# Understanding Physical Dynamics with Counterfactual World Modeling

Rahul Venkatesh<sup>1\*</sup>, Honglin Chen<sup>1\*</sup>, Kevin Feigelis<sup>1\*</sup>, Daniel M. Bear<sup>1</sup>, Khaled Jedoui<sup>1</sup>, Klemen Kotar<sup>1</sup>, Felix Binder<sup>3</sup>, Wanhee Lee<sup>1</sup>, Sherry Liu<sup>1</sup>, Kevin A. Smith<sup>2</sup>, Judith E. Fan<sup>1</sup>, and Daniel L. K. Yamins<sup>1</sup>

 $^1$  Stanford  $^2$  MIT  $^3$  UC San Diego

**Abstract.** The ability to understand physical dynamics is critical for agents to act in the world. Here, we use Counterfactual World Modeling (CWM) to extract vision structures for dynamics understanding. CWM uses a temporally-factored masking policy for masked prediction of video data without annotations. This policy enables highly effective "counterfactual prompting" of the predictor, allowing a spectrum of visual structures to be extracted from a single pre-trained predictor without finetuning on annotated datasets. We demonstrate that these structures are useful for physical dynamics understanding, allowing CWM to achieve the state-of-the-art performance on the Physion benchmark. Code is available at https://neuroailab.github.io/cwm-physics/.

# 1 Introduction

Physical dynamics understanding involves predicting the effects of physical interactions with objects (e.g. predicting the trajectory of a thrown ball [26], or the direction of a falling stacked block tower [9]). This remains a critical challenge for autonomous agents such as robots and self-driving cars interacting with the world [22]. Existing computer vision algorithms significantly lag behind humans in physical dynamics understanding [11].

One class of existing methods relies on intermediate vision structures such as 2D object segmentations and 3D particle graphs [3,7,8,41,47,48,57,60,61,63,74]. These vision structures are highly useful for accurate dynamics prediction because they abstract away irrelevant details. However, these ground-truth structures are only available in simulated or manually annotated datasets. Scaling these approaches to unlabelled real-world video data remains challenging.

A contrasting class of approaches avoids the use of intermediate structures by learning to predict raw pixels of future video frames [2,4,23,24,32,33,55,72]. While these approaches are directly applicable to real-world videos, learning to predict future frame pixels poses many challenges due to the high-dimensionality of image pixels and the stochasticity of real-world physical dynamics. These unstructured methods substantially underperform approaches with direct access to ground-truth intermediates, especially 3D particles [11].

<sup>&</sup>lt;sup>\*</sup> Equal contribution

2 Venkatesh et al.



Fig. 1: Overview of the approach. Given an input video of a physical scenario, we extract feature representations and vision structures such as keypoints, optical flow, and segments. These structures are extracted from a single pre-trained CWM predictor without finetuning on annotated datasets. We use the extracted features and structures for dynamics understanding - detecting a past collision or predicting a future collision.

Beyond task-specific methods for physical dynamics prediction, a promising alternative is self-supervised learning of task-agnostic visual representations that transfer well to downstream vision tasks [12,21,35,53,65]. Methods such as DINO [12,53], masked autoencoder (MAE) [35], and VideoMAE [21,65,67] could potentially learn representations useful for dynamics understanding. An additional promise is the emergence of semantic segmentation structure in DINO [12,53], which could potentially improve dynamics understanding. However, these models are mostly used in a transfer learning or fine-tuning paradigm, which requires annotations. It remains unclear whether they can be prompted to extract meaningful structures without finetuning on annotated datasets.

Therefore, a key research question is designing methods that pre-train on realworld video data without annotations and support extraction of structures for dynamics understanding. In this work, we use a simple and powerful framework, called Counterfactual World Modeling (CWM) [10]. CWM allows extraction of structures useful for understanding dynamics. Figure 1 provides an overview of our approach. We summarize the contributions of CWM below:

(a) We show that using a *temporally-factored masking policy* during pretraining enables powerful prompting abilities. As in VideoMAE, we train a masked predictor on real-world video data. Unlike VideoMAE, in CWM, the predictor only takes in a few patches of the last frame and fully visible preceding frames as inputs, and predicts the remaining patches in the last frame. This temporally-factored masking policy encourages the predictor to concentrate information about transformations between frame pairs into the embeddings of a small number of patch tokens. This in turn enables the predictor to support effective prompting via simple interventions on those few key tokens, allowing the system to answer hypotheticals, such as what will the next frame look like if an object in an image is moved to the right.

(b) We demonstrate that CWM can be prompted to extract multiple vision structures useful for understanding dynamics. As a result of the masking policy, we can extract structures by feeding CWM different prompts. These structures are extracted from a single predictor without being supervised on annotated datasets. Utilizing the extracted structures, CWM achieves state-of-the-art performance on the challenging Physion benchmark [11].

CWM can be understood in the context of Pearl's Ladder of Causation [28], describing how counterfactual reasoning can be built up from statistical models. The first rung of the Ladder is *Association*, in which a model of the predictable statistical relations between observed events over time is constructed. In CWM, this role is played by the world model itself, the large pretrained predictor which absorbs correlations from observed video inputs. The second rung is Intervention, in which at key junctions of the statistical model, observational data are replaced by specific fixed choices ("interventions") intended to produce some desired outcome. In CWM, this role is played by patch-level prompting, whose utility is greatly enabled by the temporally-factored training of the underlying predictor. The third rung of the Ladder is *Counterfactual*, in which the results of interventions are compared to alternative futures to identify true causes of events. In CWM, the comparison between outcomes of counterfactual interventions (prompts) and alternative futures (observed ground truth or observed predictions) are used for structure extraction, which – since they better capture core underlying causes of physical events – end up being useful for improved physical prediction.

In what follows, we review the literature on related works, and describe the core concepts of the CWM framework. We then demonstrate that the extracted structures of CWM are highly useful for physical dynamics understanding. Lastly, we provide an analysis of the quality of the extracted structures and ablation studies of CWM.

# 2 Related Works

**Structured dynamics prediction** Researchers have made substantial progress in physical dynamics prediction using structured particle representations as inputs [3, 7, 8, 34, 41, 47, 48, 60, 61, 63]. These approaches simulate large systems of particle-based representations by constructing interaction graphs and propagating information between graph nodes. Besides particle graphs, alternative object structures such as entity locations [74] and keypoints [39] are useful for physical dynamics prediction. However, these methods rely on ground-truth object structures, which are only available in simulated or manually annotated datasets. The scalability of these methods on real-world unlabelled data remains limited.

Video prediction One class of approaches learns physics understanding by predicting the pixels of future video frames [2,4,23,24,32,33,55,72]. These methods are directly applicable to real-world videos without depending on ground-truth object structures, which are difficult to obtain in general scenarios. Recent video diffusion models [16,38,45,66] and transformer-based prediction models [31,76] have made progress towards more realistic pixel prediction of future video frames. However, learning to predict future frame pixels poses many challenges due to

the high-dimensionality of image pixels and the stochasticity of real-world physical dynamics. Existing state-of-the-art methods are prone to creating physically implausible motions in the predicted video frames [44].

Self-supervised visual representation learning Beyond task-specific methods for physical dynamics prediction, a promising alternative is self-supervised learning of task-agnostic visual representations from large-scale unlabeled image or video data. These methods learn to generate visual features that transfer well to downstream vision tasks. One school of works leverages different pretext tasks for pre-training [18,27,50,54,68,77]. Another class of works models image similarity and dissimilarity between augmented views of an image [12,15,36,51,53,71] and different clips of a video [17,58,79] via constrastive learning. The most recent family of masked visual modeling approaches learns effective visual representations via masking and reconstruction of visual tokens. iGPT [14] and ViT [19] pioneer this direction by training transformers on pixel or patch tokens and exploring masked prediction with patches. MAE [35] introduces autoencoding with an asymmetric encoder-decoder architecture and empirically shows that a high masking ratio is crucial for image tasks. VideoMAE [21,65] extends to the video domains and shows that an even higher masking ratio leads to strong performance for activity recognition tasks. V-JEPA [6] explores feature prediction as an objective for unsupervised learning from video and achieves state-of-the-art results on activity recognition task in the Something-Something V2 dataset [29]. However, the usefulness of these representations for physical dynamics understanding remains unexplored. Furthermore, these models are mostly used in a transfer learning or fine-tuning paradigm, which requires ground-truth annotations. It remains unclear whether they can be prompted to extract meaningful structures without additional training on annotated data.

# 3 Method

We discuss in generality the three concepts of CWM by climbing Pearl's Ladder of Causation [28]: (1) temporally-factored masked predictor for learning associations, (2) prompting as interventions and (3) structure extractions using counterfactuals. We will discuss the application of CWM to physical dynamics understanding in Section 4.

### 3.1 Temporally-factored masked predictor for learning associations

Masked predictor Following MAE [35] and VideoMAE [21,65], we train an encoder-decoder architecture to reconstruct masked observations of video frames. The input video frames are first divided into non-overlapping spatiotemporal square patches. Then a subset of the patches is masked, and only the remaining visible patches are passed as inputs into the encoder. Finally, the embedded tokens from the encoder and learnable mask tokens, with added positional embedding on all the tokens, are passed as inputs into a shallow decoder to reconstruct the masked patches. The predictor is trained with the mean squared error (MSE)



Fig. 2: Climbing the Ladder of Causation with the CWM framework: (a) Temporally-factored masked predictor for association learning. Given a frame pair input, the predictor takes in dense visible patches from the first frame and only a sparse subset of patches from the second frame as inputs, and learns to predict the masked patches. This policy encourages the model to concentrate scene dynamics into embeddings of a few patches. (b) Prompting as interventations. As a result of the temporally-factored masking, we can intervene by modifying one or a few visual patches in the prompt and steer the outcome of the predictor. (c) Structure extraction using counterfactuals. Multiple vision structures can be extracted by comparing the results of interventions to alternative futures (e.g. observed ground truth or observed predictions).

loss between the reconstructed patches and the original masked patches. The predictor learns the associations between spatiotemporal patches of observed video inputs.

**Temporally-factored masking** Unlike VideoMAE [21, 65], which randomly samples "tubes" or "cubes" of spatiotemporal patches to be masked, we use a temporally-factored masking policy for video inputs. Without loss of generality, we discuss the masking policy with a frame pair  $x_1, x_2 \in \mathbb{R}^{3 \times H \times W}$  as input. Given the input frame pair, we train a predictor  $\Psi$ :

$$\Psi(x_1^{\alpha}, x_2^{\beta}) = \tilde{x}_2 \tag{1}$$

which takes in first frame  $x_1$  and second frame  $x_2$  with masking ratio  $\alpha, \beta \in [0, 1]$ . The predictor  $\Psi$  predicts the masked patches of  $x_2$ , and minimizes the MSE loss between the reconstructed patches  $\tilde{x}_2$  and the masked patches of  $x_2$ . Figure 2a illustrates this masking policy.

Here, we set the masking ratio  $\alpha$  to 0 and  $\beta$  to 0.90, a highly asymmetric masking policy. As a result of this high masking ratio, the predictor  $\Psi$  learns to complete the second frame given only a few patches of it, along with the fully visible first frame; the only way it can do this is by inferring scene transformations from a few second-frame patches, then applying these transformations to the first-frame patches to complete the second frame [10]. This implies that the predictor learns to concentrate transformations between frame pairs into the embeddings of a few visible patches. Consequently, modifying the contents of a

few patches, which represent transformations, can exert meaningful control over the next-frame predictions.

### **3.2** Prompting as interventions

With a pre-trained predictor, at inference time we can replace empirical data observations with interventions intended to produce some desired outcome [28]. As a result of the temporally-factored masking policy, we can modify the original inputs at a few patch locations to generate alternative outcomes using the predictor. To formalize the procedure of intervention, we first define a prompt p as a set of video frames that is given as input to the predictor:

$$p = \{x_1, x_2 \mid x_1, x_2 \in \mathbb{R}^{3 \times H \times W}\}$$
(2)

where  $x_2$  has a small number of visible patches that specify scene transformations. An intervention  $\bar{p}$  is defined as an input to the predictor that has been modified from the initial prompt p. We use two basic types of interventions: (a) appearance prompts, which involve modifications to the first frame  $x_1$ , and (b) motion prompts, which involve modifications to the second frame  $x_2$ . Given a intervention  $\bar{p}$ , the associated prediction is the outcome of the predictor  $\Psi(\bar{p})$  [10]. Figure 3a shows the predictions for a series of motion prompts. These prompts use a single image,  $x_1$  and construct  $x_2$  by revealing only a few patches in the input image and translating them by a small offset.

### 3.3 Structure extraction using counterfactuals

The observation of the previous section shows how it is possible to generate counterfactual object motion by modifying the positions of a small number of patches. Next, we discuss how different structure extractions can be specified as counterfactuals [28] by comparing the outcomes of the interventions with alternative futures (e.g. observed ground-truth data or predictions).

**Keypoints** have been previously defined by manual category-specific annotations [20, 42, 75]. CWM provides a general category-agnostic definition of keypoints as patch locations in  $x_2$  that, when revealed to the predictor, yield the lowest error in the reconstruction  $\Psi(p)$  [10]. Let  $\mathcal{I}$  be a set of patch locations of an image. The set of keypoints is defined as:

$$K(x_1, x_2, n) = \underset{k \subset \mathcal{I}, |k|=n}{\operatorname{arg min}} \mathcal{L}(\Psi(p), \Psi(\bar{p}))$$
where  $\bar{p} = \{x_1, x_2^m \mid x_2^m \text{ is visible at } k\}$ 
(3)

Here, the intervention  $\bar{p}$  is the modification of the original input  $p = \{x_1, x_2\}$ , where the second frame  $x_2^m$  is masked everywhere except at keypoint locations. This construction defines a set of dynamical RGB keypoints on  $x_2$ . For large values of n, this is in general an intractable optimization problem. In practice, we thus first start with an empty set and add keypoints one at a time to greedily reduce the reconstruction error until n keypoints have been obtained. We show examples of extracted keypoints in the top panel of Figure 3b.



Fig. 3: Counterfactual predictions and structure extraction. (a) Counterfactual predictions. A small number of visual patches exert meaningful control of scene dynamics. Each panel shows a prompt consisting of the input image (left), a few patches copied from the input image (middle), and the resulting predictions (right). A red patch is copied into the same location as its source, simulating the appearance of holding an object fixed. A green patch is copied into a different location at an offset from the source location, simulating the appearance of an apparent object motion. (b) Structure extraction for keypoints, flows, and segments

**Optical flow** is the task of estimating per-pixel motion between video frames [64]. To estimate per-pixel motion, we introduce an appearance intervention that adds a small perturbation to the pixel in the first frame. We can estimate the pixel motion by localizing the perturbation response in  $\Psi(\bar{p})$  [10].

More specifically, given a prompt  $p = \{x_1, x_2^{\beta}\}$  and a pixel location (i, j), we construct an intervention  $\bar{p} = \{x_1+\delta_{ij}, x_2^{\beta}\}$ , which adds a small perturbation  $\delta_{ij}$  to the first frame at the pixel location. This creates a perturbed first frame by modifying its appearance at a pixel location. For this reason, we call this an appearance intervention. With a perturbed first frame, the predictor propagates the perturbation in the next frame, under the original scene transformations specified by  $x_2^{\beta}$ . The corresponding pixel location in the next frame can be localized by finding the peak of the perturbation response. The perturbation response in the next frame can be computed as the absolute difference between the counterfactual prediction  $\Psi(\bar{p})$  and the observed prediction  $\Psi(p)$ . Then, we locate the peak of the perturbation response by taking an argmax over the set of patch locations  $\mathcal{I}$ . The flow at pixel location (i, j) is then defined as the spatial displacement between (i, j) and the peak of perturbation response:

$$F_{i,j}(x_1, x_2) = \underset{\mathcal{I}}{\arg\max} |\Psi(\bar{p}) - \Psi(p)| - (i, j)$$

$$(4)$$

This algorithm is simple and often effective, as shown in the middle panel of Figure 3b, but it might fail in two ways. First, one of the revealed patches in  $x_2^{\beta}$  may cover the place where the perturbation at location (i, j) is expected to move. This can be remedied by running the above procedure for multiple random choices of  $x_2^{\beta}$  and taking their average perturbation responses [10].

A second potential failure mode is that the intervention  $\bar{p}$  might be out of distribution for  $\Psi$ , which could happen when the perturbation  $\delta_{ij}$  is too large [10]. On the other hand, if the perturbation is too small, it might not be detected and moved accurately. This can be naturally addressed by using infinitesimal perturbations. We normalize the magnitude of the perturbation response by the magnitude of the perturbation as the limit goes to zero. This is exactly the derivative of the  $\Psi$  at location (i, j):

$$\lim_{\delta_{ij}\to 0} \frac{|\Psi(\bar{p}) - \Psi(p)|}{|\delta_{ij}|} = \nabla_x \Psi \Big|_{(i,j)}$$
(5)

To simultaneously estimate optical flow at all locations of an input frame, we can compute the Jacobian of  $\Psi$ . This is a tensorial operation that can be computed once at all pixels using PyTorch autograd [1]. We describe more details about the procedures of extracting flow in the supplementary material.

**Segmentation** is defined as a collection of physical stuff that moves together under the application of everyday physical actions [13]. This is inspired by the notion of Spelke object in infant object recognition: infants tend to group scene elements that move together as a single object [62]. CWM extracts segmentation of objects by motion interventions, which simulate object motion at a pixel location, followed by grouping parts of the image that move together.

Given a single image x as input, we define an intervention  $\bar{p} = \{x, \bar{x}^m\}$ . These prompts produces the second frame  $\bar{x}^m$  by revealing only a few patches in the input image and translating those patches by a small offset. With a temporallyfactored masked predictor, moving a few patches in the prompt will cause the entire object to move in the resultant counterfactual predictions  $\Psi(\bar{p})$ . Segments can be extracted by thresholding the flows between the input image x and  $\Psi(\bar{p})$ :

$$S(x) = F(x, \Psi(\bar{p})) > 0 \tag{6}$$

Once a segment is extracted, we iterate the procedure above to refine the segment by revealing more patches within the segment region into  $\bar{x}^m$  and translating patches in the same direction. We set the number of iterations as 3. To automatically discover multiple objects in a single image, we iteratively extract segments at pixel locations that are not part of a discovered object. Once an object segment is discovered, we reveal patches that are not within the segment

region and repeat the procedure to discover the next object. We show examples of extracted segments in the bottom panel of Figure 3b. We discuss more details in the supplementary material.

#### 4 Experiments

Section 4.1 first investigates the usefulness of the extracted structures for downstream physical dynamics understanding tasks. Section 4.2 evaluates the quality of counterfactual motion predictions and extracted visual structures on realworld datasets. Section 4.3 discusses ablations studies on the CWM design.

#### 4.1**Physical Dynamics Understanding**

Physion benchmark consists of realistic simulations of diverse physical scenarios where objects are manipulated in a variety of configurations to test different types of physical reasoning such as stability, rolling motion, object linkage, etc. We use the latest version of Physion [11], referred to as Physion  $v1.5^1$ , which has improved rendering quality and more physically plausible simulations.

In the ideal scenario, we would evaluate CWM on a real-world physicsunderstanding benchmark, but such benchmarks are not available. Recent works have shown that simulated data can be highly valuable [43, 64, 78]. Physion is a challenging benchmark as it contains diverse physical phenomena, object dynamics and realistic 3D simulations. This makes it a preferable choice when compared to other benchmarks such as ShapeStacks [30] and IntPhys [59] which contain very limited object dynamics, or Phyre [5] which only operates in 2D environments. Existing video models still significantly lag behind human performance on the Physion benchmark [11]. Moreover, the CWM model is trained on real-world videos from Kinetics-400 dataset [40] and tested on Physion, and is thus a strong generalization test.

The benchmark consists of two tasks: (a) Object contact prediction (OCP), which tests the model's ability to *predict* whether two objects *will* contact at some point in the future given a context video, and (b) *Object contact detection* (OCD) which tests the model's ability to *detect* if two objects have already come into contact in the observed video. The video stimuli are generated in such a way that the model needs to have an understanding of the physical dynamics in order to answer the contact-related question correctly. Figure 4 shows example stimuli for the two tasks. For both tasks, the two objects of interest are rendered with red and yellow texture to cue the model.

**Evaluation protocol** We follow the three-step evaluation protocol of the Physion benchmark [11]. First, we extract features from the last layer of a frozen pre-trained encoder on a training set of 5,600 videos for OCP and OCD tasks, respectively. For image-based methods, features are extracted from 4 frames that are 150 ms apart. For video-based methods, the input frames are fed to the

<sup>&</sup>lt;sup>1</sup> https://physion-benchmark.github.io



Fig. 4: Physion v1.5 evaluation protocol. We evaluate on two physical dynamics understanding tasks - (a) Object contact prediction where the model is asked to predict contact events in the future and (b) Object contact detection where the model is asked to reason about contact events that occur in the observed video stimulus. The objects of interest for which we want to ask the contact question are rendered with red and yellow texture to cue the model.

model at the specific frame rate used during their training. Second, the extracted features are used to train a logistic regression model to predict the contact label for the given video stimulus. Lastly, the trained classifier is evaluated on a test set of 1,000 videos across different physical scenarios.

**Baseline methods** We compare CWM with five classes of baseline approaches: (a) video prediction models including MCVD [66], R3M [49], FitVid [4], and TECO [73], (b) self-supervised representation learning methods on images including DINO [12], DINOv2 [53], and MAE [35], (c) self-supervised representation learning methods on videos including VideoMAE [65], VideoMAEv2 [67] and the recent state-of-the-art method V-JEPA [6] (d) vision-language models like GPT4-V [52] and lastly (e) ground truth 3D particle-based dynamics prediction models such as SGNN [34].

**Results** In Table 1 we report results on the two Physion tasks for both CWM with ViT-B and ViT-L architectures and other baseline methods discussed above. We evaluate CWM with both features and extracted vision structures input to the linear classifier. We find that video prediction models (such as MCVD [66] and TECO [73]) perform poorly especially on the Physion tasks. Self-supervised image representation models on the other hand, are better but they saturate around 72% and 87% for OCP and OCD respectively with the ViT-B architecture. It is interesting to note that CWM outperforms methods such as DINOv2 ViT-g and MAE ViT-H which have 13 and 7 times more parameters. When scaled up to ViT-L, CWM achieves superior performance on OCP.

We find that CWM exhibits superior performance compared to both Video-MAE [65] and VideoMAEv2 [67]. To ensure a fair evaluation, we train a variant of VideoMAE, denoted as VideoMAE<sup>\*</sup>, that matches CWM in terms of the number of frames and patch size, and include comparable structure extractions from the model for linear probing. Our findings indicate that CWM performs better than VideoMAE<sup>\*</sup>. Furthermore, CWM surpasses the recently released V-JEPA [6], a state-of-the-art model for video representation, despite being trained on a considerably smaller dataset. Furthermore, we find on OCP, CWM achieves a performance that closely approaches that of ground truth 3D particle-based Table 1: State-of-the-art accuracy on Physion v1.5. We compare CWM to five classes of baseline methods across different architectures on the OCP and OCD tasks. We evaluate CWM with both features and extracted structures and find that it achieves state-of-the-art performance on these tasks. Original VideoMAE [65] uses 16 input frames and a patch size of 16. We trained VideoMAE\* with 3 input frames, a patch size of 8, and include extracted vision structures from the model for a strictly fair comparison with CWM.

method	training data	arch	param	$\operatorname{OCP} \uparrow$	$OCD\uparrow$
supervised ground t	ruth 3D particle-b	ased model			
SGNN	Physion v1.5	GNN	GNN 23 M		98.8
video prediction mo	odels				
MCVD [66] R3M [49] FitVid [4] TECO [73]	$egin{array}{c} \mathrm{K400+Ego4D}\\ \mathrm{K400+Ego4D}\\ \mathrm{K400+Ego4D}\\ \mathrm{K600} \end{array}$	UNet Res50 VAE vq-gan	251 M 38 M 303 M 160 M	$63.4 \\ 67.6 \\ 64.3 \\ 69.3$	80.8 78.1 59.5 80.9
self-supervised imag	ge representation	models			
DINO [12] DINOv2 [53] DINOv2 [53] DINOv2 [53] MAE [35] MAE [35] MAE [35] MAE [35] MAE [35]	IN-1K LVD-142M LVD-142M IN-142M IN-1K IN-1K IN-1K IN-1K IN-4.5M IN-4.5M	ViT-B ViT-B ViT-L ViT-g ViT-B ViT-L ViT-H ViT-B ViT-L	86M 86 M 304 M 1.1 B 86 M 304 M 632 M 86 M 304 M	$\begin{array}{c} 72.1 \\ 72.2 \\ 72.2 \\ 72.7 \\ 72.6 \\ 71.6 \\ 73.3 \\ 72.1 \\ 72.6 \end{array}$	85.4 87.1 85.5 84.6 81.6 82.3 80.8 81.7 81.9
self-supervised vide	o representation n	nodels			
VideoMAE [65] VideoMAE* VideoMAE [65] VideoMAE [65] VideoMAEv2 [67] V-JEPA [6]	K400 K400 K400 K400 U-Hybrid VideoMix2M	ViT-B ViT-B ViT-L ViT-H ViT-g ViT-L	86 M 86 M 304 M 632 M 1.1B 304M	$72.1 \\73.2 \\73.6 \\73.5 \\72.2 \\73.4$	85.7 86.2 86.1 87.5 85.0 87.0
vision-language mo	dels				
GPT4-V [52]	-	-	-	52.9	54.7
CWM CWM	K400 K400	ViT-B ViT-L	86 M 304 M	75.9 <b>76.1</b>	<b>89.1</b> 88.7

simulation models (i.e SGNN [34]) learned on Physion, despite being trained on Kinetics-400 [40] – a considerably different real world dataset.

We also evaluate GPT4-V [52] on Physion v1.5 tasks by providing it with a single composite image with a sequence of four video frames sampled at a gap of 150ms. The model is prompted with questions similar to those in Figure 4 (see supplementary for more details about the specific prompts used). We find GPT4-V scores nearly at chance on OCP and slightly above chance on OCD, which highlights a considerable limitation in the ability of large-scale vision-language models to understand physical scene dynamics.



Fig. 5: Qualitative comparison of counterfactual motion prediction and structure extraction on real-world datasets. We find that when we apply our extraction procedures described in Section. 3.3 on VideoMAE, the model fails to generate counterfactual motion and extracts less meaningful structures than CWM. Segments cannot be extracted from VideoMAE due to the failure of counterfactual predictions, and hence not shown in this comparison. This shows the importance of the temporally-factored masking policy during pre-training.

### 4.2 Analysis of the extracted structures

We analyze the quality of structures extracted by CWM. Although not all baseline methods can perform counterfactual predictions or structure extractions, we apply our procedures to the baseline methods for a fair comparison with CWM. We show CWM yields more meaningful predictions, keypoints, and flows than VideoMAE, enabled by the temporally-factored masking policy. We also show that the quality of segments extracted by CWM is close to the state-of-the-art method CutLER [70], which extracts segments from DINO [12].

**Counterfactual prediction** We compare CWM and VideoMAE on the quality of counterfactual predictions in Table 2a and Figure 5. We generate counterfactual motions using input images from the DAVIS dataset [56]. The quality of generation is measured by the Fréchet Inception Distance (FID) [37]. CWM significantly outperforms VideoMAE. For a strictly fair comparison, we train another VideoMAE model (referred to as VideoMAE\*) with the same number of frames and patch size as CWM. Although the model achieves a slightly lower FID relative to VideoMAE, the reconstructions are still quite blurry without accurate object motions. This illustrates the importance of temporally-factored masking in generating plausible counterfactual predictions.

**Keypoints** Existing keypoint datasets are generally created with manually specified templates for certain object categories [20, 42, 75]. Therefore, these datasets do not provide suitable quantitative evaluations of CWM keypoint, which are category-agnostic. Figure 5 shows that CWM can extract more meaningful dynamic keypoints as compared to VideoMAE.

Table 2: Quantitative comparison of counterfactual motions, flow and segment extraction on real-world datasets. In (a) we compare to VideoMAE on counterfactual motions and flow. For a strictly fair comparison, we also evaluate Video-MAE\* which we trained with the same patch size and number of frames as CWM. In (b) we compare CWM to CutLER [70], which extracts segmentations from DINO [12], and FreeSOLO [69] on the quality of segmentations.

(a) Counterfac	tual motion (C	M) and Flow	(b) Segmen	(b) Segments extraction			
Methods	CM (FID $\downarrow$ )	Flow (F1 $\downarrow)$	Methods	Segment (AP $\uparrow)$			
VideoMAE [65] VideoMAE* CWM	213.4 166.3 <b>25.4</b>	56.3 54.9 <b>46.8</b>	FreeSOLO [69] CutLER [70] CWM	4.3 <b>8.4</b> 8.2			

**Optical flow** We evaluate the quality of optical flows on the SPRING benchmark [46] using the F1 metric [25]. We find that CWM is better compared to both VideoMAE and VideoMAE\* (See Table. 2a). This is also supported by the qualitative results shown in Figure. 5. We include more qualitative comparisons and additional implementation details in the supplementary.

Segments We extract segments on images from COCO train2017 [42] using CWM. We follow the same procedures in CutLER [70] to learn a detector using the extracted segments as self-supervision. We train CutLER on COCO training images for a fair comparison. We compare CWM with FreeSOLO [69] and CutLER [70] in Table 2b and Figure 5. CWM outperforms FreeSOLO [69] significantly and attains similar performance to the current state-of-the-art approach CutLER [70]. Although Spelke objects are segment-like structures, the definition of Spelke objects is not exactly aligned with the definition of instance segmentations in the COCO datasets.

### 4.3 Ablation studies

We ablate the CWM design with the default backbone of ViT-B. Each ablated model is trained for 800 epochs on the Kinetics-400 dataset. Results of the ablation study are reported in Table. 3.

Vision structures We study the importance of each visual structure in understanding dynamics. Adding patch features at keypoint locations improves the OCP accuracy from 73.6% to 74.4%. Enriching these patch features with optical flow patches further improves the accuracy to 75.5%. Finally, including segments achieves a score of 75.9%.

**Training schedule** We find that a model trained with a longer training schedule of 1600 epochs achieves an OCP score of 75.9% – a relatively small improvement over an 800 epoch trained model (75.4%).

Masking Policy We study the importance of temporal factoring by training a model with a random tube masking strategy, which was originally proposed in VideoMAE [65]. The temporally-factored mask policy is essential for extraction of meaningful vision structures, improving the OCP accuracy improves from 73.2% with tube masking to 75.9% with temporally-factored masking.

**Table 3: CWM ablation studies.** The best setting is shown in the first row. We investigate the importance of different vision structures, masking policy, training epochs, masking ratio, context frames and patch size.

Ablations fa	temporal	Input to the classifier			mask	patch	context	training	Metrics		
	factoring	feat.	keyp	. flow	segm.	ratio	size	frames	epochs	$OCP\uparrow$	$OCD\uparrow$
Best setting	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	0.90	8	2	1600	75.9	89.1
Structures	$\checkmark$ $\checkmark$	√ √ √	√ √ ×	√ × ×	× × ×	$0.90 \\ 0.90 \\ 0.90$	8 8 8	$2 \\ 2 \\ 2$	$     1600 \\     1600 \\     1600 $	75.5 74.4 73.6	88.5 89.1 89.1
Training epoc	hs √	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	0.90	8	2	800	75.4	88.9
Masking polic	y <b>X</b>	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	0.90	8	2	800	73.2	86.2
Masking ratio	$\checkmark$ $\checkmark$	\$ \$ \$	$\checkmark$ $\checkmark$	$\checkmark$ $\checkmark$	√ √ √	$0.85 \\ 0.95 \\ 0.99$	8 8 8	2 2 2	800 800 800	75.0 74.6 72.5	88.9 88.3 86.6
Context frame	es √	√ √	√ √	√ √	√ √	$0.90 \\ 0.90$	8 8	$\frac{1}{4}$	800 800	$71.0 \\ 68.5$	85.2 79.9
Patch size	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	$\checkmark$	0.90	16	2	800	74.2	88.8

Mask ratio We observe that a high ratio on the last frame (90%) during model training achieves good performance on both the OCP and OCD tasks. This trend aligns with our aforementioned hypothesis that the dynamics between frame pairs at a short timescale has a low-dimensional causal structure, which can be concentrated into a small number of tokens.

**Context length** We compare the performance of CWM with different numbers of context frames. CWM with 2 context frames during pre-training performs better as compared to using 1 context frame. However, including 4 context frames degrades the performance.

**Patch size** Our analysis indicates that the patch size used for training the model can influence the performance; specifically, a patch size of 8 yields a superior OCP accuracy of 75.9%, compared to a patch size of 16, which results in a lower accuracy of 74.2%.

# 5 Conclusion

In this work, we show that a simple temporally-factored masking policy during pre-training enables powerful prompting abilities. As a result, we can use counterfactual prompts and their associated predictions to extract vision structures, which abstract away irrelevant details and thus end up being useful for improved dynamics understanding. As compared to random masking, temporally-factored masking policy allows more meaningful and useful structures to be extracted from the pre-trained predictor. CWM achieves state-of-the-art results on the challenging Physion benchmark as compared to previous self-supervised methods, approaching the performance of the best supervised methods in terms of object contact prediction accuracy. Acknowledgements This work was supported by the following awards: To D.L.K.Y.: Simons Foundation grant 543061, National Science Foundation CA-REER grant 1844724, Office of Naval Research grant S5122, ONR MURI 00010802 and ONR MURI S5847. We also thank the Google TPU Research Cloud team for computing support.

# References

- 1. Pytorch autograd, https://pytorch.org/tutorials/beginner/blitz/autograd\_ tutorial.html, accessed: March 3, 2024
- Agrawal, P., Nair, A.V., Abbeel, P., Malik, J., Levine, S.: Learning to poke by poking: Experiential learning of intuitive physics. Advances in neural information processing systems 29 (2016)
- Ajay, A., Bauza, M., Wu, J., Fazeli, N., Tenenbaum, J.B., Rodriguez, A., Kaelbling, L.P.: Combining physical simulators and object-based networks for control. In: 2019 International Conference on Robotics and Automation (ICRA). pp. 3217–3223. IEEE (2019)
- Babaeizadeh, M., Saffar, M.T., Nair, S., Levine, S., Finn, C., Erhan, D.: Fitvid: Overfitting in pixel-level video prediction. arXiv preprint arXiv:2106.13195 (2021)
- Bakhtin, A., van der Maaten, L., Johnson, J., Gustafson, L., Girshick, R.: Phyre: A new benchmark for physical reasoning. Advances in Neural Information Processing Systems **32** (2019)
- Bardes, A., Garrido, Q., Ponce, J., Rabbat, M., LeCun, Y., Assran, M., Ballas, N.: Revisiting feature prediction for learning visual representations from video. arXiv preprint (2024)
- Bates, C.J., Yildirim, I., Tenenbaum, J.B., Battaglia, P.: Modeling human intuitions about liquid flow with particle-based simulation. PLoS computational biology 15(7), e1007210 (2019)
- Battaglia, P., Pascanu, R., Lai, M., Jimenez Rezende, D., et al.: Interaction networks for learning about objects, relations and physics. Advances in neural information processing systems 29 (2016)
- Battaglia, P.W., Hamrick, J.B., Tenenbaum, J.B.: Simulation as an engine of physical scene understanding. Proceedings of the National Academy of Sciences 110(45), 18327–18332 (2013)
- Bear, D.M., Feigelis, K., Chen, H., Lee, W., Venkatesh, R., Kotar, K., Durango, A., Yamins, D.L.: Unifying (machine) vision via counterfactual world modeling. arXiv preprint arXiv:2306.01828 (2023)
- Bear, D.M., Wang, E., Mrowca, D., Binder, F.J., Tung, H.Y.F., Pramod, R., Holdaway, C., Tao, S., Smith, K., Sun, F.Y., et al.: Physion: Evaluating physical prediction from vision in humans and machines. arXiv preprint arXiv:2106.08261 (2021)
- Caron, M., Touvron, H., Misra, I., Jégou, H., Mairal, J., Bojanowski, P., Joulin, A.: Emerging properties in self-supervised vision transformers. In: Proceedings of the International Conference on Computer Vision (ICCV) (2021)
- Chen, H., Venkatesh, R., Friedman, Y., Wu, J., Tenenbaum, J.B., Yamins, D.L., Bear, D.M.: Unsupervised segmentation in real-world images via spelke object inference. In: European Conference on Computer Vision. pp. 719–735. Springer (2022)
- Chen, M., Radford, A., Child, R., Wu, J., Jun, H., Luan, D., Sutskever, I.: Generative pretraining from pixels. In: International conference on machine learning. pp. 1691–1703. PMLR (2020)

- 16 Venkatesh et al.
- Chen, T., Kornblith, S., Norouzi, M., Hinton, G.: A simple framework for contrastive learning of visual representations. In: International conference on machine learning. pp. 1597–1607. PMLR (2020)
- Chen, X., Wang, Y., Zhang, L., Zhuang, S., Ma, X., Yu, J., Wang, Y., Lin, D., Qiao, Y., Liu, Z.: Seine: Short-to-long video diffusion model for generative transition and prediction. arXiv preprint arXiv:2310.20700 (2023)
- Dave, I., Gupta, R., Rizve, M.N., Shah, M.: Tclr: Temporal contrastive learning for video representation. Computer Vision and Image Understanding 219, 103406 (2022)
- Doersch, C., Gupta, A., Efros, A.A.: Unsupervised visual representation learning by context prediction. In: Proceedings of the IEEE international conference on computer vision. pp. 1422–1430 (2015)
- Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al.: An image is worth 16x16 words: Transformers for image recognition at scale. arXiv preprint arXiv:2010.11929 (2020)
- Everingham, M., Van Gool, L., Williams, C.K., Winn, J., Zisserman, A.: The pascal visual object classes (voc) challenge. International journal of computer vision 88, 303–338 (2010)
- Feichtenhofer, C., Li, Y., He, K., et al.: Masked autoencoders as spatiotemporal learners. Advances in neural information processing systems 35, 35946–35958 (2022)
- Finn, C., Goodfellow, I., Levine, S.: Unsupervised learning for physical interaction through video prediction. Advances in neural information processing systems 29 (2016)
- Finn, C., Levine, S.: Deep visual foresight for planning robot motion. In: 2017 IEEE International Conference on Robotics and Automation (ICRA). pp. 2786– 2793. IEEE (2017)
- 24. Fragkiadaki, K., Agrawal, P., Levine, S., Malik, J.: Learning visual predictive models of physics for playing billiards. arXiv preprint arXiv:1511.07404 (2015)
- Geiger, A., Lenz, P., Stiller, C., Urtasun, R.: Vision meets robotics: The kitti dataset. The International Journal of Robotics Research 32(11), 1231–1237 (2013)
- Gerstenberg, T., Peterson, M.F., Goodman, N.D., Lagnado, D.A., Tenenbaum, J.B.: Eye-tracking causality. Psychological science 28(12), 1731–1744 (2017)
- 27. Gidaris, S., Singh, P., Komodakis, N.: Unsupervised representation learning by predicting image rotations. arXiv preprint arXiv:1803.07728 (2018)
- Goldberg, L.R.: The book of why: The new science of cause and effect: by judea pearl and dana mackenzie, basic books (2018). isbn: 978-0465097609. (2019)
- Goyal, R., Ebrahimi Kahou, S., Michalski, V., Materzynska, J., Westphal, S., Kim, H., Haenel, V., Fruend, I., Yianilos, P., Mueller-Freitag, M., et al.: The" something something" video database for learning and evaluating visual common sense. In: Proceedings of the IEEE international conference on computer vision. pp. 5842– 5850 (2017)
- Groth, O., Fuchs, F.B., Posner, I., Vedaldi, A.: Shapestacks: Learning vision-based physical intuition for generalised object stacking. In: Proceedings of the european conference on computer vision (eccv). pp. 702–717 (2018)
- Gupta, A., Tian, S., Zhang, Y., Wu, J., Martín-Martín, R., Fei-Fei, L.: Maskvit: Masked visual pre-training for video prediction. arXiv preprint arXiv:2206.11894 (2022)
- 32. Hafner, D., Lillicrap, T., Ba, J., Norouzi, M.: Dream to control: Learning behaviors by latent imagination. arXiv preprint arXiv:1912.01603 (2019)

- Hafner, D., Lillicrap, T., Fischer, I., Villegas, R., Ha, D., Lee, H., Davidson, J.: Learning latent dynamics for planning from pixels. In: International conference on machine learning. pp. 2555–2565. PMLR (2019)
- Han, J., Huang, W., Ma, H., Li, J., Tenenbaum, J.B., Gan, C.: Learning physical dynamics with subequivariant graph neural networks. In: Thirty-Sixth Conference on Neural Information Processing Systems (2022)
- He, K., Chen, X., Xie, S., Li, Y., Dollár, P., Girshick, R.: Masked autoencoders are scalable vision learners. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 16000–16009 (2022)
- He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.: Momentum contrast for unsupervised visual representation learning. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 9729–9738 (2020)
- 37. Heusel, M., Ramsauer, H., Unterthiner, T., Nessler, B., Hochreiter, S.: Gans trained by a two time-scale update rule converge to a local nash equilibrium. Advances in neural information processing systems **30** (2017)
- Höppe, T., Mehrjou, A., Bauer, S., Nielsen, D., Dittadi, A.: Diffusion models for video prediction and infilling. arXiv preprint arXiv:2206.07696 (2022)
- Janny, S., Baradel, F., Neverova, N., Nadri, M., Mori, G., Wolf, C.: Filtered-cophy: Unsupervised learning of counterfactual physics in pixel space. arXiv preprint arXiv:2202.00368 (2022)
- Kay, W., Carreira, J., Simonyan, K., Zhang, B., Hillier, C., Vijayanarasimhan, S., Viola, F., Green, T., Back, T., Natsev, P., et al.: The kinetics human action video dataset. arXiv preprint arXiv:1705.06950 (2017)
- 41. Li, Y., Wu, J., Tedrake, R., Tenenbaum, J.B., Torralba, A.: Learning particle dynamics for manipulating rigid bodies, deformable objects, and fluids. arXiv preprint arXiv:1810.01566 (2018)
- Lin, T.Y., Maire, M., Belongie, S., Hays, J., Perona, P., Ramanan, D., Dollár, P., Zitnick, C.L.: Microsoft coco: Common objects in context. In: Computer Vision– ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6-12, 2014, Proceedings, Part V 13. pp. 740–755. Springer (2014)
- Liu, R., Wu, R., Van Hoorick, B., Tokmakov, P., Zakharov, S., Vondrick, C.: Zero-1-to-3: Zero-shot one image to 3d object. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 9298–9309 (2023)
- 44. Liu, Y., Zhang, K., Li, Y., Yan, Z., Gao, C., Chen, R., Yuan, Z., Huang, Y., Sun, H., Gao, J., et al.: Sora: A review on background, technology, limitations, and opportunities of large vision models. arXiv preprint arXiv:2402.17177 (2024)
- 45. Lu, H., Yang, G., Fei, N., Huo, Y., Lu, Z., Luo, P., Ding, M.: Vdt: An empirical study on video diffusion with transformers. arXiv preprint arXiv:2305.13311 (2023)
- 46. Mehl, L., Schmalfuss, J., Jahedi, A., Nalivayko, Y., Bruhn, A.: Spring: A high-resolution high-detail dataset and benchmark for scene flow, optical flow and stereo. In: Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2023)
- 47. Mottaghi, R., Bagherinezhad, H., Rastegari, M., Farhadi, A.: Newtonian scene understanding: Unfolding the dynamics of objects in static images. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 3521– 3529 (2016)
- Mrowca, D., Zhuang, C., Wang, E., Haber, N., Fei-Fei, L.F., Tenenbaum, J., Yamins, D.L.: Flexible neural representation for physics prediction. Advances in neural information processing systems **31** (2018)
- 49. Nair, S., Rajeswaran, A., Kumar, V., Finn, C., Gupta, A.: R3m: A universal visual representation for robot manipulation. arXiv preprint arXiv:2203.12601 (2022)

- 18 Venkatesh et al.
- Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: European conference on computer vision. pp. 69–84. Springer (2016)
- Oord, A.v.d., Li, Y., Vinyals, O.: Representation learning with contrastive predictive coding. arXiv preprint arXiv:1807.03748 (2018)
- 52. OpenAI: Gpt-4 for vision (chatgpt with image input) (2023), https://openai.com/, accessed: October 27, 2023
- 53. Oquab, M., Darcet, T., Moutakanni, T., Vo, H.V., Szafraniec, M., Khalidov, V., Fernandez, P., Haziza, D., Massa, F., El-Nouby, A., Howes, R., Huang, P.Y., Xu, H., Sharma, V., Li, S.W., Galuba, W., Rabbat, M., Assran, M., Ballas, N., Synnaeve, G., Misra, I., Jegou, H., Mairal, J., Labatut, P., Joulin, A., Bojanowski, P.: Dinov2: Learning robust visual features without supervision (2023)
- Pathak, D., Girshick, R., Dollár, P., Darrell, T., Hariharan, B.: Learning features by watching objects move. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2701–2710 (2017)
- Piloto, L.S., Weinstein, A., Battaglia, P., Botvinick, M.: Intuitive physics learning in a deep-learning model inspired by developmental psychology. Nature human behaviour 6(9), 1257–1267 (2022)
- Pont-Tuset, J., Perazzi, F., Caelles, S., Arbelaez, P., Sorkine-Hornung, A., Gool, L.V.: The 2017 davis challenge on video object segmentation. arXiv preprint arXiv:1704.00675 (2017)
- Qi, H., Wang, X., Pathak, D., Ma, Y., Malik, J.: Learning long-term visual dynamics with region proposal interaction networks. arXiv preprint arXiv:2008.02265 (2020)
- Qian, R., Meng, T., Gong, B., Yang, M.H., Wang, H., Belongie, S., Cui, Y.: Spatiotemporal contrastive video representation learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 6964– 6974 (2021)
- Riochet, R., Castro, M.Y., Bernard, M., Lerer, A., Fergus, R., Izard, V., Dupoux, E.: Intphys 2019: A benchmark for visual intuitive physics understanding. IEEE Transactions on Pattern Analysis and Machine Intelligence 44(9), 5016–5025 (2021)
- 60. Sanchez-Gonzalez, A., Heess, N., Springenberg, J.T., Merel, J., Riedmiller, M., Hadsell, R., Battaglia, P.: Graph networks as learnable physics engines for inference and control. In: International Conference on Machine Learning. pp. 4470–4479. PMLR (2018)
- Smith, K., Mei, L., Yao, S., Wu, J., Spelke, E., Tenenbaum, J., Ullman, T.: Modeling expectation violation in intuitive physics with coarse probabilistic object representations. Advances in neural information processing systems **32** (2019)
- 62. Spelke, E.S.: Principles of object perception. Cognitive science 14(1), 29-56 (1990)
- Tacchetti, A., Song, H.F., Mediano, P.A., Zambaldi, V., Rabinowitz, N.C., Graepel, T., Botvinick, M., Battaglia, P.W.: Relational forward models for multi-agent learning. arXiv preprint arXiv:1809.11044 (2018)
- Teed, Z., Deng, J.: Raft: Recurrent all-pairs field transforms for optical flow. In: Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part II 16. pp. 402–419. Springer (2020)
- Tong, Z., Song, Y., Wang, J., Wang, L.: Videomae: Masked autoencoders are dataefficient learners for self-supervised video pre-training. Advances in neural information processing systems 35, 10078–10093 (2022)

- Voleti, V., Jolicoeur-Martineau, A., Pal, C.: Mcvd-masked conditional video diffusion for prediction, generation, and interpolation. Advances in Neural Information Processing Systems 35, 23371–23385 (2022)
- 67. Wang, L., Huang, B., Zhao, Z., Tong, Z., He, Y., Wang, Y., Wang, Y., Qiao, Y.: Videomae v2: Scaling video masked autoencoders with dual masking. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14549–14560 (2023)
- Wang, X., Gupta, A.: Unsupervised learning of visual representations using videos. In: Proceedings of the IEEE international conference on computer vision. pp. 2794–2802 (2015)
- Wang, X., Yu, Z., De Mello, S., Kautz, J., Anandkumar, A., Shen, C., Alvarez, J.M.: Freesolo: Learning to segment objects without annotations. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 14176–14186 (2022)
- Wang, X., Girdhar, R., Yu, S.X., Misra, I.: Cut and learn for unsupervised object detection and instance segmentation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 3124–3134 (2023)
- Wu, Z., Xiong, Y., Yu, S.X., Lin, D.: Unsupervised feature learning via nonparametric instance discrimination. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3733–3742 (2018)
- Wu, Z., Dvornik, N., Greff, K., Kipf, T., Garg, A.: Slotformer: Unsupervised visual dynamics simulation with object-centric models. arXiv preprint arXiv:2210.05861 (2022)
- 73. Yan, W., Hafner, D., James, S., Abbeel, P.: Temporally consistent transformers for video generation (2023)
- Ye, Y., Singh, M., Gupta, A., Tulsiani, S.: Compositional video prediction. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 10353–10362 (2019)
- 75. You, Y., Lou, Y., Li, C., Cheng, Z., Li, L., Ma, L., Lu, C., Wang, W.: Keypointnet: A large-scale 3d keypoint dataset aggregated from numerous human annotations. arXiv preprint arXiv:2002.12687 (2020)
- Yu, L., Cheng, Y., Sohn, K., Lezama, J., Zhang, H., Chang, H., Hauptmann, A.G., Yang, M.H., Hao, Y., Essa, I., et al.: Magvit: Masked generative video transformer. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 10459–10469 (2023)
- Zhang, R., Isola, P., Efros, A.A.: Colorful image colorization. In: Computer Vision– ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part III 14. pp. 649–666. Springer (2016)
- Zheng, Y., Harley, A.W., Shen, B., Wetzstein, G., Guibas, L.J.: Pointodyssey: A large-scale synthetic dataset for long-term point tracking. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 19855–19865 (2023)
- Zhuang, C., She, T., Andonian, A., Mark, M.S., Yamins, D.: Unsupervised learning from video with deep neural embeddings. In: Proceedings of the ieee/cvf conference on computer vision and pattern recognition. pp. 9563–9572 (2020)