHAV+HS+SURR | ✓ | ✗ | ✗ | ✗ | ✗ | 64.8 | 88.7 |
| SUR+EHOI | ✗ | ✓ | ✗ | ✗ | ✗ | 65.0 | 88.7 |
| HS | ✗ | ✗ | ✓ | ✗ | ✗ | 63.9 | 88.2 |
| SUR+ES | ✗ | ✗ | ✗ | ✓ | ✗ | 64.1 | 88.4 |
| EHOI | ✗ | ✗ | ✗ | ✗ | ✓ | 64.7 | 88.6 |
| best combination | ✓ | ✓ | ✗ | ✗ | ✓ | **65.5** | **89.0** |
| All | ✓ | ✓ | ✓ | ✓ | ✓ | 65.1 | 88.8 |

(d) **Dataset-Task Agreement**



Table 4. **Ablations.** We show (a) The Role of Prompts and Tuning. (b) Effect of Synthetic Data Size. (c) Contribution of Auxiliary Tasks. (d) Dataset-Task Agreement. A polygon represents a real video dataset, and the closer a vertex is to the circle border, the greater the gain from applying that synthetic task. The gains are scaled for comparison.

ferent synthetic tasks help real datasets. Figure 4d illustrates the gain for real datasets when trained on individual auxiliary tasks[3]. It can be seen that the datasets are roughly clustered into two sets: (i) SomthingElse and Ego4D, which benefit more from Depth and Boxes. These datasets indeed contain hands interacting with *objects* within close range of the camera and therefore having clearly expressed *depth*. (ii) AVA and Diving48, which benefit more from Normals and Poses. These datasets generally consist of zoomed-out frames with mostly *full human bodies* in scenes containing *solid surfaces* (for example, pools, walls, etc.). For more details, see Section A.2 in the supplementary.

**Domain Gap Between Synthetic and Real Data**. In this work, we show that training PViT on synthetic data leads to improved performance on real data. However, as synthetic and real data come from different domains, it is not apriori clear why the former should aid the latter. We hypothesize that our synthetic tasks are mostly low-level (e.g., depth/normal maps, segm. masks), and for these, there may be a smaller gap between synthetic and real domains (See [17, 96]). To illustrate this, we use our learned task heads to predict labels on real data. Recall that these heads are learned only on synthetic data. Figure 3 shows results for this prediction, and it can be seen that the synthetic prompts predict well also on real data. This demonstrates that the synthetic tasks learned are also usable on real data.

## 5. Discussion and Limitations

Semantic understanding of videos is a key element of human visual perception, but its modeling is still challenging for machine vision. In this work, we propose a new method for exploiting various types of scene-level data to improve the performance of video understanding tasks. We present a multi-task prompt learning approach for video transformers, where a shared transformer backbone is enhanced with task-specific prompts. The use of task-specific prompts allows the model to benefit from task-related information, among different domains. We demonstrate improved performance on several video understanding benchmarks, highlighting the effectiveness of the proposed approach. However, the multi-task prompt learning method is not necessarily limited to synthetic scene data, and thus we leave to future research the challenge of extending the work to train the method on real data as well as improving other downstream tasks in addition to video understanding.

---

[3]The plot excludes segm. since it contributes equally to all datasets.

# References

[1] Iro Armeni, Sasha Sax, Amir Roshan Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *ArXiv*, abs/1702.01105, 2017. 2

[2] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer, 2021. 2, 3, 6, 7

[3] Anurag Arnab, Chen Sun, and Cordelia Schmid. Unified graph structured models for video understanding. In *ICCV*, 2021. 2, 3

[4] Akari Asai, Mohammadreza Salehi, Matthew E. Peters, and Hannaneh Hajishirzi. Attentional mixtures of soft prompt tuning for parameter-efficient multi-task knowledge sharing. *ArXiv*, abs/2205.11961, 2022. 2

[5] Elad Ben Avraham, Roei Herzig, Karttikeya Mangalam, Amir Bar, Anna Rohrbach, Leonid Karlinsky, Trevor Darrell, and Amir Globerson. Bringing image scene structure to video via frame-clip consistency of object tokens. In *Thirty-Sixth Conference on Neural Information Processing Systems*, 2022. 2, 6, 7

[6] Jimmy Ba, Jamie Ryan Kiros, and Geoffrey E. Hinton. Layer normalization. *ArXiv*, abs/1607.06450, 2016. 3

[7] Amir Bar, Roei Herzig, Xiaolong Wang, Anna Rohrbach, Gal Chechik, Trevor Darrell, and A. Globerson. Compositional video synthesis with action graphs. In *ICML*, 2021. 3

[8] Fabien Baradel, Natalia Neverova, Christian Wolf, Julien Mille, and Greg Mori. Object level visual reasoning in videos. In *ECCV*, pages 105–121, 2018. 3

[9] Peter W Battaglia, Jessica B Hamrick, Victor Bapst, Alvaro Sanchez-Gonzalez, Vinicius Zambaldi, Mateusz Malinowski, Andrea Tacchetti, David Raposo, Adam Santoro, Ryan Faulkner, et al. Relational inductive biases, deep learning, and graph networks. *arXiv preprint arXiv:1806.01261*, 2018. 3

[10] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? In *Proceedings of the International Conference on Machine Learning (ICML)*, July 2021. 3, 5, 6, 7

[11] Deblina Bhattacharjee, Tong Zhang, Sabine Süsstrunk, and Mathieu Salzmann. Muit: An end-to-end multitask learning transformer. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12021–12031, 2022. 3

[12] David Brüggemann, Menelaos Kanakis, Stamatios Georgoulis, and Luc Van Gool. Automated search for resource-efficient branched multi-task networks. *ArXiv*, abs/2008.10292, 2020. 3

[13] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision – ECCV 2020*, pages 213–229, 2020. 2

[14] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6299–6308, 2017. 6

[15] Rich Caruana. Multitask learning. In *Encyclopedia of Machine Learning and Data Mining*, 1998. 3

[16] L. Caruso, R. Russo, and S. Savino. Microsoft kinect v2 vision system in a manufacturing application. *Robotics and Computer-Integrated Manufacturing*, 48:174–181, 2017. 4

[17] Yuhua Chen, Wen Li, Xiaoran Chen, and Luc Van Gool. Learning semantic segmentation from synthetic data: A geometrically guided input-output adaptation approach. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1841–1850, 2019. 8

[18] Yujin Chen, Matthias Nießner, and Angela Dai. 4dcontrast: Contrastive learning with dynamic correspondences for 3d scene understanding. *ArXiv*, abs/2112.02990, 2021. 2

[19] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *ECCV*, 2020. 3

[20] Bowen Cheng, Maxwell D. Collins, Yukun Zhu, Ting Liu, Thomas S. Huang, Hartwig Adam, and Liang-Chieh Chen. Panoptic-deeplab: A simple, strong, and fast baseline for bottom-up panoptic segmentation. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12472–12482, 2020. 3

[21] Vasileios Choutas, Philippe Weinzaepfel, Jérôme Revaud, and Cordelia Schmid. Potion: Pose motion representation for action recognition. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7024–7033, 2018. 2

[22] Michael Crawshaw. Multi-task learning with deep neural networks: A survey. *ArXiv*, abs/2009.09796, 2020. 3

[23] CR De Souza, A Gaidon, Y Cabon, and AM Lopez Pena. Procedural generation of videos to train deep action recognition networks. In *CVPR*, 2017. 5

[24] César Roberto de Souza, Adrien Gaidon, Yohann Cabon, and Antonio Manuel López Peña. Procedural generation of videos to train deep action recognition networks. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2594–2604, 2017. 3

[25] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021. 2, 3

[26] Yong Du, Wei Wang, and Liang Wang. Hierarchical recurrent neural network for skeleton based action recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1110–1118, 2015. 2

[27] Haoqi Fan, Bo Xiong, Karttikeya Mangalam, Yanghao Li, Zhicheng Yan, Jitendra Malik, and Christoph Feichtenhofer. Multiscale vision transformers. In *ICCV*, 2021. 2, 3, 5, 7

[28] C. Feichtenhofer, H. Fan, J. Malik, and K. He. Slowfast networks for video recognition. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6201–6210, 2019. 6, 7

[29] Chuang Gan, Jeremy Schwartz, Seth Alter, Martin Schrimpf, James Traer, Julian De Freitas, Jonas Kubilius, Abhishek Bhandwaldar, Nick Haber, Megumi Sano, Kuno Kim, Elias Wang, Damian Mrowca, Michael Lingelbach, Aidan Curtis, Kevin T. Feigelis, Daniel Bear, Dan Gutfreund, David Cox, James J. DiCarlo, Josh H. McDermott, Joshua B. Tenenbaum, and Daniel L. K. Yamins. Threedworld: A platform for interactive multimodal physical simulation. *ArXiv*, abs/2007.04954, 2021. 2, 3

[30] Chen Gao, Jiarui Xu, Yuliang Zou, and Jia-Bin Huang. Drg: Dual relation graph for human-object interaction detection. *ArXiv*, abs/2008.11714, 2020. 3

[31] Yuan Gao, Haoping Bai, Zequn Jie, Jiayi Ma, Kui Jia, and Wei Liu. Mtl-nas: Task-agnostic neural architecture search towards general-purpose multi-task learning. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11540–11549, 2020. 3

[32] Yuan Gao, Qi She, Jiayi Ma, Mingbo Zhao, W. Liu, and Alan Loddon Yuille. Nddr-cnn: Layerwise feature fusing in multi-task cnns by neural discriminative dimensionality reduction. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3200–3209, 2019. 3

[33] Rohit Girdhar, Joao Carreira, Carl Doersch, and Andrew Zisserman. Video action transformer network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 244–253, 2019. 2, 3

[34] Rohit Girdhar, Deva Ramanan, Abhinav Kumar Gupta, Josef Sivic, and Bryan C. Russell. Actionvlad: Learning spatio-temporal aggregation for action classification. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3165–3174, 2017. 2, 3

[35] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The" something something" video database for learning and evaluating visual common sense. In *ICCV*, page 5, 2017. 5

[36] Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Girdhar, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan, Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray, Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Cartillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano Fragomeni, Qichen Fu, Christian Fuegen, Abrham Gebreselasie, Cristina Gonzalez, James Hillis, Xuhua Huang, Yifei Huang, Wenqi Jia, Weslie Khoo, Jachym Kolar, Satwik Kottur, Anurag Kumar, Federico Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Merey Ramazanova, Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo, Yuchen Wang, Xindi Wu, Takuma Yagi, Yunyi Zhu, Pablo Arbelaez, David Crandall, Dima Damen, Giovanni Maria Farinella, Bernard Ghanem, Vamsi Krishna Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva, Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba, Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. Ego4d: Around the World in 3,000 Hours of Egocentric Video. *CoRR*, abs/2110.07058, 2021. 2, 5, 6

[37] Chunhui Gu, Chen Sun, David A. Ross, Carl Vondrick, Caroline Pantofaru, Yeqing Li, Sudheendra Vijayanarasimhan, George Toderici, Susanna Ricco, Rahul Sukthankar, Cordelia Schmid, and Jitendra Malik. AVA: A video dataset of spatio-temporally localized atomic visual actions. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 6047–6056. IEEE Computer Society, 2018. 5, 7

[38] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR*, pages 770–778, 2016. 3

[39] Roei Herzig, Amir Bar, Huijuan Xu, Gal Chechik, Trevor Darrell, and Amir Globerson. Learning canonical representations for scene graph to image generation. In *European Conference on Computer Vision*, 2020. 3

[40] Roei Herzig, Elad Ben-Avraham, Karttikeya Mangalam, Amir Bar, Gal Chechik, Anna Rohrbach, Trevor Darrell, and Amir Globerson. Object-region video transformers. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. 2, 3

[41] Roei Herzig, Elad Levi, Huijuan Xu, Hang Gao, Eli Brosh, Xiaolong Wang, Amir Globerson, and Trevor Darrell. Spatio-temporal action graph networks. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 0–0, 2019. 2, 3

[42] Roei Herzig, Moshiko Raboh, Gal Chechik, Jonathan Berant, and Amir Globerson. Mapping images to scene graphs with permutation-invariant structured prediction. In *Advances in Neural Information Processing Systems (NIPS)*, 2018. 3

[43] Ronghang Hu and Amanpreet Singh. Unit: Multimodal multitask learning with a unified transformer. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 1419–1429, 2021. 3

[44] Hochul Hwang, Cheongjae Jang, Geonwoo Park, Junghyun Cho, and Ig-Jae Kim. Eldersim: A synthetic data generation platform for human action recognition in eldercare applications, 2020. 5

[45] Achiya Jerbi, Roei Herzig, Jonathan Berant, Gal Chechik, and Amir Globerson. Learning object detection from captions via textual scene attributes. *ArXiv*, abs/2009.14558, 2020. 3

[46] Jingwei Ji, Ranjay Krishna, Li Fei-Fei, and Juan Carlos Niebles. Action genome: Actions as composition of spatio-temporal scene graphs. *arXiv preprint arXiv:1912.06992*, 2019. 2, 3

[47] Yanli Ji, Feixiang Xu, Yang Yang, Fumin Shen, Heng Tao Shen, and Wei-Shi Zheng. A large-scale varying-view rgb-d action dataset for arbitrary-view human action recognition, 2019. 5

[48] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *European Conference on Computer Vision (ECCV)*, 2022. 2, 6

[49] Justin Johnson, Agrim Gupta, and Li Fei-Fei. Image generation from scene graphs. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1219–1228, 2018. 3

[50] Vicky S. Kalogeiton, Philippe Weinzaepfel, Vittorio Ferrari, and Cordelia Schmid. Action tubelet detector for spatio-temporal action localization. *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 4415–4423, 2017. 2

[51] Keizo Kato, Yin Li, and Abhinav Gupta. Compositional learning for human object interaction. In *ECCV*, 2018. 3

[52] 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. *arXiv preprint arXiv:1705.06950*, 2017. 2, 5, 6

[53] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 7482–7491, 2018. 3

[54] Ranjay Krishna, Ines Chami, Michael S. Bernstein, and Li Fei-Fei. Referring relationships. *ECCV*, 2018. 3

[55] Rosario Leonardi, Francesco Ragusa, Antonino Furnari, and Giovanni Maria Farinella. Egocentric human-object interaction detection exploiting synthetic data, 2022. 5

[56] Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic, Nov. 2021. Association for Computational Linguistics. 2

[57] Liunian Harold Li, Mark Yatskar, Da Yin, Cho-Jui Hsieh, and Kai-Wei Chang. Visualbert: A simple and performant baseline for vision and language. *ArXiv*, abs/1908.03557, 2019. 3

[58] Xiujun Li, Xi Yin, Chunyuan Li, Xiaowei Hu, Pengchuan Zhang, Lei Zhang, Lijuan Wang, Houdong Hu, Li Dong, Furu Wei, Yejin Choi, and Jianfeng Gao. Oscar: Object-semantics aligned pre-training for vision-language tasks. *ECCV 2020*, 2020. 3

[59] Yingwei Li, Yi Li, and Nuno Vasconcelos. Resound: Towards action recognition without representation bias. In *Proceedings of the European Conference on Computer Vision (ECCV)*, September 2018. 5

[60] Yanghao Li, Chao-Yuan Wu, Haoqi Fan, Karttikeya Mangalam, Bo Xiong, Jitendra Malik, and Christoph Feichtenhofer. Mvitv2: Improved multiscale vision transformers for classification and detection. In *CVPR*, 2022. 2, 3, 5, 6, 7, 8

[61] Junwei Liang, Lu Jiang, Juan Carlos Niebles, Alexander G. Hauptmann, and Li Fei-Fei. Peeking into the future: Predicting future person activities and locations in videos. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5718–5727, 2019. 3

[62] Tianwei Lin, Xiao Liu, Xin Li, Errui Ding, and Shilei Wen. Bmn: Boundary-matching network for temporal action proposal generation. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 3888–3897, 2019. 6

[63] Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. In *NeurIPS*, 2021. 3

[64] Liyang Liu, Yi Li, Zhanghui Kuang, Jing-Hao Xue, Yimin Chen, Wenming Yang, Qingmin Liao, and Wayne Zhang. Towards impartial multi-task learning. In *ICLR*, 2021. 3

[65] Mengyuan Liu, Hong Liu, and Chen Chen. Enhanced skeleton visualization for view invariant human action recognition. *Pattern Recognition*, 68:346–362, 2017. 2

[66] Shikun Liu, Edward Johns, and Andrew J. Davison. End-to-end multi-task learning with attention. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1871–1880, 2019. 3

[67] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. *arXiv preprint arXiv:2106.13230*, 2021. 3

[68] Matthew Loper, Naureen Mahmood, Javier Romero, Gerard Pons-Moll, and Michael J. Black. Smpl: a skinned multi-person linear model. *ACM Trans. Graph.*, 34:248:1–248:16, 2015. 5

[69] Joanna Materzynska, Tete Xiao, Roei Herzig, Huijuan Xu, Xiaolong Wang, and Trevor Darrell. Something-else: Compositional action recognition with spatial-temporal interaction networks. In *proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2020. 5, 6, 7

[70] Hiroaki Mikami, Kenji Fukumizu, Shogo Murai, Shuji Suzuki, Yuta Kikuchi, Taiji Suzuki, Shin ichi Maeda, and Kohei Hayashi. A scaling law for synthetic-to-real transfer: How much is your pre-training effective?, 2021. 3

[71] Samarth Mishra, Rameswar Panda, Cheng Perng Phoo, Chun-Fu Chen, Leonid Karlinsky, Kate Saenko, Venkatesh Saligrama, and Rogério Schmidt Feris. Task2sim: Towards effective pre-training and transfer from synthetic data. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9184–9194, 2022. 3

[72] Tushar Nagarajan, Yanghao Li, Christoph Feichtenhofer, and Kristen Grauman. Ego-topo: Environment affordances from egocentric video. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 160–169, 2020. 3

[73] Mandela Patrick, Dylan Campbell, Yuki M. Asano, Ishan Misra Florian Metze, Christoph Feichtenhofer, Andrea Vedaldi, and Joao F. Henriques. Keeping your eye on the ball: Trajectory attention in video transformers, 2021. 3, 5, 6

[74] Xingchao Peng, Baochen Sun, Karim Ali, and Kate Saenko. Learning deep object detectors from 3d models. *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 1278–1286, 2015. 3

[75] Aayush Prakash, Shaad Boochoon, Mark Brophy, David Acuna, Eric Cameracci, Gavriel State, Omer Shapira, and Stan Birchfield. Structured domain randomization: Bridging the reality gap by context-aware synthetic data. *2019 International Conference on Robotics and Automation (ICRA)*, pages 7249–7255, 2019. 3

[76] Albert Pumarola, Enric Corona, Gerard Pons-Moll, and Francesc Moreno-Noguer. D-nerf: Neural radiance fields for dynamic scenes. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10313–10322, 2021. 2

[77] Siyuan Qiao, Yukun Zhu, Hartwig Adam, Alan Loddon Yuille, and Liang-Chieh Chen. Vip-deeplab: Learning visual perception with depth-aware video panoptic segmentation. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3996–4007, 2021. 3

[78] Moshiko Raboh, Roei Herzig, Gal Chechik, Jonathan Berant, and Amir Globerson. Differentiable scene graphs. In *WACV*, 2020. 3

[79] Hamid Rezatofighi, Nathan Tsoi, JunYoung Gwak, Amir Sadeghian, Ian Reid, and Silvio Savarese. Generalized intersection over union. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019. 5

[80] Mike Roberts, Jason Ramapuram, Anurag Ranjan, Atulit Kumar, Miguel Angel Bautista, Nathan Paczan, Russ Webb, and Joshua M. Susskind. Hypersim: A photorealistic synthetic dataset for holistic indoor scene understanding. In *International Conference on Computer Vision (ICCV) 2021*, 2021. 5, 2

[81] Germán Ros, Laura Sellart, Joanna Materzynska, David Vázquez, and Antonio M. López. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3234–3243, 2016. 3

[82] Victor Sanh, Albert Webson, Colin Raffel, Stephen Bach, Lintang Sutawika, Zaid Alyafeai, Antoine Chaffin, Arnaud

Stiegler, Arun Raja, Manan Dey, M Saiful Bari, Canwen Xu, Urmish Thakker, Shanya Sharma Sharma, Eliza Szczechla, Taewoon Kim, Gunjan Chhablani, Nihal Nayak, Debajyoti Datta, Jonathan Chang, Mike Tian-Jian Jiang, Han Wang, Matteo Manica, Sheng Shen, Zheng Xin Yong, Harshit Pandey, Rachel Bawden, Thomas Wang, Trishala Neeraj, Jos Rozen, Abheesht Sharma, Andrea Santilli, Thibault Fevry, Jason Alan Fries, Ryan Teehan, Teven Le Scao, Stella Biderman, Leo Gao, Thomas Wolf, and Alexander M Rush. Multitask prompted training enables zero-shot task generalization. In *International Conference on Learning Representations*, 2022. 2

[83] Adam Santoro, David Raposo, David G Barrett, Mateusz Malinowski, Razvan Pascanu, Peter Battaglia, and Timothy Lillicrap. A simple neural network module for relational reasoning. In I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017. 3

[84] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi Parikh, and Dhruv Batra. Habitat: A platform for embodied ai research. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 9338–9346, 2019. 2, 3

[85] Amir Shahroudy, Jun Liu, Tian-Tsong Ng, and Gang Wang. Ntu rgb+d: A large scale dataset for 3d human activity analysis. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1010–1019, 2016. 5

[86] Chen Sun, Abhinav Shrivastava, Carl Vondrick, Kevin Murphy, Rahul Sukthankar, and Cordelia Schmid. Actor-centric relation network. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 318–334, 2018. 3

[87] Guolei Sun, Thomas Probst, Danda Pani Paudel, Nikola Popovic, Menelaos Kanakis, Jagruti R. Patel, Dengxin Dai, and Luc Van Gool. Task switching network for multi-task learning. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8271–8280, 2021. 3

[88] Ximeng Sun, Rameswar Panda, and Rogério Schmidt Feris. Adashare: Learning what to share for efficient deep multi-task learning. *ArXiv*, abs/1911.12423, 2020. 3

[89] Xu Sun, Tongwei Ren, Yuan Zi, and Gangshan Wu. Video visual relation detection via multi-modal feature fusion. *Proceedings of the 27th ACM International Conference on Multimedia*, 2019. 3

[90] Hao Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *EMNLP*, 2019. 3

[91] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In *International Conference on Machine Learning*, pages 10347–10357. PMLR, 2021. 3

[92] Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Mti-net: Multi-scale task interaction networks for multi-task learning. In *ECCV*, 2020. 3

[93] Gül Varol, Ivan Laptev, Cordelia Schmid, and Andrew Zisserman. Synthetic humans for action recognition from unseen viewpoints. *ArXiv*, abs/1912.04070, 2021. 3

[94] Gül Varol, Ivan Laptev, Cordelia Schmid, and Andrew Zisserman. Synthetic humans for action recognition from unseen viewpoints. In *IJCV*, 2021. 5

[95] Tu Vu, Brian Lester, Noah Constant, Rami Al-Rfou', and Daniel Cer. SPoT: Better frozen model adaptation through soft prompt transfer. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 5039–5059, Dublin, Ireland, May 2022. Association for Computational Linguistics. 2

[96] Tuan-Hung Vu, Himalaya Jain, Maxime Bucher, Mathieu Cord, and Patrick Pérez. Dada: Depth-aware domain adaptation in semantic segmentation. In *ICCV*, 2019. 8

[97] Xiaolong Wang and Abhinav Gupta. Videos as space-time region graphs. In *ECCV*, 2018. 2, 3

[98] Zhonghao Wang, Mo Yu, Yunchao Wei, Rogério Schmidt Feris, Jinjun Xiong, Wen mei W. Hwu, Thomas S. Huang, and Humphrey Shi. Differential treatment for stuff and things: A simple unsupervised domain adaptation method for semantic segmentation. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 12632–12641, 2020. 3

[99] Zifeng Wang, Zizhao Zhang, Sayna Ebrahimi, Ruoxi Sun, Han Zhang, Chen-Yu Lee, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, et al. Dualprompt: Complementary prompting for rehearsal-free continual learning. *European Conference on Computer Vision*, 2022. 2

[100] Zifeng Wang, Zizhao Zhang, Chen-Yu Lee, Han Zhang, Ruoxi Sun, Xiaoqi Ren, Guolong Su, Vincent Perot, Jennifer Dy, and Tomas Pfister. Learning to prompt for continual learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 139–149, 2022. 2

[101] Yo whan Kim, SouYoung Jin, Rameswar Panda, Hilde Kuehne, Leonid Karlinsky, Samarth Mishra, Venkatesh Saligrama, Kate Saenko, Aude Oliva, and Rogério Schmidt Feris. How transferable are video representations based on synthetic data?, 2022. 3

[102] Chao-Yuan Wu and Philipp Krähenbühl. Towards Long-Form Video Understanding. In *CVPR*, 2021. 2, 3

[103] 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. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13577–13587, 2022. 2

[104] F. Xia, Amir Roshan Zamir, Zhi-Yang He, Alexander Sax, Jitendra Malik, and Silvio Savarese. Gibson env: Real-world perception for embodied agents. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9068–9079, 2018. 3

[105] Bingjie Xu, Yongkang Wong, Junnan Li, Qi Zhao, and M. Kankanhalli. Learning to detect human-object interactions with knowledge. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2019–2028, 2019. 3

[106] Keyulu Xu, Jingling Li, Mozhi Zhang, Simon S. Du, Ken ichi Kawarabayashi, and Stefanie Jegelka. What can neural networks reason about? In *International Conference on Learning Representations*, 2020. 3

[107] Yangyang Xu, Xiangtai Li, Haobo Yuan, Yibo Yang, Jing Zhang, Yunhai Tong, Lefei Zhang, and Dacheng Tao. Multi-task learning with multi-query transformer for dense prediction. *ArXiv*, abs/2205.14354, 2022. 3

[108] Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. *ArXiv*, abs/2001.06782, 2020. 3

[109] Vinicius Zambaldi, David Raposo, Adam Santoro, Victor Bapst, Yujia Li, Igor Babuschkin, Karl Tuyls, David Reichert, Timothy Lillicrap, Edward Lockhart, et al. Relational deep reinforcement learning. *arXiv preprint arXiv:1806.01830*, 2018. 3

[110] Amir Roshan Zamir, Alexander Sax, Bokui (William) Shen, Leonidas J. Guibas, Jitendra Malik, and Silvio Savarese. Taskonomy: Disentangling task transfer learning. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3712–3722, 2018. 2

[111] Yubo Zhang, Pavel Tokmakov, Martial Hebert, and Cordelia Schmid. A structured model for action detection. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9967–9976, 2019. 3

[112] Jiaojiao Zhao, Yanyi Zhang, Xinyu Li, Hao Chen, Shuai Bing, Mingze Xu, Chunhui Liu, Kaustav Kundu, Yuanjun Xiong, Davide Modolo, Ivan Marsic, Cees G. M. Snoek, and Joseph Tighe. Tuber: Tubelet transformer for video action detection. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 13588–13597, 2022. 2

[113] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Conditional prompt learning for vision-language models. *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16795–16804, 2022. 2

[114] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language models. *Int. J. Comput. Vis.*, 130:2337–2348, 2022. 2

[115] Lingli Zhou, Zhen Cui, Chunyan Xu, Zhenyu Zhang, Chaoqun Wang, Tong Zhang, and Jian Yang. Pattern-structure diffusion for multi-task learning. *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4513–4522, 2020. 3