# Understanding Video Transformers via Universal Concept Discovery

Matthew Kowal[1,3*]    Achal Dave[3]    Rares Ambrus[3]
Adrien Gaidon[3]    Konstantinos G. Derpanis[1,2]    Pavel Tokmakov[3]
[1]York University, [2]Samsung AI Centre Toronto, [3]Toyota Research Institute
yorkucvil.github.io/VTCD

## Abstract

*This paper studies the problem of concept-based interpretability of transformer representations for videos. Concretely, we seek to explain the decision-making process of video transformers based on high-level, spatiotemporal concepts that are automatically discovered. Prior research on concept-based interpretability has concentrated solely on image-level tasks. Comparatively, video models deal with the added temporal dimension, increasing complexity and posing challenges in identifying dynamic concepts over time. In this work, we systematically address these challenges by introducing the first Video Transformer Concept Discovery (VTCD) algorithm. To this end, we propose an efficient approach for unsupervised identification of units of video transformer representations - concepts, and ranking their importance to the output of a model. The resulting concepts are highly interpretable, revealing spatiotemporal reasoning mechanisms and object-centric representations in unstructured video models. Performing this analysis jointly over a diverse set of supervised and self-supervised representations, we discover that some of these mechanism are universal in video transformers. Finally, we demonstrate that VTCD can be used to improve model performance for fine-grained tasks.*

## 1. Introduction

Understanding the hidden representations within neural networks is essential for addressing regulatory concerns [11, 31], preventing harms during deployment [5, 30], and can aid innovative model designs [13]. This problem has been studied extensively in the image world, both for convolutional neural networks (CNNs) [4, 22, 26, 35] and, more recently, for transformers [51, 64], resulting in a number of key insights. For example, image classification models extract low-level positional and texture cues at early layers and gradually combine them into higher-level, semantic
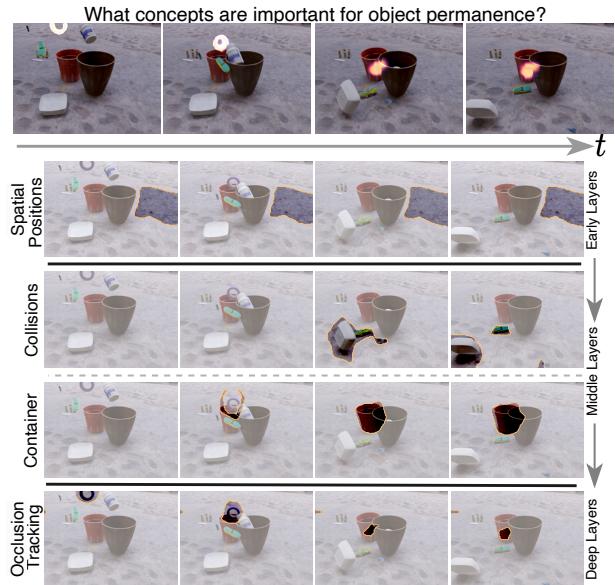
Figure 1. Heatmap predictions of the TCOW model [62] for tracking through occlusions (top), together with concepts discovered by our VTCD (bottom). We can see that the model encodes positional information in early layers, identifies containers and collisions events in mid-layers and tracks through occlusions in late layers. Only one video is shown, but the discovered concepts are shared between many dataset samples (see video for full results).

concepts at later layers [4, 24, 45].

However, while video transformers do share their overall architecture with image-level ViTs, the insights obtained in existing works do very little to explain their inner mechanisms. Consider, for example, the recent approach for tracking occluded objects [62] shown in Figure 1 (top). To accurately reason about the trajectory of the invisible object inside the pot, texture or semantic cues alone would not suffice. What, then, are the *spatiotemporal* mechanisms used by this approach? And are any of these mechanisms *universal* across video models trained for different tasks?

To answer these questions, in this work we propose the

Video Transformer Concept Discovery algorithm (VTCD) - the first concept-discovery methodology for interpreting the representations of deep video transformers. We focus on concept-based interpretability [21, 22, 26, 69] due to its capacity to explain the decision-making process of a complex model's distributed representations in high-level, intuitive terms. Our goal is to decompose a representation at any given layer into human-interpretable 'concepts' without any labelled data (i.e., concept discovery) and then rank them in terms of their importance to the model output.

Concretely, we first group *model features* at a given layer into spatio-temporal tubelets via SLIC clustering [1], which serve as a basis for our analysis (Section 3.1.1). Next, we cluster these tubelets across videos to discover high-level concepts [14, 21, 22, 39, 69] (Section 3.1.2) . The resulting concepts for an occluded object tracking method [62] are shown in Figure 1 (bottom) and span a broad range of cues, including spatiotamporal ones that detect events, like collisions, or track the containers.

To better understand the decision-making mechanisms of video transformers, we then quantify the importance of concepts for the model's predictions. Inspired by previous work on saliency maps [49], we propose a novel, noise-robust approach to estimate concept importance (Section 3.2). Unlike existing techniques that rely on gradients [35], or concept occlusion [21], our approach effectively handles redundancy in self-attention heads in transformer architectures.

Next, we use VTCD to study whether there are any universal mechanisms in video transformer models, that emerge irrespective of their training objective. To this end, we extend the recent work by Dravid et al. [16] to automatically identify *important* concepts that are shared between several models in Section 4.1. We then analyze a diverse set of representations (*e.g.* supervised, self-supervised, or video-language) and make a number of discoveries: (i) many concepts are indeed shared between models trained for different tasks; (ii) early layers tend to form a spatiotemporal basis that underlines the rest of the information processing; (iii) later layers form object-centric video representations even in models trained in a self-supervised way.

Finally, VTCD can be used to turn a pre-trained video transformer into an efficient and effective fine-grained recognition model by pruning least important units. We demonstrate in Section 5.4 that removing one third of the heads from an action classification model results in a 4.3% increase in accuracy while reducing computation by 33%.

## 2. Related work

Our work proposes a novel *concept-based interpretability* algorithm that focuses on *transformer-based representations* for *video understanding*. Below, we review the most relevant works in each of these fields.

**Concept-based interpretability** is a family of neural network interpretability methods used to understand, post-hoc, the representations that a model utilizes for a given task. Closed-world interpretability operates under the premise of having a labeled dataset of concepts [4, 35]. However, for videos, it is unclear what concepts may exist and also difficult to densely label videos even if they were known a priori.

In contrast, unsupervised concept discovery makes no assumptions on the existence of semantic concepts and uses clustering to partition data into interpretable components within the model's feature space. ACE [26] and CRAFT [22] segment input images into superpixels and random crops, before applying clustering at a given layer. In videos, however, the potential tubelets far outnumber image crops, prompting us to introduce a more efficient method for segmenting videos into proposals in Section 3.1.1.

A necessary component of concept-based interpretability is measuring the importance (*i.e.* fidelity) of the discovered concepts to the model. However, the aforementioned methods [21, 22, 26, 35, 69] were developed for CNNs, and are not readily applicable to transformers. The main challenge of ranking concepts in attention heads is due to the transformers' robustness to minor perturbations in self-attention layers. To address this limitation, we introduce a new algorithm to rank the significance of any architectural unit, covering both heads and intra-head concepts in Section 3.2.

Recent work [16] identifies neurons that produce similar activation maps across various image models (including transformers). However, neurons are unable to explain a full extent of a models' representation due to it being distributed across many dimensions [18]. In contrast, our method works on any dimensional features and is applied to video-based models.

**Interpretability of transformers** has received a significant amount of attention recently, due to the success of this architecture in variety of computer vision problems. Early work [51] contrasted vision transformer representations with CNNs (representational differences per layer, receptive fields, localization of information, etc). Other work aims to generate saliency heatmaps based on attention maps of a model [8, 9]. Later works focused on understanding the impact of different training protocols [47, 64] (*e.g.* self-supervised learning (SSL) vs. supervised) and robustness [50, 70]. The features of a specific SSL vision transformer, DINO [2], were explored in detail and shown to have surprising utility for part-based segmentation tasks. However, none of these works address concept-based interpretability or study video representations.

Independently, studies in natural language processing (NLP) have analyzed self-attention layers [17, 63] and found that heads are often specialized to capture different linguistic or grammatical phenomenon. This is qualitatively seen in works that shows dissimilar attention maps for different self-attention heads [12, 34]. Moreover, other NLP
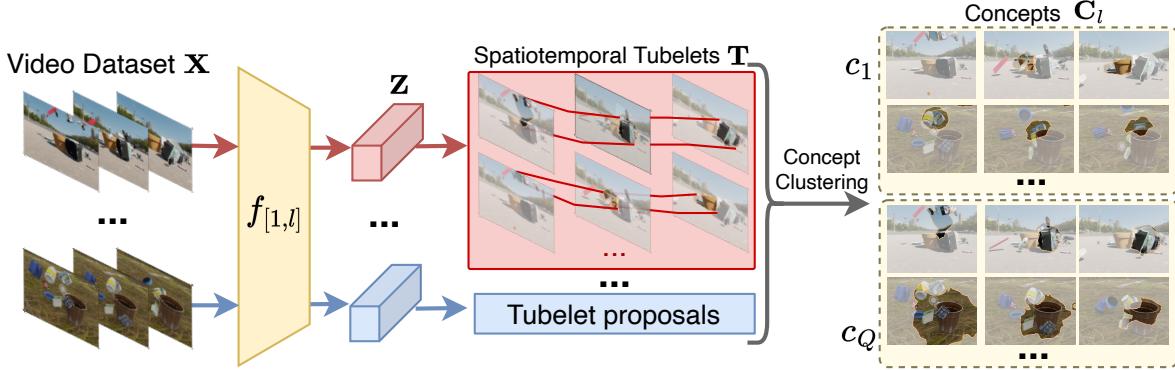
Figure 2. Video Transformer Concept Discovery (VTCD) takes a dataset of videos, $\mathbf{X}$, as input and passes them to a model, $f_{[1,l]}$ (shown in yellow). The set of video features, $\mathbf{Z}$, are then parsed into spatiotemporal tubelet proposals, $\mathbf{T}$ (shown in red), via SLIC clustering in the feature space. Finally, tubelets are clustered across the videos to discover high-level units of network representation - concepts, $\mathbf{C}$ (right).

works [44, 63] explore the impact of removing heads and find that only a small number need to be kept to produce similar performance. Our findings agree with evidence from these works, and in Section 5.4 we further demonstrate that pruning unimportant heads from video transformers can actually *improve* model's performance.

**Video model interpretability** is an under-explored area of research considering the recent successes of deep learning models in action recognition [36], video object segmentation [40, 48, 59, 62], or self-supervised approaches [20, 52, 58, 60, 65]. Efforts have used proxy tasks to measure the degree to which models use dynamic information [25, 29, 33] or scene bias [10, 41, 42]. One method quantifies the static and dynamic information contained in a video model's intermediate representation [37, 38]. However, these methods can only measure one or two predefined concepts (*i.e.* static, dynamic, or scene information) while our approach is not restricted to a subset of concepts. Another work visualizes videos that activate single neurons (or filters) in 3D CNN's via activation maximization with temporal regularization [19]. While this method has no restrictions on what a neuron can encode, it only applies to 3D CNNs and does not truly capture 'concepts' across distributed representations (*i.e.* feature space directions that generalize across videos).

## 3. Video transformer concept discovery

We study the problem of decomposing a video representation into a set of high-level open-world concepts and ranking their importance for the model's predictions. We are given a set of (RGB) videos, $\mathbf{X} \in \mathbb{R}^{N \times 3 \times T \times H \times W}$, where $N$, $T$, $H$, and $W$ denote the dataset size, time, height, and width, respectively, and an $L$ layer pretrained model, $f$. Let $f_{[r,l]}$ denote the model from layer $r$ to $l$, with $f_{[1,l]}(\mathbf{X}) = \mathbf{Z}_l \in \mathbb{R}^{N \times C \times T' \times H' \times W'}$ being the intermediate representation at layer $l$. To decompose $\mathbf{Z}_l$ into a set of human-interpretable concepts, $\mathbf{C}_l = \{c_1, ..., c_Q\}$, existing, image-level approaches [22, 26] first parse the $N$ feature maps into a set of $M$ proposals, $\mathbf{T} \in \mathbb{R}^{M \times C}$ $(M > N)$, where each $T_m$ corresponds to a region of the input image. These proposals are then clustered into $Q << M$ concepts in the feature space of the model to form an assignment matrix $W \in \mathbb{R}^{M \times Q}$. Finally, the importance of each concept $c_i$ to the model's prediction is quantified by a score $s_i \in [0, 1]$. Performing this analysis over all the layers in $f$ produces the entire set of concepts for a model, $\mathbf{C} = \{\mathbf{C}_1, ..., \mathbf{C}_L\}$, together with their corresponding importance scores.

However, existing approaches are not immediately applicable to video transformers because they do not scale well and are focused on 2D CNN architectures. In this work, we extend concept-based interpretability to video representations. To this end, we first describe a computationally tractable proposal generation method (Section 3.1.1) that operates over space-time feature volumes and outputs spatiotemporal tublets. Next (Section 3.1.2), we adapt existing concept clustering techniques to video transformer representations. Finally, in Section 3.2 we propose CRIS - a novel concept importance estimation approach applicable to any architecture units, including transformer heads.

### 3.1. Concept discovery

#### 3.1.1 Tubelet proposals

Previous methods [21, 22, 26] use superpixels or crops in RGB space to propose segments; however, the number of possible segments is exponentially greater for videos. Moreover, proposals in color space are unrestricted and may not align with the model's encoded information, leading to many irrelevant or noisy segments. To address these drawbacks, we instantiate proposals in *feature space*, which naturally partitions a video based on the information contained within each layer (shown in Figure 2, left).

More specifically, we construct tubelets per video via Simple Linear Iterative Clustering [1, 32] (SLIC) on the
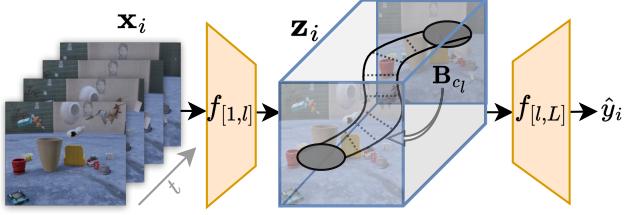
3

Figure 3. A visual representation of concept masking for a single concept. Given a video $\mathbf{x_i}$ and a concept, $c_l$, we mask the tokens of the intermediate representation $\mathbf{z_i} = f_{[1,l]}(\mathbf{x_i})$ with the concepts' binary support masks, $\mathbf{B}_{c_l}$, to obtain the perturbed prediction, $\hat{y}_i$.

spatiotemporal features via

$$\mathbf{T} = \text{GAP}(\mathbf{B} \odot \mathbf{Z}) = \text{GAP}(\text{SLIC}(\mathbf{Z}) \odot \mathbf{Z}), \quad (1)$$

where $\mathbf{T} \in \mathbb{R}^{M \times C}$ is the set of tubelets for the dataset, $\mathbf{B} \in \{0,1\}^{C \times M \times N \times T' \times H' \times W'}$ are spatiotemporal binary support masks obtained from the SLIC clustering, $M$ is the total number of tubelets for all $N$ videos ($M >> N$), and GAP is a global average pooling operation over the space and time dimensions.

SLIC is an extension of the K-Means algorithm that controls a trade-off between cluster support regularity and adaptability, and also constrains cluster support masks to be connected. Together these properties produce non-disjoint tubelets that are easier to interpret for humans because they reduce the need to attend to multiple regions in a video at a time. Further, the pruning step in SLIC makes it more robust to the hyperparameter that controls the desired number of clusters, as it automatically prunes spurious, disconnected tubelets. Next, we describe our approach for grouping individual tubelets into higher-level concept clusters.

### 3.1.2 Concept clustering

Recent work [21, 22, 69] has used Non-Negative Matrix Factorization (NMF) [14] to cluster proposals into concepts. Given a non-negative data matrix, $\mathbf{T}^+ \in \mathbb{R}^{M \times C}$, NMF aims to find two non-negative matrices, $\mathbf{W}^+ \in \mathbb{R}^{M \times Q}$ and $\mathbf{C}^+ \in \mathbb{R}^{Q \times C}$, such that $\mathbf{T}^+ = \mathbf{W}^+ \mathbf{C}^+$, where $\mathbf{W}^+$ is the cluster assignment matrix. Unfortunately, NMF cannot be applied to transformers as they use GeLU non-linearities, rather than ReLU, resulting in negative activations.

We solve this problem by leveraging Convex Non-negative Matrix Factorization [14] (CNMF). Despite the name, CNMF extends NMF and allows for negative input values. This is achieved by constraining the factorization such that the columns of $\mathbf{W}$ are convex combinations of the columns of $\mathbf{T}$, i.e. each column of $\mathbf{W}$ is a weighted average of the columns of $\mathbf{T}$. This constraint can be written as

$$\mathbf{W} = \mathbf{TG}, \quad (2)$$

where $\mathbf{G} \in [0,1]^{C \times Q}$ and $\sum_j \mathbf{G}_{i,j} = 1$. To cluster a set of tublets, $\mathbf{T}$, into corresponding concepts, we optimize

$$(\mathbf{G}^*, \mathbf{C}^*) = \underset{\mathbf{C} > 0, \mathbf{G} > 0}{\arg\min} ||\mathbf{T} - \mathbf{TGC}||^2, \quad (3)$$

where the final set of concepts are represented by the rows of the matrix $\mathbf{C}$, i.e. concept centroid $c_i$ is the $i^{th}$ row of $\mathbf{C}$ (Figure 2, right).

### 3.2. Concept importance

Given a set of discovered concepts, we now aim to quantify their impact on model performance. One approach, shown in Figure 3, is to mask out each concept independently and rank the importance based on the drop in performance [21]. Formally, let $c_l$ be a single target concept, and $\mathbf{B}_{c_l} \in \{0,1\}^{C \times M \times N \times T' \times H' \times W'}$ the corresponding binary support masks over $\mathbf{X}$. It can then be masked in layer $l$ via

$$\hat{y} = f_{[l,L]}(\mathbf{Z}_l \odot (1 - \mathbf{B}_{c_l})). \quad (4)$$

However, while this approach works well for CNNs [21], transformers are robust to small perturbations within self-attention layers [44, 63]. Therefore, single concept masking has little effect on performance (shown by results in Figure 4). Instead, we mask *a high percentage* of sampled concepts in parallel (across all layers and heads) and then empirically validate in Section 5.1 that averaging the results over thousands of samples produces valid concept rankings.

Formally, we propose **C**oncept **R**andomized **I**mportance **S**ampling (CRIS), a robust method to compute importance for any unit of interest. To this end, we first randomly sample $K$ different concept sets, such that each $\mathbf{C}^k \subset \mathbf{C}$. We then define $\mathbf{C}_l^k$ as the set of concepts in $\mathbf{C}^k$ discovered at layer $l$, with $\mathbf{B}_{\mathbf{C}_l^k}$ denoting the corresponding binary support masks. We mask out every concept at every layer of the model via

$$\hat{y}_k = g(\tilde{\mathbf{B}}_{\mathbf{c}_L^k} \odot f_{[L-1,L]}(\cdots(\tilde{\mathbf{B}}_{\mathbf{C}_1^k} \odot f_{[0,1]}(\mathbf{X})))), \quad (5)$$

where $g(\cdot)$ is the prediction head (e.g. an MLP) and $\tilde{\mathbf{B}}$ denotes the inverse mask (i.e. $1 - \mathbf{B}$). Finally, we calculate the importance of each concept, $c_i$, via

$$s_i = \frac{1}{K} \sum_k^K (\mathbb{D}(\tilde{y}, y) - \mathbb{D}(\hat{y}_k, y)) \mathbb{1}_{c_i \in \mathbf{C}^k}, \quad (6)$$

where $\tilde{y}$ is the original prediction without any masking and $\mathbb{D}$ is a metric quantifying performance (e.g. accuracy).

## 4. Understanding transformers with VTCD

Our algorithm facilitates the identification of concepts within any unit of a model and quantifying their significance in the final predictions. However, this is not enough

to fully represent the computations performed by a video transformer. It is also crucial to understand how these concepts are employed in the model's information flow.

As several recent works have shown [17, 46], the residual stream of a transformer serves as the backbone of the information flow. Each self-attention block then reads information from the residual stream with a linear projection, performs self-attention operations to process it, and finally writes the results back into the residual stream. Crucially, self-attention processing is performed individually for each head and several studies have shown, both in vision [2, 15, 34] and NLP [63], that different self-attention heads capture distinct information. In other word, heads form the basis of the transformer representation.

A closer analysis of the concepts found in the heads of TCOW with VTCD allows us to identify several patterns in the information processing of that model. In particular, Figure 1 shows that the heads in early layers group input tokens based on their spatiotemporal positions. This information is then used to track objects and identify events in mid-layers, and later layers utilize mid-layer representations to reason about occlusions. Next, we study whether any of these mechanisms are *universal* across video transformers trained on different datasets with varying objectives.

### 4.1. Rosetta concepts

Inspired by [16], we propose to mine for *Rosetta concepts* that are shared between models and represent the same information. The key to identifying Rosetta units is a robust metric, $R$, where a higher R-score corresponds to the two units having a larger amount of shared information. Previous work [16] focused on finding such neurons in image models based on correlating their activation maps. We instead propose to measure the similarity between concepts (*i.e.* distributed representations) via the mean Intersection over Union (mIoU) of the concepts' support.

Formally, we mine Rosetta concepts by first applying VTCD to a set of $D$ models $\{f^1, ..., f^D\}$, resulting in discovered concepts, $\mathbf{C}^j = \{c_1^j, ..., c_i^j\}$, and importance scores, $\mathbf{S}^j = \{s_1^j, ..., s_i^j\}$, for each model $f^j$. We then aim to measure the similarity between all concept $D$-tuples from the models. Given a set of $D$ concepts, $\{c_i^1, ..., c_i^D\}$ and corresponding binary support masks, $\{\mathbf{B}_i^1, ..., \mathbf{B}_i^D\}$, we define the similarity score of these concepts as

$$R_i^D = \frac{|\mathbf{B}_i^1 \cap \cdots \cap \mathbf{B}_i^D|}{|\mathbf{B}_i^1 \cup \cdots \cup \mathbf{B}_i^D|}. \qquad (7)$$

Naively computing the similarity between all $D$-tuples results in an exponential number of computations and is intractable for even small $D$. To mitigate this issues, we exclude two types of concepts: (i) unimportant ones and (ii) those with a low R-score among $d$-tuples where $d < D$. More specifically, we only consider the most important

$\epsilon\%$ of concepts from each model. We then iterate over $d \in \{2, ..., D\}$ and filter out any concept that has an R-score less than $\delta$ for all d-tuples in which it participates. Formally, the filtered Rosetta d-concept scores are defined as

$$\mathbf{R}_{\epsilon, \delta}^d = \{R_i^d \mid R_i^d > \delta \, \forall R_i^d \in \mathbf{R}_\epsilon^d\}, \qquad (8)$$

where $\mathbf{R}_\epsilon^d$ is the set of all R-scores among $d$ concepts after the $\epsilon$ importance filtering. This results in a significantly smaller pool of candidates for the next stage $d + 1$, reducing the overall computational complexity of the algorithm. Finally, as some concepts may reside in a subset of the models but are still interesting to study, we examine the union of all important and confident Rosetta d-concepts corresponding to R-scores $\mathbf{R}_{\epsilon, \delta}^2 \cup \cdots \cup \mathbf{R}_{\epsilon, \delta}^D$.

## 5. Experiments

We evaluate the quality of our concept discovery algorithm quantitatively and qualitatively across a variety of models trained for different tasks.

**Datasets.** We use two datasets in our experiments: TCOW Kubric [62] and Something-Something-v2 (SSv2) [27]. The former is a synthetic, photo-realistic dataset of 4,000 videos with randomized object location and motion, based on the Kubric synthetic video generator [28]. This dataset is intended for semi-supervised video object segmentation (semi-VOS) through occlusions. SSv2 contains 220,847 real videos intended for finegrained action recognition. Each sample is a crowdsourced video of a person-object interaction (*i.e.* doing something to something). Unlike many other video classification benchmarks [6, 56], temporal reasoning is fundamental for distinguishing SSv2 actions, making it an ideal choice for analyzing spatiotemporal mechanisms in transformers.

**Models.** We evaluate four models with public pretrained checkpoints: (i) TCOW [62] trained on Kubric for semi-VOS, (ii) VideoMAE [60] trained on SSv2 for action classification (Supervised VideoMAE), (iii) VideoMAE self-supervised on SSv2 (SSL VideoMAE), and (iv) Intern-Video [66], a video-text foundation model trained contrastively on 12M video clips from eight video datasets and 100M image-text pairs from LAION-400M [54]. As TCOW requires a segmentation mask as input, when applying it to SSv2, we hand label the most salient object in the initial frame to use as the query. We focus our analysis on the first two models, and use the last two to validate the universality of our Rosetta concepts.

**Implementation details.** For all experiments, we run VTCD for 30 randomly sampled videos, and discover concepts from every head at every model layer. Prior work [2] shows Keys produce the most meaningful clusters in self-attention heads, so we focus here on Keys and present results with Queries and Values, together with the rest of the
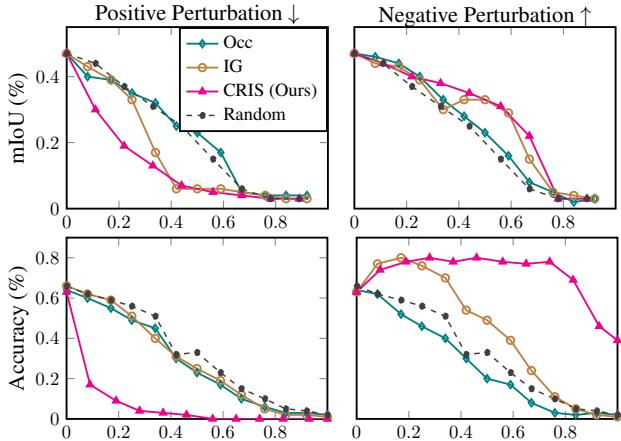
5

Figure 4. Concept attribution curves for every layer of TCOW trained on Kubric (top) and VideoMAE trained on SSv2 (bottom). We remove concepts ordered from most-to-least (left) or least-to-most important (right). CRIS produces better concept importance rankings than methods based on occlusion (Occ) or gradients (IG).

hyperparameters, in the appendix. Our code for the VTCD toolbox allowing to analyze any video transformer representation will be released.

**VTCD target metrics.** To discover and rank concepts, VTCD requires a target evaluation metric. For TCOW Kubric, we use the intersection-over-union (IoU) between the predicted and the groundtruth masks. For SSv2, we use classification accuracy for a target class.

## 5.1. Quantitative concept evaluation

To quantitatively confirm the effectiveness of our method, we follow the standard evaluation protocol for concept-based interpretability, and measure the *fidelity* of the discovered concepts [22, 26, 69]. To this end, we calculate attribution curves [21, 22, 26], where concepts (across all layers) are removed from a model in either most-to-least order (*positive* perturbation), or least-to-most (*negative* perturbation). The intuition is that concepts with higher fidelity and more accurate importance scores will have a steeper performance decrease when removing the most important concepts, and vice-versa for the reverse order.

In Figure 4, we plot concept attribution curves of our method for TCOW for Kubric (top) and Supervised Video-MAE for SSv2 (bottom), targeting 10 randomly sampled classes and averaging results (bottom). In addition, we report results for several baselines: (i) concept removal in a random order, (ii) standard, occlusion-based [21] concept importance estimation, and (iii) a gradient based approach [21, 35]. In all cases, CRIS produces a more viable importance ranking, dramatically outperforming both random ordering and the occlusion baseline. Integrated gradients approach performs similarly to ours for TCOW, but is significantly worse for the action classification VideoMAE model.

Notably, we observe that the performance actually *increases* for VideoMAE when up to 70% of the least important concepts are removed. Recall that SSv2 is a fine-grained action classification dataset. Our method removes concepts that are irrelevant for the given class, hence increasing the robustness of the model's predictions. In Section 5.4, we elaborate on this observation and demonstrate how VTCD can be used improve performance and efficiency of any fine-grained video transformer.

## 5.2. Qualitative analysis

We have seen that the importance assigned to concepts discovered by VTCD aligns well with the accuracy of the model. We now turn to assessing the concepts themselves qualitatively. To this end, in Figure 5 we show two representative videos for the top three most important concepts for the TCOW and VideoMAE models for the class *dropping something into something*.

For TCOW, the most important concept occurs in layer five and tracks the target object. Interestingly, the same concept highlights objects with similar appearance and 2D position to the target. This suggests that the model solves the disambiguation problem by first identifying possible distractors in mid-layers (*i.e.* five) and then using this information to more accurately track the target in final layers. In fact, the second most important concept, occurring in layer nine, tracks the target object throughout the video.

For VideoMAE, the most important concepts highlights the object being dropped until the dropping event, at which point both the object and the container are highlighted. The second most important concept clearly captures the container being dropped into, notably not capturing the object itself and making a ring-like shape. These concepts identify an important mechanism of the model that helps it differentiating between similar classes (*e.g.* dropping something *into/behind/in-front* of something).

The third most important concept for each model capture similar information, occurring in early layers: a temporally invariant, spatial support. This observation corroborates research [2, 24] suggesting that positional information processing occurs early in the model and acts as a reference frame between the semantic information and the tokens themselves. We now set out to discover concepts that are *shared* across multiple models, termed Rosetta concepts.

## 5.3. Rosetta concepts

We begin by applying VTCD to the four models, targeting two classes chosen due to their dynamic nature: *rolling something across a flat surface* and *dropping something behind something*. We then mine for Rosetta concepts using the method described in Section 4.1 and setting $\delta = 0.15$ and $\epsilon = 15\%$ in all experiments. The resulting set of Rosetta 4-concepts contains 40 tuples with an average $R$
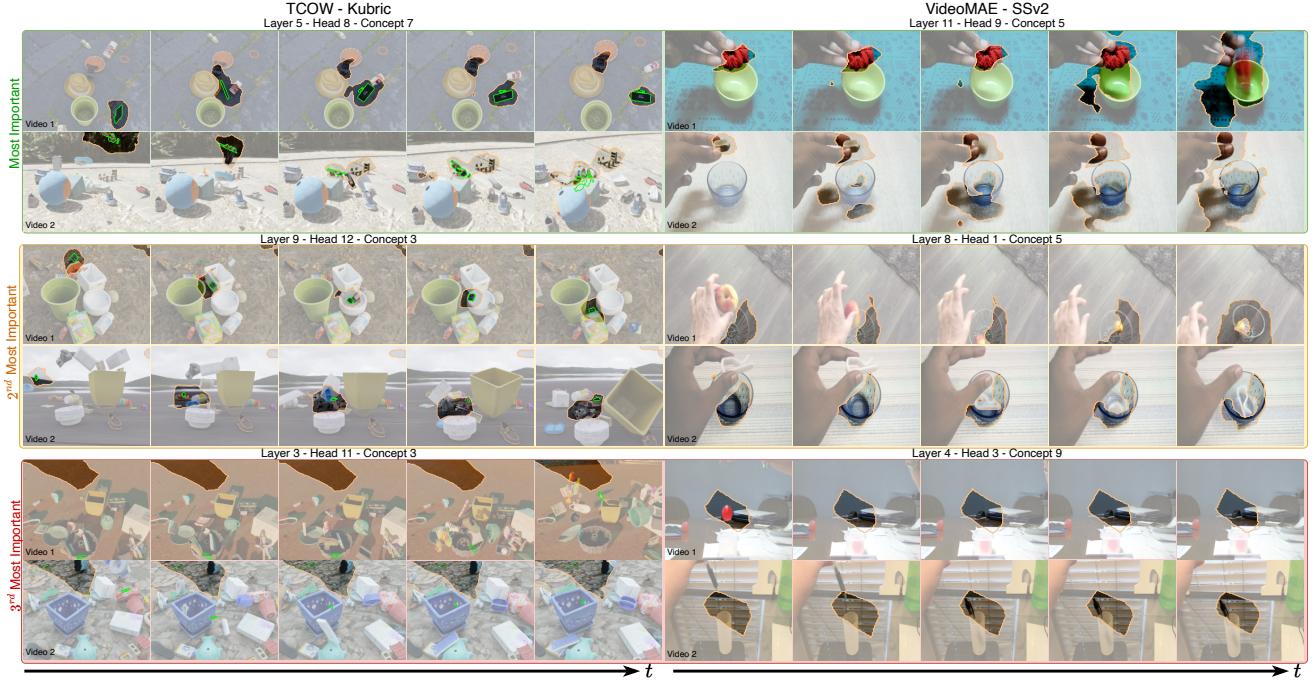
Figure 5. The top-3 most important concepts for the TCOW model trained on Kubric (left) and VideoMAE trained on SSv2 for the target class *dropping something into something* (right). Two videos are shown for each concept and the query object is denoted with a green border in Kubric. For TCOW, the $1^{st}$ and $2^{nd}$ (top-left, middle-left) most important concepts track multiple objects including the target and the distractors. For VideoMAE, the top concept (top-right) captures the object and dropping event (*i.e.* hand, object and container) while the $2^{nd}$ most important concept (middle-right) captures solely the container. Interestingly, for both models and tasks, the third most important concept (bottom) is a temporally invariant tubelet. See further discussion in Section 5.3 (and video for full results).
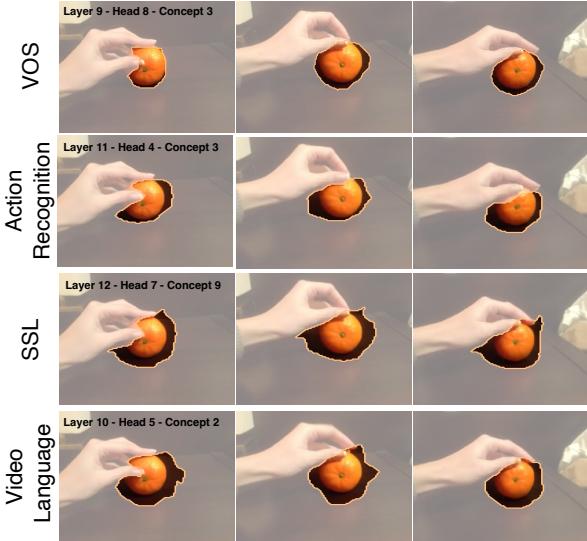


Figure 6. A sample Rosetta concept found in four models trained for different tasks. Interestingly, we find object-centric representations in all the models (see video for full results).

score (Equation 7) of 17.1. For comparison, the average $R$ score between all possible 4-concepts is 0.6, indicating the significance of the selected matches. We visualize one of the mined Rosetta 4-concepts in Figure 6, which captures an object tracking representation in the late layers. This re-

sult demonstrates that universal representations indeed exist between all four models. We show many more examples of shared concepts on the project web page in video format.

Next, we qualitatively analyze all Rosetta d-concept with $d \in \{2, 3, 4\}$ at various layers, and show a representative sample in Figure 7. Firstly, we observe that in early layers the models learn spatiotemporal basis representations (Figure 7, left). That is, they decompose the space-time volume of a video into connected regions that facilitate higher-level reasoning in later layers. This is consistent with prior works that showed spatial position is encoded in the early layers of image transformers [2, 24].

In the mid-layers (Figure 7, middle), we find that, among other things, all the models learn to localize and track individual objects. This result introduces a new angle to the recently developed field of object-centric representation learning [3, 23, 43, 55]: it invites us to explore how specialized approaches contribute, given that object concepts naturally emerge in video transformers. In addition, all models, except for synthetically trained TCOW, develop hand tracking concepts, confirming the importance of hands for action recognition from a bottom-up perspective [57, 68].

Finally, in deeper layers we find concepts that build on top of an object-centric representation to capture specific spatiotemporal events. For example, three out of the four
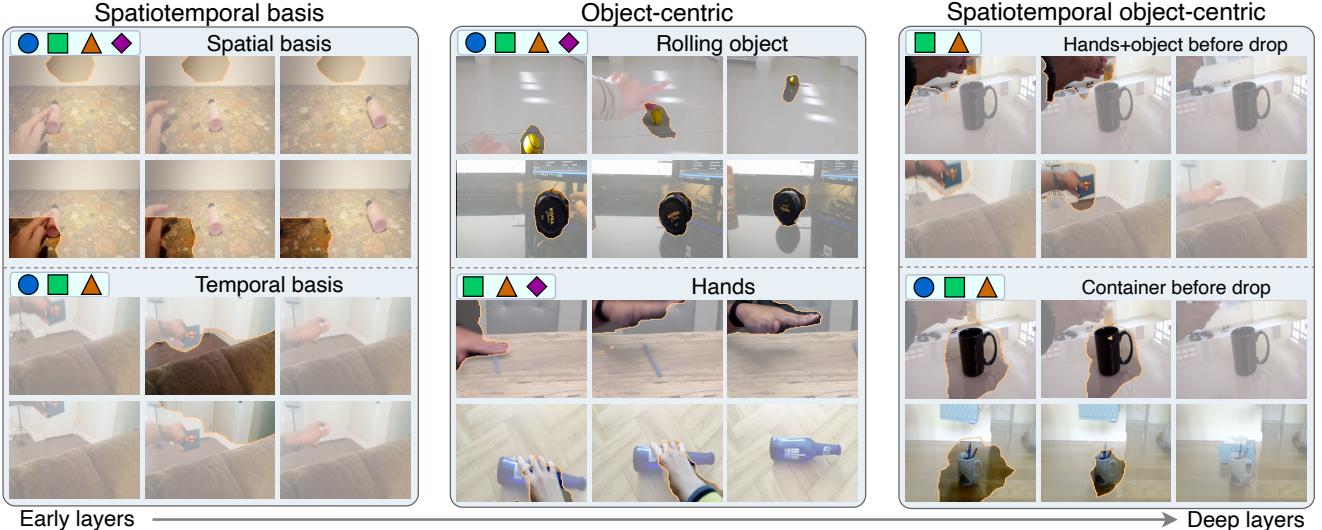
Figure 7. Universal concepts emerge in video transformers despite being trained for different tasks. Early layers encode spatiotemporal positional information. Middle layers track various objects. Deep layers capture fine-grained spatiotemporal concepts, e.g. related to occlusion reasoning (see video for full results). Legend: ● TCOW, ■ Supervised VideoMAE, ▲ SSL VideoMAE, ◆ InternVideo.

| Model | Accuracy ↑ | GFLOPs ↓ |
|---|---|---|
| Baseline | 37.1 | 180.5 |
| VTCD 33% Pruned | **41.4** | 121.5 |
| VTCD 50% Pruned | 37.8 | **91.1** |

Table 1. Pruning unimportant heads in VideoMAE results in improved efficiency and accuracy when targeting a subset of classes. Here, we target the six SSv2 classes containing types of spills.

models learn to identify containers before an object has been dropped into them and two models track the object in the hand until it gets dropped. One notable exception here is InternVideo [66], which is primarily trained on static images and has a limited spatiotemporal modeling capacity. On the other hand, we emphasize that these concepts are also found in the self-supervised VideoMAE [60] that was never trained to reason about object-container relationships. This intriguing observation raises the question: can intuitive physics models [7, 67] be learned via large-scale training of video representations?

## 5.4. Application: Concept pruning with VTCD

Consider a video transformer model trained for $K$ classes, but during deployment, the user is only interested in a subset of classes, $K_d \subset K$. Typically, they can either utilize the model as is, with many parameters dedicated to classes outside $K_d$, or enhance performance by finetuning the model solely on the subset. Motivated by the findings that different self-attention heads focus on class-specific concepts (Section 5.1), we instead propose a simple method leveraging VTCD to prune unimportant heads to $K_d$ and improve model performance and efficiency without finetuning.

We choose the VideoMAE model trained on the SSv2 dataset and focus on the six classes where a 'spill' takes place (listed in the appendix). We then use CRIS to rank all the heads in VideoMAE according to their effect on performance on those six classes using the training set and report the results for pruning the least important of them on the validation set in Table 1. Pruning the 33% least important heads actually *improves* the accuracy by 4.3% while reducing FLOPS from 180.5 to 121.5. Further removing 50% of the heads retains the full performance of the original model (+0.7%) and reduces FLOPs to 91.1. These results suggest that one can control the trade-off between performance and computation by choosing the number of heads to remove.

## 6. Conclusion

In this work, we introduced VTCD, the first algorithm for concept discovery in video transformers. We experimentally demonstrated that it is capable of extracting human-interpretable concepts from video understanding models and quantifying their importance for the final predictions. Using VTCD, we discovered shared concepts among several models with varying objectives, revealing common processing patterns like a spatiotemporal basis in early layers. In later layers, useful, higher-level representations universally emerge, such as those responsible for object tracking. Large-scale video representation learning is an active area of research at the moment and our approach can serve as a key to unlocking its full potential.

# Understanding Video Transformers via Universal Concept Discovery

# Appendix

In this appendix, we report additional results, visualizations and implementation details. Note that we include additional video results and corresponding discussions on the project web page. We begin by ablating the tubelet generation component of VTCD and comparing to an existing concept discovery approach in Section 7.1. We then provide statistics of concepts importance distribution between layers in Section 7.2. Next, in Section 7.3, we provide further discussion and qualitative results showing how different concepts are captured in different self-attention heads from the same layer. Finally, we provide further implementation details in Section 8.

## 7. Additional results

### 7.1. Tubelet validation

Recall that, unlike previous methods that partition the inputs into proposals in the pixel space, VTCD generates the proposals via SLIC [1] clustering in the model's feature space. We ablate this design choice by comparing VTCD with CRAFT [22] - a recent concept discovery approach that uses random cropping for proposal generation, in Table 2.

In particular, we report concept attribution results for the TCOW [62] and VideoMAE [60] models for both our method and CRAFT. In all cases, VTCD result in concepts that are more faithful to the model's representation. To further isolate the effect of proposals on the performance, we then equip CRAFT with our concept importance estimation approach (shown as 'CRAFT [22] + CRIS' in the table). The results confirm our observation from Figure 4 in the main paper that CRIS is superior to the occlusion-based masking used in CRAFT. However, VTCD still outperforms this strong baseline in all settings, validating that generating tubelet proposals in the feature space of the transformer indeed results in concepts that are more faithful to the model's representation.

### 7.2. Quantitative analysis of per-layer concept importance

We now quantify the importance of each model layer for the two target models analyzed in Section 5.2 in the main paper. To this end, we calculate the average concept importance ranking per-layer and then normalize this value, which results in a $[0 - 1]$ score, where higher values indicate more important layers, and plot the results in Figure 8.

We immediately see similarities and differences between the two models. For example, the first two layers are less important than mid layers for both models. For Video-

| Model | TCOW | | VideoMAE | |
| | Positive ↓ | Negative ↑ | Positive ↓ | Negative ↑ |
|---|---|---|---|---|
| CRAFT [22] | 0.174 | 0.274 | 0.240 | 0.300 |
| CRAFT [22] + CRIS | 0.166 | 0.284 | 0.157 | 0.607 |
| VTCD (Ours) | **0.102** | **0.288** | **0.094** | **0.625** |

Table 2. Ablation of our tubelet proposal approach via comparison to CRAFT [22] for both TCOW [62] and VideoMAE [60]. Our tubelets result in concepts that are more faithful to the model's representations even when the baseline is equipped with our concept scoring algorithm (CRAFT [22] + CRIS).
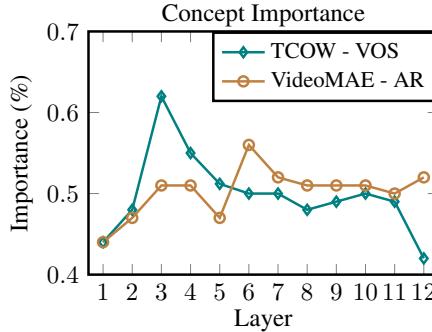
Figure 8. The average concept importance over all model layers for a VOS model (TCOW) and an action recognition model (VideoMAE). Interestingly, while VideoMAE encodes important concepts both in the middle and late in the model, TCOW encodes most important concepts at layer three and the least important in the final layer.

MAE, the middle (6) and end layer (12) are the most important. Interestingly, for TCOW, the most important layer by far is layer 3, while the final layer is the least important. This makes intuitive sense since TCOW is an object tracking model, hence it most heavily utilizes spatiotemporal positional information and object-centric representations in early-to-mid layers. In contrast, VideoMAE is trained for action classification, which requires fine-grained, spatiotemporal concepts in the last layers.

### 7.3. Uniqueness of head concepts

As discussed in Section 4 in the main paper, we qualitatively visualize the concepts from the same layer but different heads of a model to demonstrate that the heads encode diverse concepts. For example, Figure 9 shows that discovered concepts in heads one and six in layer five of the TCOW [62] model encode unrelated concepts (*e.g.* positional and falling objects). This corroborates existing work [2, 17, 34, 46, 63] that heads capture independent information and are therefore a necessary unit of study using VTCD.
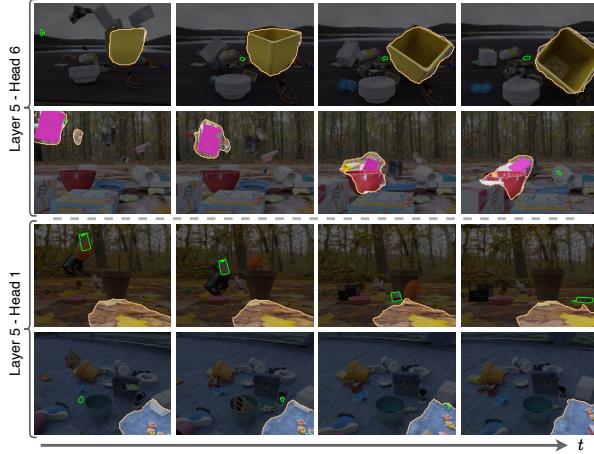
Figure 9. Different heads in the same layer capture independent concepts. In layer 5 of TCOW [62], head 6 (top two rows) highlights falling objects, while head 1 (bottom two rows) captures spatial position.

## 8. Implementation details

**Concept discovery.** When generating tubelets (Section 3.1.1), we use 12 segments and set all other hyperparameters to the Scikit-Image [61] defaults, except for the compactness parameter, which is tuned on a held-out set for each model to the following values: TCOW - 0.01, Video-MAE - 0.1, SSL-VideoMAE - 0.1, InternVideo - 0.15. When clustering concepts using CNMF (Section 3.1.2) we follow the same protocol as [2] and use the Elbow method, with the Silhouette metric [53] as the distance, to select the number of clusters with a threshold of 0.9.

**Concept importance.** For all importance rankings using CRIS, we use the original loss the models were trained with. For InternVideo, we use logit value for the target class by encoding the class name with the text encoder, and then taking the dot product between the text and video features. We use 4,000 masking iterations for all models, except for TCOW [62], where we empirically observe longer convergence times and use 8,000 masks.

**Concept pruning with VTCD.** The six classes targeted in the concept pruning application (Section 5.4 in the main paper) are as follows:

1. Pouring something into something until it overflows
2. Spilling something behind something
3. Spilling something next to something
4. Spilling something onto something
5. Tipping something with something in it over, so something in it falls out
6. Trying to pour something into something, but missing so it spills next to it

2

# References

[1] Radhakrishna Achanta, Appu Shaji, Kevin Smith, Aurelien Lucchi, Pascal Fua, and Sabine Süsstrunk. SLIC superpixels compared to state-of-the-art superpixel methods. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(11):2274–2282, 2012. 2, 3, 1

[2] Shir Amir, Yossi Gandelsman, Shai Bagon, and Tali Dekel. Deep VIT features as dense visual descriptors. In *ECCV Workshops*, 2022. 2, 5, 6, 7, 1

[3] Zhipeng Bao, Pavel Tokmakov, Allan Jabri, Yu-Xiong Wang, Adrien Gaidon, and Martial Hebert. Discovering objects that can move. In *CVPR*, 2022. 7

[4] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *CVPR*, 2017. 1, 2

[5] Joy Buolamwini and Timnit Gebru. Gender shades: Intersectional accuracy disparities in commercial gender classification. In *ACM FAccT*, 2018. 1

[6] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? A new model and the kinetics dataset. In *CVPR*, 2017. 5

[7] Michael Chang, Tomer Ullman, Antonio Torralba, and Joshua Tenenbaum. A compositional object-based approach to learning physical dynamics. In *ICLR*, 2016. 8

[8] Hila Chefer, Shir Gur, and Lior Wolf. Generic attention-model explainability for interpreting bi-modal and encoder-decoder transformers. In *ICCV*, 2021. 2

[9] Hila Chefer, Shir Gur, and Lior Wolf. Transformer interpretability beyond attention visualization. In *CVPR*, 2021. 2

[10] Jinwoo Choi, Chen Gao, C. E. Joseph Messou, and Jia-Bin Huang. Why can't I dance in the mall? Learning to mitigate scene bias in action recognition. In *NeurIPS*, 2019. 3

[11] European Commision. Laying down harmonised rules on artificial intelligence (artificial intelligence act) and amending certain union legislative acts. *European Commision*, 2021. 1

[12] Jean-Baptiste Cordonnier, Andreas Loukas, and Martin Jaggi. On the relationship between self-attention and convolutional layers. In *ICLR*, 2020. 2

[13] Timothée Darcet, Maxime Oquab, Julien Mairal, and Piotr Bojanowski. Vision transformers need registers. *arXiv preprint arXiv:2309.16588*, 2023. 1

[14] Chris HQ Ding, Tao Li, and Michael I Jordan. Convex and semi-nonnegative matrix factorizations. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 32(1):45–55, 2008. 2, 4

[15] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. In *ICLR*, 2021. 5

[16] Amil Dravid, Yossi Gandelsman, Alexei A Efros, and Assaf Shocher. Rosetta neurons: Mining the common units in a model zoo. In *ICCV*, 2023. 2, 5

[17] Nelson Elhage, Neel Nanda, Catherine Olsson, Tom Henighan, Nicholas Joseph, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Nova DasSarma, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. A mathematical framework for transformer circuits. *Transformer Circuits Thread*, 2021. https://transformer-circuits.pub/2021/framework/index.html. 2, 5, 1

[18] Nelson Elhage, Tristan Hume, Catherine Olsson, Nicholas Schiefer, Tom Henighan, Shauna Kravec, Zac Hatfield-Dodds, Robert Lasenby, Dawn Drain, Carol Chen, Roger Grosse, Sam McCandlish, Jared Kaplan, Dario Amodei, Martin Wattenberg, and Christopher Olah. Toy models of superposition. *Transformer Circuits Thread*, 2022. https://transformer-circuits.pub/2022/toy_model/index.html. 2

[19] Christoph Feichtenhofer, Axel Pinz, Richard P Wildes, and Andrew Zisserman. Deep insights into convolutional networks for video recognition. *International Journal of Computer Vision*, 128:420–437, 2020. 3

[20] Christoph Feichtenhofer, Yanghao Li, Kaiming He, et al. Masked autoencoders as spatiotemporal learners. *NeurIPS*, 2022. 3

[21] Thomas Fel, Victor Boutin, Mazda Moayeri, Rémi Cadène, Louis Bethune, Mathieu Chalvidal, Thomas Serre, et al. A holistic approach to unifying automatic concept extraction and concept importance estimation. *NeurIPS*, 2023. 2, 3, 4, 6

[22] Thomas Fel, Agustin Picard, Louis Bethune, Thibaut Boissin, David Vigouroux, Julien Colin, Rémi Cadène, and Thomas Serre. Craft: Concept recursive activation factorization for explainability. In *CVPR*, 2023. 1, 2, 3, 4, 6

[23] Ruohan Gao, Dinesh Jayaraman, and Kristen Grauman. Object-centric representation learning from unlabeled videos. In *ACCV*, 2017. 7

[24] Amin Ghiasi, Hamid Kazemi, Eitan Borgnia, Steven Reich, Manli Shu, Micah Goldblum, Andrew Gordon Wilson, and Tom Goldstein. What do vision transformers learn? a visual exploration. *arXiv preprint arXiv:2212.06727*, 2022. 1, 6, 7

[25] Amir Ghodrati, Efstratios Gavves, and Cees G. M. Snoek. Video time: Properties, encoders and evaluation. In *BMVC*, 2018. 3

[26] Amirata Ghorbani, James Wexler, James Y Zou, and Been Kim. Towards automatic concept-based explanations. In *NeurIPS*, 2019. 1, 2, 3, 6

[27] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, Florian Hoppe, Christian Thurau, Ingo Bax, and Roland Memisevic. The "something something" video database for learning and evaluating visual common sense. In *ICCV*, 2017. 5

[28] Klaus Greff, Francois Belletti, Lucas Beyer, Carl Doersch, Yilun Du, Daniel Duckworth, David J Fleet, Dan Gnanapragasam, Florian Golemo, Charles Herrmann, et al. Kubric: A scalable dataset generator. In *CVPR*, 2022. 5

[29] Isma Hadji and Richard P Wildes. A new large scale dynamic

texture dataset with application to convnet understanding. In *ECCV*, 2018. 3

[30] Sven Ove Hansson, Matts-Åke Belin, and Björn Lundgren. Self-driving vehicles-An ethical overview. *Philosophy & Technology*, pages 1–26, 2021. 1

[31] The White House. President biden issues executive order on safe, secure, and trustworthy artificial intelligence. *The White House*, 2023. 1

[32] Filip Ilic and Axel Pinz. Representing objects in video as space-time volumes by combining top-down and bottom-up processes. In *WACV*, 2020. 3

[33] Filip Ilic, Thomas Pock, and Richard P Wildes. Is appearance free action recognition possible? In *ECCV*, 2022. 3

[34] Rezaul Karim, He Zhao, Richard P. Wildes, and Mennatullah Siam. MED-VT: Multiscale encoder-decoder video transformer with application to object segmentation. In *CVPR*, 2023. 2, 5, 1

[35] Been Kim, Martin Wattenberg, Justin Gilmer, Carrie Cai, James Wexler, and Fernanda Viegas. Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (TCAV). In *ICML*, 2018. 1, 2, 6

[36] Yu Kong and Yun Fu. Human action recognition and prediction: A survey. *International Journal of Computer Vision*, 130(5):1366–1401, 2022. 3

[37] Matthew Kowal, Mennatullah Siam, Md Amirul Islam, Neil DB Bruce, Richard P Wildes, and Konstantinos G Derpanis. A deeper dive into what deep spatiotemporal networks encode: Quantifying static vs. dynamic information. In *CVPR*, 2022. 3

[38] Matthew Kowal, Mennatullah Siam, Md Amirul Islam, Neil DB Bruce, Richard P Wildes, and Konstantinos G Derpanis. Quantifying and learning static vs. dynamic information in deep spatiotemporal networks. *arXiv preprint arXiv:2211.01783*, 2022. 3

[39] Daniel D Lee and H Sebastian Seung. Learning the parts of objects by non-negative matrix factorization. *Nature*, 401 (6755):788–791, 1999. 2

[40] Fuxin Li, Taeyoung Kim, Ahmad Humayun, David Tsai, and James M Rehg. Video segmentation by tracking many figure-ground segments. In *ICCV*, 2013. 3

[41] Yi Li and Nuno Vasconcelos. Repair: Removing representation bias by dataset resampling. In *CVPR*, 2019. 3

[42] Yingwei Li, Yi Li, and Nuno Vasconcelos. Resound: Towards action recognition without representation bias. In *ECCV*, 2018. 3

[43] Francesco Locatello, Dirk Weissenborn, Thomas Unterthiner, Aravindh Mahendran, Georg Heigold, Jakob Uszkoreit, Alexey Dosovitskiy, and Thomas Kipf. Object-centric learning with slot attention. *NeurIPS*, 2020. 7

[44] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *NeurIPS*, 2019. 3, 4

[45] Chris Olah, Nick Cammarata, Ludwig Schubert, Gabriel Goh, Michael Petrov, and Shan Carter. Zoom in: An introduction to circuits. *Distill*, 2020. https://distill.pub/2020/circuits/zoom-in. 1

[46] Catherine Olsson, Nelson Elhage, Neel Nanda, Nicholas Joseph, Nova DasSarma, Tom Henighan, Ben Mann, Amanda Askell, Yuntao Bai, Anna Chen, Tom Conerly, Dawn Drain, Deep Ganguli, Zac Hatfield-Dodds, Danny Hernandez, Scott Johnston, Andy Jones, Jackson Kernion, Liane Lovitt, Kamal Ndousse, Dario Amodei, Tom Brown, Jack Clark, Jared Kaplan, Sam McCandlish, and Chris Olah. In-context learning and induction heads. *Transformer Circuits Thread*, 2022. https://transformer-circuits.pub/2022/in-context-learning-and-induction-heads/index.html. 5, 1

[47] Namuk Park, Wonjae Kim, Byeongho Heo, Taekyung Kim, and Sangdoo Yun. What do self-supervised vision transformers learn? In *ICLR*, 2023. 2

[48] Federico Perazzi, Jordi Pont-Tuset, Brian McWilliams, Luc Van Gool, Markus Gross, and Alexander Sorkine-Hornung. A benchmark dataset and evaluation methodology for video object segmentation. In *CVPR*, 2016. 3

[49] Vitali Petsiuk, Abir Das, and Kate Saenko. RISE: Randomized input sampling for explanation of black-box models. In *BMVC*, 2018. 2

[50] Yao Qin, Chiyuan Zhang, Ting Chen, Balaji Lakshminarayanan, Alex Beutel, and Xuezhi Wang. Understanding and improving robustness of vision transformers through patch-based negative augmentation. *NeurIPS*, 2022. 2

[51] Maithra Raghu, Thomas Unterthiner, Simon Kornblith, Chiyuan Zhang, and Alexey Dosovitskiy. Do vision transformers see like convolutional neural networks? *NeurIPS*, 2021. 1, 2

[52] Kanchana Ranasinghe, Muzammal Naseer, Salman Khan, Fahad Shahbaz Khan, and Michael S Ryoo. Self-supervised video transformer. In *CVPR*, 2022. 3

[53] Peter J Rousseeuw. Silhouettes: A graphical aid to the interpretation and validation of cluster analysis. *Journal of Computational and Applied Mathematics*, 20:53–65, 1987. 2

[54] Christoph Schuhmann, Richard Vencu, Romain Beaumont, Robert Kaczmarczyk, Clayton Mullis, Aarush Katta, Theo Coombes, Jenia Jitsev, and Aran Komatsuzaki. Laion-400m: Open dataset of clip-filtered 400 million image-text pairs. *arXiv preprint arXiv:2111.02114*, 2021. 5

[55] Maximilian Seitzer, Max Horn, Andrii Zadaianchuk, Dominik Zietlow, Tianjun Xiao, Carl-Johann Simon-Gabriel, Tong He, Zheng Zhang, Bernhard Schölkopf, Thomas Brox, et al. Bridging the gap to real-world object-centric learning. In *ICLR*, 2023. 7

[56] Khurram Soomro, Amir Roshan Zamir, and Mubarak Shah. UCF101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012. 5

[57] Chen Sun, Abhinav Shrivastava, Carl Vondrick, Kevin Murphy, Rahul Sukthankar, and Cordelia Schmid. Actor-centric relation network. In *ECCV*, 2018. 7

[58] Xinyu Sun, Peihao Chen, Liangwei Chen, Changhao Li, Thomas H Li, Mingkui Tan, and Chuang Gan. Masked motion encoding for self-supervised video representation learning. In *CVPR*, 2023. 3

[59] Pavel Tokmakov, Jie Li, and Adrien Gaidon. Breaking the" object" in video object segmentation. In *CVPR*, 2023. 3

[60] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked autoencoders are data-efficient learners

for self-supervised video pre-training. *NeurIPS*, 2022. 3, 5, 8, 1

[61] Stefan Van der Walt, Johannes L Schönberger, Juan Nunez-Iglesias, François Boulogne, Joshua D Warner, Neil Yager, Emmanuelle Gouillart, and Tony Yu. scikit-image: image processing in python. *PeerJ*, 2:e453, 2014. 2

[62] Basile Van Hoorick, Pavel Tokmakov, Simon Stent, Jie Li, and Carl Vondrick. Tracking through containers and occluders in the wild. In *CVPR*, 2023. 1, 2, 3, 5

[63] Elena Voita, David Talbot, Fedor Moiseev, Rico Sennrich, and Ivan Titov. Analyzing multi-head self-attention: Specialized heads do the heavy lifting, the rest can be pruned. In *ACL*, 2019. 2, 3, 4, 5, 1

[64] Matthew Walmer, Saksham Suri, Kamal Gupta, and Abhinav Shrivastava. Teaching matters: Investigating the role of supervision in vision transformers. In *CVPR*, 2023. 1, 2

[65] Limin Wang, Bingkun Huang, Zhiyu Zhao, Zhan Tong, Yinan He, Yi Wang, Yali Wang, and Yu Qiao. VideoMAE v2: Scaling video masked autoencoders with dual masking. In *CVPR*, 2023. 3

[66] Yi Wang, Kunchang Li, Yizhuo Li, Yinan He, Bingkun Huang, Zhiyu Zhao, Hongjie Zhang, Jilan Xu, Yi Liu, Zun Wang, et al. InternVideo: General video foundation models via generative and discriminative learning. *arXiv preprint arXiv:2212.03191*, 2022. 5, 8

[67] Jiajun Wu, Ilker Yildirim, Joseph J Lim, Bill Freeman, and Josh Tenenbaum. Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. *NeurIPS*, 2015. 8

[68] Chuhan Zhang, Ankush Gupta, and Andrew Zisserman. Helping hands: An object-aware ego-centric video recognition model. In *ICCV*, 2023. 7

[69] Ruihan Zhang, Prashan Madumal, Tim Miller, Krista A Ehinger, and Benjamin IP Rubinstein. Invertible concept-based explanations for cnn models with non-negative concept activation vectors. In *Proceedings of the AAAI Conference on Artificial Intelligence*, 2021. 2, 4, 6

[70] Daquan Zhou, Zhiding Yu, Enze Xie, Chaowei Xiao, Animashree Anandkumar, Jiashi Feng, and Jose M Alvarez. Understanding the robustness in vision transformers. In *ICML*, 2022. 2