

# Learning To Recognize Procedural Activities with Distant Supervision

Xudong Lin<sup>1\*</sup> Fabio Petroni<sup>2</sup> Gedas Bertasius<sup>3</sup>

Marcus Rohrbach<sup>2</sup> Shih-Fu Chang<sup>1</sup> Lorenzo Torresani<sup>2,4</sup>

<sup>1</sup>Columbia University <sup>2</sup>Facebook AI Research <sup>3</sup>UNC Chapel Hill <sup>4</sup>Dartmouth

## Abstract

*In this paper we consider the problem of classifying fine-grained, multi-step activities (e.g., cooking different recipes, making disparate home improvements, creating various forms of arts and crafts) from long videos spanning up to several minutes. Accurately categorizing these activities requires not only recognizing the individual steps that compose the task but also capturing their temporal dependencies. This problem is dramatically different from traditional action classification, where models are typically optimized on videos that span only a few seconds and that are manually trimmed to contain simple atomic actions. While step annotations could enable the training of models to recognize the individual steps of procedural activities, existing large-scale datasets in this area do not include such segment labels due to the prohibitive cost of manually annotating temporal boundaries in long videos. To address this issue, we propose to automatically identify steps in instructional videos by leveraging the distant supervision of a textual knowledge base (wikiHow) that includes detailed descriptions of the steps needed for the execution of a wide variety of complex activities. Our method uses a language model to match noisy, automatically-transcribed speech from the video to step descriptions in the knowledge base. We demonstrate that video models trained to recognize these automatically-labeled steps (without manual supervision) yield a representation that achieves superior generalization performance on four downstream tasks: recognition of procedural activities, step classification, step forecasting and egocentric video classification.*

## 1. Introduction

Imagine being in your kitchen, engaged in the preparation of a sophisticated dish that involves a sequence of complex steps. Fortunately, your J.A.R.V.I.S.<sup>1</sup> comes to your rescue. It actively recognizes the task that you are trying to accomplish and guides you step-by-step in the successful

execution of the recipe. The dramatic progress witnessed in activity recognition [9, 11, 51, 53] over the last few years has certainly made these fictional scenarios a bit closer to reality. Yet, it is clear that in order to attain these goals we must extend existing systems beyond atomic-action classification in trimmed clips to tackle the more challenging problem of understanding procedural activities in long videos spanning several minutes. Furthermore, in order to classify the procedural activity, the system must not only recognize the individual semantic steps in the long video but also model their temporal relations, since many complex activities share several steps but may differ in the order in which these steps appear or are interleaved. For example, “beating eggs” is a common step in many recipes which, however, are likely to differ in the preceding and subsequent steps.

In recent years, the research community has engaged in the creation of several manually-annotated video datasets for the recognition of procedural, multi-step activities. However, in order to make detailed manual annotations possible at the level of both segments (step labels) and videos (task labels), these datasets have been constrained to have a narrow scope or a relatively small scale. Examples include video benchmarks that focus on specific domains, such as recipe preparation or kitchen activities [11, 27, 40, 62], as well as collections of instructional videos manually-labeled for step and task recognition [50, 63]. Due to the large cost of manually annotating temporal boundaries, these datasets have been limited to a small size both in terms of number of tasks (about a few hundreds activities at most) as well as amount of video examples (about 10K samples, for roughly 400 hours of video). While these benchmarks have driven early progress in this field, their limited size and narrow scope prevent the training of modern large-capacity video models for recognition of general procedural activities.

On the other end of the scale/scope spectrum, the HowTo100M dataset [34] stands out as an exceptional resource. It is over 3 orders of magnitude bigger than prior benchmarks in this area along several dimensions: it includes over 100M clips showing humans performing and narrating more than 23,000 complex tasks for a total duration of 134K hours of video. The downside of this massive

\*Research done while XL was an intern at Facebook AI Research.

<sup>1</sup>A fictional AI assistant in the Marvel Cinematic Universe.

amount of data is that its scale effectively prevents manual annotation. In fact, all videos in HowTo100M are unverified by human annotators. While this benchmark clearly fulfills the size and scope requirements needed to train large-capacity video models, its lack of segment annotations and the unvalidated nature of the videos impedes the training of accurate step or task classifiers.

In this paper we present a novel approach for training models to recognize procedural steps in instructional video *without* any form of manual annotation, thus enabling optimization on large-scale unlabeled datasets, such as HowTo100M. We propose a *distant supervision* framework that leverages a textual knowledge base as a guidance to automatically identify segments corresponding to different procedural steps in video. Distant supervision has been used in Natural Language Processing [35, 39, 42] to mine relational examples from noisy text corpora using a knowledge base. In our setting, we are also aiming at relation extraction, albeit in the specific setting of identifying video segments relating to semantic steps. The knowledge base that we use is wikiHow [2]—a crowdsourced multimedia repository containing over 230,000 “how-to” articles describing and illustrating steps, tips, warnings and requirements to accomplish a wide variety of tasks. Our system uses language models to compare segments of narration automatically transcribed from the videos to the textual descriptions of steps in wikiHow. The matched step descriptions serve as distant supervision to train a video understanding model to learn step-level representations. Thus, our system uses the knowledge base to mine step examples from the noisy, large-scale unlabeled video dataset. To the best of our knowledge, this is the first attempt at learning a step video representation with distant supervision.

We demonstrate that video models trained to recognize these pseudo-labeled steps in a massive corpus of instructional videos, provide a general video representation transferring effectively to four different downstream tasks on new datasets. Specifically, we show that we can apply our model to represent a long video as a sequence of step embeddings extracted from the individual segments. Then, a shallow sequence model (a single Transformer layer [52]) is trained on top of this sequence of embeddings to perform temporal reasoning over the step embeddings. Our experiments show that such an approach yields state-of-the-art results for classification of procedural tasks on the labeled COIN dataset, outperforming the best reported numbers in the literature by more than 16%. Furthermore, we use this method to make additional insightful observations:

1. Step labels assigned with our distant supervision framework yield better downstream results than those obtained by using the unverified task labels of HowTo100M.
2. Our distantly-supervised video representation outper-

forms *fully-supervised* video features trained with action labels on the large-scale Kinetics-400 dataset [9].

3. Our step assignment procedure produces better downstream results than a representation learned by directly matching video to the ASR narration [33], thus showing the value of the distant supervision framework.

We also evaluate the performance of our system for classification of procedural activities on the Breakfast dataset [27]. Furthermore, we present transfer learning results on three additional downstream tasks on datasets different from that used to learn our representation (HowTo100M): step classification and step forecasting on COIN, as well as categorization of egocentric videos on EPIC-KITCHENS-100 [10]. On all of these tasks, our distantly-supervised representation achieves higher accuracy than previous works, as well as additional baselines that we implement based on training with full supervision. These results provide further evidence of the generality and effectiveness of our unsupervised representation for understanding complex procedural activities in videos. We will release the code and the automatic annotations provided by our distant supervision<sup>2</sup>.

## 2. Related Work

During the past decade, we have witnessed dramatic progress in action recognition. However, the benchmarks in this field consist of brief videos (usually, a few seconds long) trimmed to contain the individual atomic action to recognize [19, 24, 28, 44]. In this work, we consider the more realistic setting where videos are untrimmed, last several minutes, and contain sequences of steps defining the complex procedural activities to recognize (e.g., a specific recipe, or a particular home improvement task).

**Understanding Procedural Videos.** Procedural knowledge is an important part of human knowledge [4, 37, 48] essentially answering “how-to” questions. Such knowledge is displayed in long procedural videos [11, 27, 34, 41, 50, 62, 63], which have attracted active research in recognition of multi-step activities [21, 23, 61]. Early benchmarks in this field contained manual annotations of steps within the videos [50, 62, 63] but were relatively small in scope and size. The HowTo100M dataset [34], on the other hand, does not contain any manual annotations but it is several orders of magnitude bigger and the scope of its “how-to” videos is very broad. An instructional or how-to video contains a human subject demonstrating and narrating how to accomplish a certain task. Early works on HowTo100M have focused on leveraging this large collection for learning models that can be transferred to other tasks, such as action recognition [3, 33, 34], video captioning [20, 32, 62], or text-video retrieval [6, 33, 56]. The problem of recognizing the task

<sup>2</sup>Please check <https://arxiv.org/abs/2201.10990> for updates.

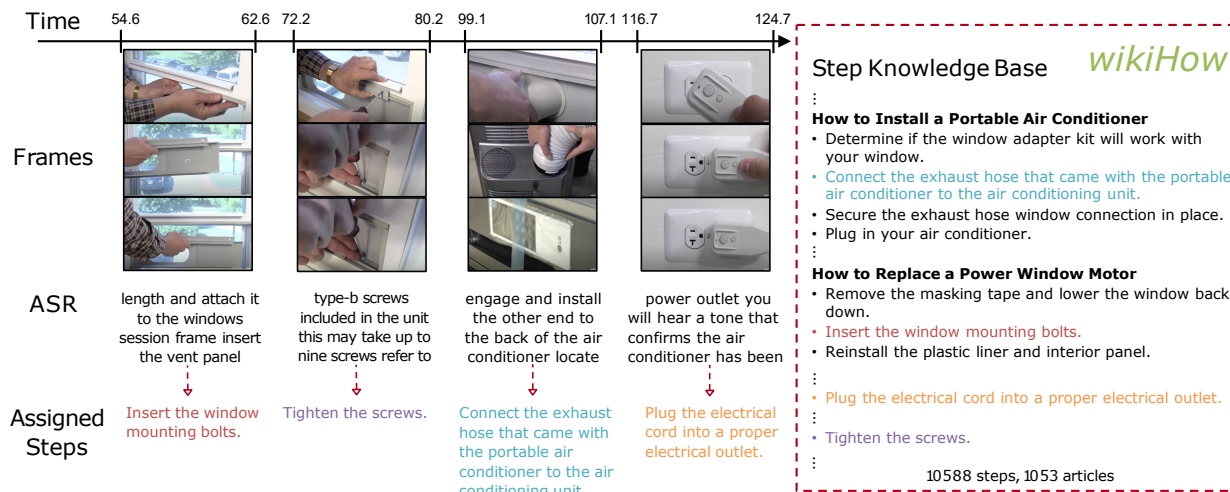


Figure 1. Illustration of our proposed framework. Given a long instructional video as input, our method generates distant supervision by matching segments in the video to steps described in a knowledge base (wikiHow). The matching is done by comparing the automatically-transcribed narration to step descriptions using a pretrained language model. This distant supervision is then used to learn a video representation recognizing these automatically annotated steps. This video is from the HowTo100M dataset. More examples are provided in the supplementary material.

performed in the instructional video has been considered by Bertasius *et al.* [7]. However, their proposed approach does not model the procedural nature of instructional videos.

**Learning Video Representations with Limited Supervision.** Learning semantic video representations [30, 36, 45, 46, 58] is a fundamental problem in video understanding research. The representations pretrained from labeled datasets are limited by the pretraining domain and the predefined ontology. Therefore, many attempts have been made to obtain video representations with less human supervision. In the unsupervised setting, supervision signal is usually constructed by augmenting videos [16, 45, 55]. For example, Wei *et al.* [55] proposed to predict the order of videos as the supervision to learn order-aware video representations. In the weakly supervised setting, the supervision signals are usually obtained from hashtags [18], ASR transcriptions [34], or meta-information extracted from the Web [17]. Miech *et al.* [34] show that ASR sentences extracted from audio can serve as a valuable information source to learn video representations. Previous works [13, 14] have also studied learning to localize keyframes using task labels as supervision. This is different from the focus of this paper, which addresses the problem of learning step-level representations from unlabeled instructional videos.

**Distant Supervision.** Distant supervision [35, 59] has been studied in natural language processing and generally refers to a training scheme where supervision is obtained by automatically mining examples from a large noisy corpus utilizing a clean and informative knowledge base. It has been shown to be very successful on the problem of relation ex-

traction. For example, Mintz *et al.* [35] leverage knowledge from Freebase [8] to obtain supervision for relation extraction. However, the concept of distant supervision has not been exploited in video understanding. Huang *et al.* [20] have proposed to use wikiHow as a textual dataset to pre-train a video captioning model but the knowledge base is not used to supervise video understanding models.

### 3. Technical Approach

Our goal is to learn a segment-level representation to express a long procedural video as a sequence of step embeddings. The application of a sequence model, such as a Transformer, on this video representation can then be used to perform temporal reasoning over the individual steps. Most importantly, we want to learn the step-level representation without manual annotations, so as to enable training on large-scale unlabeled data. The key insight leveraged by our framework is that knowledge bases, such as wikiHow, provide detailed textual descriptions of the steps for a wide range of tasks. In this section, we will first describe how to obtain distant supervision from wikiHow, then discuss how the distant supervision can be used for step-level representation learning, and finally, we will introduce how our step-level representation is leveraged to solve several downstream problems.

#### 3.1. Extracting Distant Supervision from wikiHow

The wikiHow repository contains high-quality articles describing the sequence of individual steps needed for the completion of a wide variety of practical tasks. For-

mally, we refer to wikiHow as a knowledge base  $\mathbb{B}$  containing textual step descriptions for  $T$  tasks:  $\mathbb{B} = \{y_1^{(1)}, \dots, y_{S_1}^{(1)}, \dots, y_1^{(T)}, \dots, y_{S_T}^{(T)}\}$ , where  $y_s^{(t)}$  represents the language-based description of step  $s$  for task  $t$ , and  $S_t$  is the number of steps involved for the execution of task  $t$ . We view an instructional video  $x$  as a sequence of  $L$  segments  $\{x_1, \dots, x_l, \dots, x_L\}$ , with each segment  $x_l$  consisting of  $F$  RGB frames having spatial resolution  $H \times W$ , i.e.,  $x_l \in \mathbb{R}^{H \times W \times 3 \times F}$ . Each video is accompanied by a paired sequence of text sentences  $\{a_1, \dots, a_l, \dots, a_L\}$  obtained by applying ASR to the audio narration. We note that the narration  $a_l$  can be quite noisy due to ASR errors. Furthermore, it may describe the step being executed only implicitly, e.g., by referring to secondary aspects. An example is given in Fig. 1, where the ASR in the second segment describes the type of screws rather than the action of tightening the screws, while the last segment refers to the the tone confirmation of the air conditioner being activated rather than the plugging of the cord into the outlet. The idea of our approach is to leverage the knowledge base  $\mathbb{B}$  to de-noise the narration  $a_l$  and to convert it into a supervisory signal that is more directly related to the steps represented in segments of the video. We achieve this goal through the framework of distant supervision, which we apply to approximate the unknown conditional distribution  $P(y_s^{(t)}|x_l)$  over the steps executed in the video, *without* any form of manual labeling. To approximate this distribution we employ a textual similarity measure  $\mathcal{S}$  between  $y_s^{(t)}$  and  $a_l$ :

$$P(y_s^{(t)}|x_l) \approx \frac{\exp(\mathcal{S}(a_l, y_s^{(t)}))}{\sum_{t,s} \exp(\mathcal{S}(a_l, y_s^{(t)}))}. \quad (1)$$

The textual similarity  $\mathcal{S}$  is computed as a dot product between language embeddings

$$\mathcal{S}(a_l, y_s^{(t)}) = e(a_l)^\top \cdot e(y_s^{(t)}) \quad (2)$$

where  $e(a_l), e(y_s^{(t)}) \in \mathbb{R}^d$  and  $d$  is the dimension of the language embedding space. The underlying intuition of our approach is that, compared to the noisy and unstructured narration  $a_l$ , the distribution  $P(y_s^{(t)}|x_l)$  provides a more salient supervisory signal for training models to recognize individual steps of procedural activities in video. The last row of Fig. 1 shows the steps in the knowledge base having highest conditional probability given the ASR text. We can see that, compared to the ASR narrations, the step sentences provide a more fitting description of the step executed in each segment. Our key insight is that we can leverage modern language models to reassign noisy and imprecise speech transcriptions into the clean and informative step descriptions of our knowledge base. Beyond this qualitative illustration (plus additional ones available in the supplementary material), our experiments provide quantitative evidence of the

benefits of training video models by using  $P(y_s^{(t)}|x_l)$  as supervision as opposed to the raw narration.

### 3.2. Learning Step Embeddings from Unlabeled Video

We use the approximated distribution  $P(y_s^{(t)}|x_l)$  as the supervision to learn a video representation  $f(x_l) \in \mathbb{R}^d$ . We consider three different training objectives for learning the video representation  $f$ : (1) step classification, (2) distribution matching, and (3) step regression.

**Step Classification.** Under this learning objective, we first train a step classification model  $\mathcal{F}_C : \mathbb{R}^{H \times W \times 3 \times F} \rightarrow [0, 1]^S$  to classify each video segment into one of the  $S$  possible steps in the knowledge base  $\mathbb{B}$ , where  $S = \sum_t S_t$ . Specifically, let  $t^*, s^*$  be the indices of the step in  $\mathbb{B}$  that best describes segment  $x_l$  according to our target distribution, i.e.,

$$t^*, s^* = \arg \max_{t,s} P(y_s^{(t)}|x_l). \quad (3)$$

Then, we use the standard cross-entropy loss to train  $\mathcal{F}_C$  to classify video segment  $x_l$  into class  $(t^*, s^*)$ :

$$\min_{\theta} -\log \left( [\mathcal{F}_C(x_l; \theta)]_{(t^*, s^*)} \right) \quad (4)$$

where  $\theta$  denotes the learning parameters of the video model. The model uses a softmax activation function in the last layer to define a proper distribution over the steps, such that  $\sum_{t,s} [\mathcal{F}_C(x_l; \theta)]_{(t,s)} = 1$ . Although here we show the loss for one segment  $x_l$  only, in practice we optimize the objective by averaging over a mini-batch of video segments sampled from the entire collection in each iteration. After learning, we use  $\mathcal{F}_C(x_l)$  as a feature extractor to capture step-level information from new video segments. Specifically, we use the second-to-last layer of  $\mathcal{F}_C(x_l)$  (before the softmax function) as the step embedding representation  $f(x_l)$  for classification of procedural activities in long videos.

**Distribution Matching.** Under the objective of Distribution Matching, we train the step classification model  $\mathcal{F}_C$  to minimize the KL-Divergence between the predicted distribution  $\mathcal{F}_C(x_l)$  and the target distribution  $P(y_s^{(t)}|x_l)$ :

$$\min_{\theta} \sum_{t,s} P(y_s^{(t)}|x_l) \log \frac{P(y_s^{(t)}|x_l)}{[\mathcal{F}_C(x_l; \theta)]_{(t,s)}}. \quad (5)$$

Due to the large step space ( $S = 10,588$ ), in order to effectively optimize this objective we empirically found it beneficial to use only the top- $K$  steps in  $P(y_s^{(t)}|x_l)$ , with the probabilities of the other steps set to zero.

**Step Regression.** Under Step Regression, we train the video model to predict the language embedding  $e(y_{s^*}^{(t^*)}) \in \mathbb{R}^d$  associated to the pseudo ground-truth step  $(t^*, s^*)$ . Thus, in this case the model is a regression function to the



language embedding space, i.e.,  $\mathcal{F}_R : \mathbb{R}^{H \times W \times 3 \times F} \rightarrow \mathbb{R}^d$ . We follow [33] and use the NCE loss as the objective:

$$\min_{\theta} - \log \frac{\exp \left( e(y_{s^*}^{t^*})^\top \mathcal{F}_R(x_l; \theta) \right)}{\sum_{(t,s) \neq (t^*, s^*)} \exp \left( e(y_s^{(t)})^\top \mathcal{F}_R(x_l; \theta) \right)} \quad (6)$$

Because  $\mathcal{F}_R(x_l)$  is trained to predict the language representation of the step, we can directly use its output as step embedding representation for new video segments, i.e.,  $f(x_l) = \mathcal{F}_R(x_l)$ .

### 3.3. Classification of Procedural Activities

In this subsection we discuss how we can leverage our learned step representation to recognize fine-grained procedural activities in long videos spanning up to several minutes. Let  $x'$  be a new input video consisting of a sequence of  $L'$  segments  $x'_l \in \mathbb{R}^{H \times W \times 3 \times F}$  for  $l = 1, \dots, L'$ . The intuition is that we can leverage our pretrained step representation to describe the video as a sequence of step embeddings. Because our step embeddings are trained to reveal semantic information about the individual steps executed in the segments, we use a transformer [52]  $\mathcal{T}$  to model dependencies over the steps and to classify the procedural activity:  $\mathcal{T}(f(x'_1), \dots, f(x'_{L'}))$ . Since our objective is to demonstrate the effectiveness of our step representation  $f$ , we choose  $\mathcal{T}$  to include a single transformer layer, which is sufficient to model sequential dependencies among the steps and avoids making the classification model overly complex. We refer to this model as the ‘‘Basic Transformer.’’

We also demonstrate that our step embeddings enable further beneficial information transfer from the knowledge base  $\mathbb{B}$  to improve the classification of procedural activities during inference. The idea is to adopt a retrieval approach to find for each segment  $x'_l$  the step  $y_{s'}^{t'} \in \mathbb{B}$  that best explains the segment according to the pretrained video model  $\mathcal{F}(x'_l; \theta)$ . For the case of Step Classification and Distribution Matching, where we learn a classification model  $\mathcal{F}_C(x'_l; \theta) \in [0, 1]^S$ , we simply select the step class yielding the maximum classification score:

$$t', s' = \arg \max_{t,s} [\mathcal{F}_C(x'_l; \theta)]_{(t,s)}. \quad (7)$$

In the case of Step Regression, since  $\mathcal{F}_R(x'_l; \theta)$  generates an output in the language space, we can choose the step that has maximum language embedding similarity:

$$t', s' = \arg \max_{t,s} e(y_s^{(t)})^\top \mathcal{F}_R(x'_l; \theta). \quad (8)$$

Let  $\hat{y}(x'_l)$  denote the step description assigned through this procedure, i.e.,  $\hat{y}(x'_l) = y_{s'}^{t'}$ .

Then, we can incorporate the knowledge retrieved from  $\mathbb{B}$  for each segment in the input provided to the transformer,

together with the step embeddings extracted from the video:

$$\mathcal{T}(f(x'_1), e(\hat{y}(x'_1)), f(x'_2), e(\hat{y}(x'_2)), \dots, f(x'_{L'}), e(\hat{y}(x'_{L'}))). \quad (9)$$

This formulation effectively trains the transformer to fuse a representation consisting of video features and step embeddings from the knowledge base to predict the class of the procedural activity. We refer to this variant as ‘‘Transformer w/ KB Transfer’’.

### 3.4. Step Forecasting

We note that we can easily modify our proposed classification model to address forecasting tasks that require long-term analysis over a sequence of steps to predict future activity. One such problem is the task of ‘‘next-step anticipation’’ which we consider in our experiments. Given as input a video spanning  $M$  segments,  $\{x_1, \dots, x_M\}$ , the objective is to predict the step executed in the *unobserved*  $(M+1)$ -th segment. To address this task we train the transformer on the sequence step embeddings extracted from the  $M$  observed segments. In the case of Transformer w/ KB Transfer, for each input segment  $x'_l$ , we include  $f(x'_l)$  but also  $e(y_{s'+1}^{t'})$ , i.e., the embedding of the step immediately after the step matched in the knowledge base. This effectively provides the transformer with information about the likely future steps according to the knowledge base.

### 3.5. Model Design

We use MPNet [43] as the language model to extract 768-dimensional language embeddings for both the ASR sentences and the step descriptions in wikiHow articles. MPNet (paraphrase-mpnet-base-v2) is at the time of writing (August, 2021) ranked first by Sentence Transformers [1], based on performance across 14 language retrieval tasks [38]. The similarity between two embedding vectors is chosen to be the dot product between the two vectors.

We choose as video model the TimeSformer architecture [7]. Starting from a configuration of ViT initialized with ImageNet-21K ViT pretraining [12], we train TimeSformer on HowTo100M using clips of 8 frames uniformly sampled from time-spans of 8 seconds. The evaluations in our experiments are carried out by learning the step representation on HowTo100M (without manual labels) and by assessing the performance of our embeddings on smaller-scale downstream datasets where task and/or step manual annotations are available. To perform classification of multi-step activities on these downstream datasets we use a single transformer layer [52] trained on top of our fixed embeddings. We use this shallow long-term model without finetuning in order to directly measure the value of the representation learned via distant supervision from the unlabeled instructional videos. We refer the reader to the supplementary material for additional implementation details.

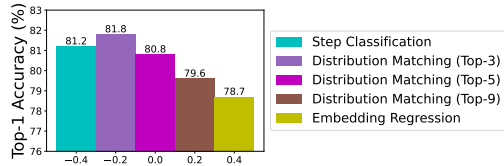


Figure 2. Accuracy of classifying procedural activities in COIN using three different distant supervision objectives.

## 4. Experiments

### 4.1. Datasets and Evaluation Metrics

**Pretraining.** HowTo100M (HT100M) [34] includes over 1M long instructional videos split into about 120M video clips in total. We use the complete HowTo100M dataset only in the final comparison with the state-of-the-art (sec. 4.3). In the ablations, in order to reduce the computational cost, we use a smaller subset corresponding to the collection of 80K long videos defined by Bertasius *et al.* [7].

**Classification of Procedural Activities.** Performance on this task is evaluated using two labeled datasets: COIN [49, 50] and Breakfast [27]. COIN contains about 11K instructional videos representing 180 tasks (i.e., classes of procedural activities). Breakfast [27] contains 1,712 videos for 10 complex cooking tasks. In both datasets, each video is manually annotated with a label denoting the task class. We use the standard splits [21, 50] for these two datasets and measure performance as task classification accuracy.

**Step Classification.** It requires classifying the step observed in a single video segment (without history), which is a good testbed to evaluate the effectiveness of our step embeddings. To evaluate methods on this problem, we use the step annotations from COIN, corresponding to a total of 778 step classes representing parts of tasks. The steps are manually annotated within each video with temporal boundaries and step class labels. Classification accuracy [50] of a linear classifier (Linear Acc) is used as the metric.

**Step Forecasting.** We also use step annotations available in COIN. The objective is to predict the class of the step in the next segment given as input the sequence of observed video segments up to that step (excluded). Note that there is a substantial temporal gap (21 seconds on average) between the end of the last observed segment and the start of the step to be predicted. This makes the problem quite challenging and representative of real-world conditions. We set the history to contain at least one step. We use classification accuracy of the predicted step as the evaluation metric.

**Egocentric Activity Recognition.** EPIC-KITCHENS-100 [10] is a large-scale egocentric video dataset. It consists of 100 hours of first-person videos, showing humans performing a wide range of procedural activities in the kitchen. The dataset includes manual annotations of 97 verbs and 300 nouns in manually-labeled video segments. We follow

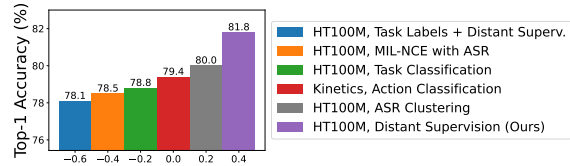


Figure 3. Accuracy of procedural activity classification on COIN using video representations learned with different supervisions.

the standard protocol [10] to train and evaluate our models.

### 4.2. Ablation Studies

We begin by studying how different design choices in our framework affect the accuracy of task classification on COIN using the basic Transformer as our long-term model.

#### 4.2.1 Different Training Objectives

Fig. 2 shows the accuracy of COIN task classification using the three distant supervision objectives presented in Sec. 3.2. Distribution Matching and Step Classification achieve similar performance, while Embedding Regression produces substantially lower accuracy. Based on these results we choose Distribution Matching (Top-3) as our learning objective for all subsequent experiments.

#### 4.2.2 Comparing Different Forms of Supervision

In Fig. 3, we compare the results of different pretrained video representations for the problem of classifying procedural activities on the COIN dataset. We include as baselines several representations learned on the same subset of HowTo100M as our step embeddings, using the same TimeSformer as video model. MIL-NCE [33] performs contrastive learning between the video and the narration obtained from ASR. The baseline (HT100M, Task Classification) is a representation learned by training TimeSformer as a classifier using as classes the task ids available in HowTo100M. The task ids are defined by the keywords used to find the video on YouTube. The baseline (HT100M, Task Labels + Distant Superv.) uses the task ids to narrow down the potential steps considered by distant supervision (only wikiHow steps corresponding to the task id of the video are considered). We also include a representation obtained by training TimeSformer on the fully-supervised Kinetics-400 dataset [9]. Finally, to show the benefits of distant supervision, we run  $k$ -means clustering on the language embeddings of ASR sentences using the same number of clusters as the steps in wikiHow (i.e.,  $k = S = 10,588$ ), and then train the video model using the cluster ids as supervision.

We observe several important results in Fig. 3. First, our distant supervision achieves an accuracy gain of 3.3% over MIL-NCE with ASR. This suggests that our distant supervision framework provides more explicit supervision to

Segment Model	Pretraining Supervision	Pretraining Dataset	Linear Acc (%)
TSN (RGB+Flow) [50]	Supervised: action labels	Kinetics	36.5*
S3D [33]	Unsupervised: MIL-NCE on ASR	HT100M	37.5*
ClipBERT [29]	Supervised: captions	COCO + Visual Genome	30.8
VideoCLIP [57]	Unsupervised: NCE on ASR	HT100M	39.4
SlowFast [15]	Supervised: action labels	Kinetics	32.9
TimeSformer [7]	Supervised: action labels	Kinetics	48.3
TimeSformer [7]	Unsupervised: $k$ -means on ASR	HT100M	46.5
<b>TimeSformer</b>	<b>Unsupervised: distant supervision (ours)</b>	HT100M	<b>54.1</b>

Table 1. Comparison to the state-of-the-art for **step classification** on the COIN dataset. \* indicates results by finetuning on COIN.

Long-term Model	Segment Model	Pretraining Supervision	Pretraining Dataset	Acc (%)
TSN (RGB+Flow) [50]	Inception [47]	Supervised: action labels	Kinetics	73.4*
Basic Transformer	S3D [33]	Unsupervised: MIL-NCE on ASR	HT100M	70.2*
Basic Transformer	ClipBERT [29]	Supervised: captions	COCO + Visual Genome	65.4
Basic Transformer	VideoCLIP [57]	Unsupervised: NCE on ASR	HT100M	72.5
Basic Transformer	SlowFast [15]	Supervised: action labels	Kinetics	71.6
Basic Transformer	TimeSformer [7]	Supervised: action labels	Kinetics	83.5
Basic Transformer	TimeSformer [7]	Unsupervised: $k$ -means on ASR	HT100M	85.3
<b>Basic Transformer</b>	<b>TimeSformer</b>	<b>Unsupervised: distant supervision (ours)</b>	HT100M	<b>88.9</b>
<b>Transformer w/ KB Transfer</b>	<b>TimeSformer</b>	<b>Unsupervised: distant supervision (ours)</b>	HT100M	<b>90.0</b>

Table 2. Comparison to the state-of-the-art for the **classification of procedural activities** on the COIN dataset.

learn step-level representations compared to using directly the ASR text. This is further confirmed by the performance of ASR Clustering, which is 1.7% lower than that obtained by leveraging the wikiHow knowledge base.

Moreover, our step-level representation outperforms by 3% the weakly-supervised task embeddings (Task Classification) and does even better (by 2.4%) than the video representation learned with full supervision from the large-scale Kinetics dataset. This is due to the fact that steps typically involve multiple atomic actions. For example, about 85% of the steps consist of at least two verbs. Thus, our step embeddings capture a higher-level representation than those based on traditional atomic action labels.

Finally, using the task ids to restrict the space of step labels considered by distant supervision produces the worst results. This indicates that the task ids are quite noisy and that our approach leveraging relevant steps from other tasks can provide more informative supervision. These results further confirm the superior performance of distantly supervised step annotations over existing task or action labels to train representations for classifying procedural activities.

### 4.3. Comparisons to the State-of-the-Art

#### 4.3.1 Step Classification

We study the problem of step classification as it directly measures whether the proposed distant supervision framework provides a useful training signal for recognizing steps in video. For this purpose, we use our distantly supervised model as a frozen feature extractor to extract step-level embeddings for each video segment and then train a linear classifier to recognize the step class in the input segment.

Table 1 shows that our distantly supervised representation achieves the best performance and yields a large gain over several strong baselines. Even on this task, our distant supervision produces better results compared to a video representation trained with fully-supervised action labels on Kinetics. The significant gain (7.6%) over ASR clustering again demonstrates the importance of using wikiHow knowledge. Finally, our model achieves strong gains over previously reported results on this benchmark based on different backbones, including results obtained by finetuning and using optical flow as an additional modality [50].

#### 4.3.2 Classification of Procedural Activities

Table 2 and Table 3 show accuracy of classifying procedural activities in long videos on the COIN and Breakfast dataset, respectively. Our model outperforms all previous works on these two benchmarks. For this problem, the accuracy gain on COIN over the representations learned with Kinetics action labels has become even larger (6.5%) compared to the improvement achieved for step classification (5.8%). This indicates that the distantly supervised representation is indeed highly suitable for recognizing long procedural activities. We also observe a substantial gain (8.8%) over the Kinetics baseline for the problem of recognizing complex cooking activities in the Breakfast dataset. As GHRM provided also the result obtained by finetuning the feature extractor on the Breakfast benchmark (89.0%), we measured the accuracy achieved by finetuning our model and observed a large gain: 91.6%. We also tried replacing the basic transformer with Timeception as the long-term model. Timeception trained on features learned with action

Long-term Model	Segment Model	Pretraining Supervision	Pretraining Dataset	Acc (%)
Timeception [21]	3D-ResNet [54]	Supervised: action labels	Kinetics	71.3
VideoGraph [22]	I3D [9]	Supervised: action labels	Kinetics	69.5
GHRM [61]	I3D [9]	Supervised: action labels	Kinetics	75.5
Basic Transformer	S3D [33]	Unsupervised: MIL-NCE on ASR	HT100M	74.4
Basic Transformer	SlowFast [15]	Supervised: action labels	Kinetics	76.1
Basic Transformer	TimeSformer [7]	Supervised: action labels	Kinetics	81.1
Basic Transformer	TimeSformer [7]	Unsupervised: $k$ -means on ASR	HT100M	81.4
<b>Basic Transformer</b>	<b>TimeSformer</b>	<b>Unsupervised: distant supervision (ours)</b>	HT100M	<b>88.7</b>
<b>Transformer w/ KB Transfer</b>	<b>TimeSformer</b>	<b>Unsupervised: distant supervision (ours)</b>	HT100M	<b>89.9</b>

Table 3. Comparison to the state-of-the-art for the problem of classifying procedural activities on the Breakfast dataset.

Long-term Model	Segment Model	Pretraining Supervision	Pretraining Dataset	Acc (%)
Basic Transformer	S3D [33]	Unsupervised: MIL-NCE on ASR	HT100M	28.1
Basic Transformer	SlowFast [15]	Supervised: action labels	Kinetics	25.6
Basic Transformer	TimeSformer [7]	Supervised: action labels	Kinetics	34.7
Basic Transformer	TimeSformer [7]	Unsupervised: $k$ -means on ASR	HT100M	34.0
<b>Basic Transformer</b>	<b>TimeSformer</b>	<b>Unsupervised: distant supervision (ours)</b>	HT100M	<b>38.2</b>
<b>Transformer w/ KB Transfer</b>	<b>TimeSformer</b>	<b>Unsupervised: distant supervision (ours)</b>	HT100M	<b>39.4</b>

Table 4. Accuracy of different methods on the step forecasting task using the COIN dataset.

Segment Model	Pretraining Supervision	Pretraining Dataset	Action (%)	Verb (%)	Noun (%)
TSN [53]	-	-	33.2	60.2	46.0
TRN [60]	-	-	35.3	65.9	45.4
TBN [25]	-	-	36.7	66.0	47.2
MoViNet [26]	-	-	<b>47.7</b>	<b>72.2</b>	57.3
TSM [31]	Supervised: action labels	Kinetics	38.3	67.9	49.0
SlowFast [15]	Supervised: action labels	Kinetics	38.5	65.6	50.0
ViViT-L [5]	Supervised: action labels	Kinetics	44.0	66.4	56.8
TimeSformer [7]	Supervised: action labels	Kinetics	42.3	66.6	54.4
<b>TimeSformer</b>	<b>Unsupervised: distant supervision (ours)</b>	HT100M	44.4	67.1	<b>58.1</b>

Table 5. Comparison to the state-of-the-art for classification of first-person videos using the EPIC-KITCHENS-100 dataset.

labels from Kinetics gives an accuracy of 79.4%. This same model trained on our step embeddings achieves an accuracy of 83.9%. The large gain confirms the superiority of our representation for this task and suggests that our features can be effectively plugged in different long-term models.

### 4.3.3 Step Forecasting

Table 4 shows that our learned representation and a shallow transformer can be used to forecast the next step very effectively. Our representation outperforms the features learned with Kinetics action labels by 3.5%. When the step order knowledge is leveraged by stacking the embeddings of the possible next steps, the gain is further improved to 4.7%. This shows once more the benefits of incorporating information from the wikiHow knowledge base.

### 4.3.4 Egocentric Video Understanding

Recognition of activities in EPIC-KITCHENS-100 [10] is a relevant testbed for our model since first-person videos in this dataset capture diverse procedural activities from daily human life. To demonstrate the generality of our distantly supervised approach, we finetune our pretrained model for the task of noun, verb, and action recognition in

egocentric videos. For comparison purposes, we also include the results of finetuning the same model pretrained on Kinetics-400 with manually annotated action labels. Table 5 shows that the finetuning of our distantly supervised model outperforms all prior works, with the only exception of MoViNet [26] which achieves higher accuracies for Action and Verb but not for Noun. This provides further evidence about the transferability of our models to other tasks.

## 5. Conclusion

In this paper, we introduce a distant supervision framework that leverages a textual knowledge base (wikiHow) to effectively learn step-level video representations from instructional videos. We demonstrate the value of the representation on step classification, long procedural video classification, and step forecasting. We further show that our distantly supervised model generalizes well to egocentric video understanding.

## Acknowledgments

Thanks to Karl Ridgeway, Michael Iuzzolino, Jue Wang, Noureldien Hussein, and Effrosyni Mavroudi for valuable discussions.



## References

- [1] Sentence Transformers. <https://www.sbert.net/>. 5
- [2] wikiHow. <https://www.wikiHow.com/>. 2
- [3] Jean-Baptiste Alayrac, Adria Recasens, Rosalia Schneider, Relja Arandjelovic, Jason Ramapuram, Jeffrey De Fauw, Lucas Smaira, Sander Dieleman, and Andrew Zisserman. Self-supervised multimodal versatile networks. *NeurIPS*, 2(6):7, 2020. 2
- [4] John R Anderson. Acquisition of cognitive skill. *Psychological review*, 89(4):369, 1982. 2
- [5] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. *arXiv preprint arXiv:2103.15691*, 2021. 8
- [6] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. *arXiv preprint arXiv:2104.00650*, 2021. 2
- [7] Gedas Bertasius, Heng Wang, and Lorenzo Torresani. Is space-time attention all you need for video understanding? *arXiv preprint arXiv:2102.05095*, 2021. 3, 5, 6, 7, 8
- [8] Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250, 2008. 3
- [9] Joao Carreira and Andrew Zisserman. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*, 2017. 1, 2, 6, 8
- [10] Dima Damen, Hazel Doughty, Giovanni Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, et al. The epic-kitchens dataset: Collection, challenges and baselines. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (01):1–1, 2020. 2, 6, 8
- [11] Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. Scaling egocentric vision: The epic-kitchens dataset. In *European Conference on Computer Vision (ECCV)*, 2018. 1, 2
- [12] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020. 5
- [13] Ehsan Elhamifar and Dat Huynh. Self-supervised multi-task procedure learning from instructional videos. In *European Conference on Computer Vision*, pages 557–573. Springer, 2020. 3
- [14] Ehsan Elhamifar and Zwe Naing. Unsupervised procedure learning via joint dynamic summarization. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6341–6350, 2019. 3
- [15] Christoph Feichtenhofer, Haoqi Fan, Jitendra Malik, and Kaiming He. Slowfast networks for video recognition. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 6202–6211, 2019. 7, 8
- [16] Christoph Feichtenhofer, Haoqi Fan, Bo Xiong, Ross Girshick, and Kaiming He. A large-scale study on unsupervised spatiotemporal representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3299–3309, 2021. 3
- [17] Chuang Gan, Chen Sun, Lixin Duan, and Boqing Gong. Webly-supervised video recognition by mutually voting for relevant web images and web video frames. In *European Conference on Computer Vision*, pages 849–866. Springer, 2016. 3
- [18] Deepti Ghadiyaram, Du Tran, and Dhruv Mahajan. Large-scale weakly-supervised pre-training for video action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 12046–12055, 2019. 3
- [19] Raghav Goyal, Samira Ebrahimi Kahou, Vincent Michalski, Joanna Materzynska, Susanne Westphal, Heuna Kim, Valentin Haenel, Ingo Fruend, Peter Yianilos, Moritz Mueller-Freitag, et al. The” something something” video database for learning and evaluating visual common sense. In *Proceedings of the IEEE international conference on computer vision*, pages 5842–5850, 2017. 2
- [20] Gabriel Huang, Bo Pang, Zhenhai Zhu, Clara Rivera, and Radu Soricut. Multimodal pretraining for dense video captioning. *arXiv preprint arXiv:2011.11760*, 2020. 2, 3
- [21] Noureldien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Timeception for complex action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 254–263, 2019. 2, 6, 8
- [22] Noureldien Hussein, Efstratios Gavves, and Arnold WM Smeulders. Videograph: Recognizing minutes-long human activities in videos. *arXiv preprint arXiv:1905.05143*, 2019. 8
- [23] Noureldien Hussein, Mihir Jain, and Babak Ehteshami Bejnordi. Timegate: Conditional gating of segments in long-range activities. *arXiv preprint arXiv:2004.01808*, 2020. 2
- [24] Will Kay, Joao Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Paul Natsev, et al. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017. 2
- [25] Evangelos Kazakos, Arsha Nagrani, Andrew Zisserman, and Dima Damen. Epic-fusion: Audio-visual temporal binding for egocentric action recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 5492–5501, 2019. 8
- [26] Dan Kondratyuk, Liangzhe Yuan, Yandong Li, Li Zhang, Mingxing Tan, Matthew Brown, and Boqing Gong. Movinets: Mobile video networks for efficient video recognition. *CoRR*, abs/2103.11511, 2021. 8
- [27] Hilde Kuehne, Ali Arslan, and Thomas Serre. The language of actions: Recovering the syntax and semantics of goal-directed human activities. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 780–787, 2014. 1, 2, 6

- [28] Hildegard Kuehne, Hueihan Jhuang, Estíbaliz Garrote, Tomaso Poggio, and Thomas Serre. Hmdb: a large video database for human motion recognition. In *Computer Vision (ICCV), 2011 IEEE International Conference on*, pages 2556–2563. IEEE, 2011. [2](#)
- [29] Jie Lei, Linjie Li, Luowei Zhou, Zhe Gan, Tamara L. Berg, Mohit Bansal, and Jingjing Liu. Less is more: Clipbert for video-and-language learning via sparse sampling. In *CVPR*, 2021. [7](#)
- [30] Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu. Hero: Hierarchical encoder for video+language omni-representation pre-training. *arXiv preprint arXiv:2005.00200*, 2020. [3](#)
- [31] Ji Lin, Chuang Gan, and Song Han. Temporal shift module for efficient video understanding. *arXiv preprint arXiv:1811.08383*, 2018. [8](#)
- [32] Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Jason Li, Taroon Bharti, and Ming Zhou. Univl: A unified video and language pre-training model for multimodal understanding and generation. *arXiv preprint arXiv:2002.06353*, 2020. [2](#)
- [33] Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. End-to-end learning of visual representations from uncurated instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9879–9889, 2020. [2](#), [5](#), [6](#), [7](#), [8](#)
- [34] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2630–2640, 2019. [1](#), [2](#), [3](#), [6](#)
- [35] Mike Mintz, Steven Bills, Rion Snow, and Dan Jurafsky. Distant supervision for relation extraction without labeled data. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 1003–1011, 2009. [2](#), [3](#)
- [36] Zhaofan Qiu, Ting Yao, and Tao Mei. Learning spatio-temporal representation with pseudo-3d residual networks. In *2017 IEEE International Conference on Computer Vision (ICCV)*, pages 5534–5542. IEEE, 2017. [3](#)
- [37] Jens Rasmussen. Skills, rules, and knowledge; signals, signs, and symbols, and other distinctions in human performance models. *IEEE transactions on systems, man, and cybernetics*, (3):257–266, 1983. [2](#)
- [38] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*. Association for Computational Linguistics, 11 2019. [5](#)
- [39] Sebastian Riedel, Limin Yao, and Andrew McCallum. Modeling relations and their mentions without labeled text. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, pages 148–163. Springer, 2010. [2](#)
- [40] Marcus Rohrbach, Sikandar Amin, Mykhaylo Andriluka, and Bernt Schiele. A database for fine grained activity detection of cooking activities. In *2012 IEEE Conference on Computer Vision and Pattern Recognition*, pages 1194–1201. IEEE, 2012. [1](#)
- [41] Marcus Rohrbach, Anna Rohrbach, Michaela Regneri, Sikandar Amin, Mykhaylo Andriluka, Manfred Pinkal, and Bernt Schiele. Recognizing fine-grained and composite activities using hand-centric features and script data. *Int. J. Comput. Vis.*, 2016. [2](#)
- [42] Rion Snow, Daniel Jurafsky, and Andrew Ng. Learning syntactic patterns for automatic hypernym discovery. In L. Saul, Y. Weiss, and L. Bottou, editors, *Advances in Neural Information Processing Systems*, volume 17. MIT Press, 2005. [2](#)
- [43] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. MpNet: Masked and permuted pre-training for language understanding. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 16857–16867. Curran Associates, Inc., 2020. [5](#)
- [44] 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. [2](#)
- [45] Nitish Srivastava, Elman Mansimov, and Ruslan Salakhudinov. Unsupervised learning of video representations using lstms. In *International conference on machine learning*, pages 843–852. PMLR, 2015. [3](#)
- [46] Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning, 2019. [3](#)
- [47] Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2818–2826, 2016. [7](#)
- [48] Hui Li Tan, Hongyuan Zhu, Joo-Hwee Lim, and Cheston Tan. A comprehensive survey of procedural video datasets. *Computer Vision and Image Understanding*, 202:103107, 2021. [2](#)
- [49] Yansong Tang, Dajun Ding, Yongming Rao, Yu Zheng, Danyang Zhang, Lili Zhao, Jiwen Lu, and Jie Zhou. Coin: A large-scale dataset for comprehensive instructional video analysis. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1207–1216, 2019. [6](#)
- [50] Yansong Tang, Jiwen Lu, and Jie Zhou. Comprehensive instructional video analysis: The coin dataset and performance evaluation. *IEEE transactions on pattern analysis and machine intelligence*, 2020. [1](#), [2](#), [6](#), [7](#)
- [51] Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. A closer look at spatiotemporal convolutions for action recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6450–6459, 2018. [1](#)
- [52] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Advances in neural*

*information processing systems*, pages 5998–6008, 2017. 2, 5

- [53] L. Wang, Y. Xiong, Z. Wang, Y. Qiao, D. Lin, X. Tang, and L. Van Gool. Temporal segment networks: Towards good practices for deep action recognition. In *ECCV*, 2016. 1, 8
- [54] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. *arXiv preprint arXiv:1711.07971*, 10, 2017. 8
- [55] Donglai Wei, Joseph J Lim, Andrew Zisserman, and William T Freeman. Learning and using the arrow of time. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8052–8060, 2018. 3
- [56] Hu Xu, Gargi Ghosh, Po-Yao Huang, Prahal Arora, Masoumeh Aminzadeh, Christoph Feichtenhofer, Florian Metze, and Luke Zettlemoyer. Vlm: Task-agnostic video-language model pre-training for video understanding. *arXiv preprint arXiv:2105.09996*, 2021. 2
- [57] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metze, Luke Zettlemoyer, and Christoph Feichtenhofer. Videoclip: Contrastive pre-training for zero-shot video-text understanding. *arXiv preprint arXiv:2109.14084*, 2021. 7
- [58] Zhongwen Xu, Yi Yang, and Alex G Hauptmann. A discriminative cnn video representation for event detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1798–1807, 2015. 3
- [59] Daojian Zeng, Kang Liu, Yubo Chen, and Jun Zhao. Distant supervision for relation extraction via piecewise convolutional neural networks. In *Proceedings of the 2015 conference on empirical methods in natural language processing*, pages 1753–1762, 2015. 3
- [60] Bolei Zhou, Alex Andonian, Aude Oliva, and Antonio Torralba. Temporal relational reasoning in videos. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 803–818, 2018. 8
- [61] Jiaming Zhou, Kun-Yu Lin, Haoxin Li, and Wei-Shi Zheng. Graph-based high-order relation modeling for long-term action recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8984–8993, 2021. 2, 8
- [62] Luowei Zhou, Chenliang Xu, and Jason J Corso. Towards automatic learning of procedures from web instructional videos. In *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018. 1, 2
- [63] Dimitri Zhukov, Jean-Baptiste Alayrac, Ramazan Gokberk Cinbis, David Fouhey, Ivan Laptev, and Josef Sivic. Cross-task weakly supervised learning from instructional videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3537–3545, 2019. 1, 2