

# Learning to Ground Instructional Articles in Videos through Narrations

Effrosyni Mavroudi\*, Triantafyllos Afouras\*, Lorenzo Torresani  
Meta AI

{emavroudi, afourast, torresani}@meta.com

## Abstract

In this paper we present an approach for localizing steps of procedural activities in narrated how-to videos. To deal with the scarcity of labeled data at scale, we source the step descriptions from a language knowledge base (wiki-How) containing instructional articles for a large variety of procedural tasks. Without any form of manual supervision, our model learns to temporally ground the steps of procedural articles in how-to videos by matching three modalities: frames, narrations, and step descriptions. Specifically, our method aligns steps to video by fusing information from two distinct pathways: i) direct alignment of step descriptions to frames, ii) indirect alignment obtained by composing steps-to-narrations with narrations-to-video correspondences. Notably, our approach performs global temporal grounding of all steps in an article at once by exploiting order information, and is trained with step pseudo-labels which are iteratively refined and aggressively filtered. In order to validate our model we introduce a new evaluation benchmark – HT-Step – obtained by manually annotating a 124-hour subset of *HowTo100M*<sup>1</sup> with steps sourced from *wikiHow* articles. Experiments on this benchmark as well as zero-shot evaluations on *CrossTask* demonstrate that our multi-modality alignment yields dramatic gains over several baselines and prior works. Finally, we show that our inner module for matching narration-to-video outperforms by a large margin the state of the art on the *HTM-Align* narration-video alignment benchmark.

## 1. Introduction

Instructional videos have emerged as a popular means for people to learn new skills and improve their abilities in executing complex procedural activities, such as cooking a recipe, performing home improvements, or fixing things. In addition to being useful teaching materials for humans, how-to videos are a promising medium for learning by ma-

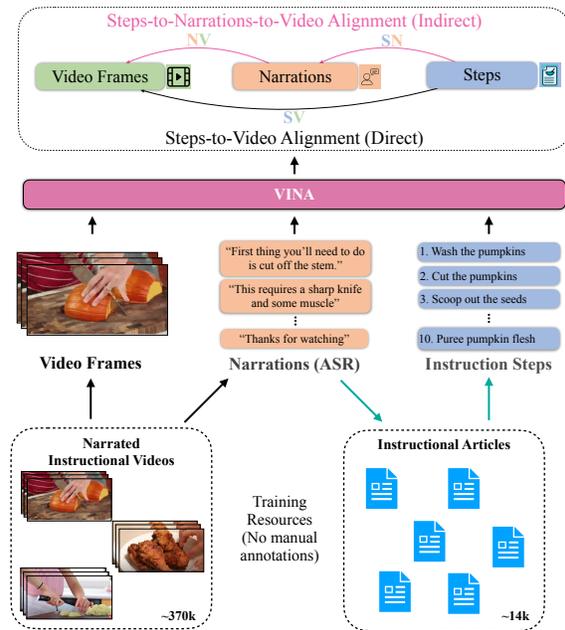


Figure 1: Our proposed Video, Instructions, and Narrations Aligner (VINA) learns to simultaneously ground narrations and instruction steps in how-to videos from an uncurated set of narrated videos and a separate knowledge base of instructional articles, *without any manual annotations*. This is contrary to prior work that learns how to align a video with a *single* sequence of sentences by leveraging *ground-truth pairs* of video-text sequences, e.g., a video and its narrations [19], or a video and an annotated, ordered list of steps demonstrated in it [14].

chines, as they provide revealing visual demonstrations of complex activities and show elaborate human-object interactions in a variety of domains. Motivated by this observation, in this work we look at the task of temporally localizing the steps of procedural activities in instructional videos. This problem is foundational to the broader goal of human-procedure understanding and advances on this task promise to enable breakthrough applications, such as AI-powered skill coaching and human-to-robot imitation learning.

Prior work has tackled procedural step localization by

\*equal contribution

<sup>1</sup>A test server is accessible at <https://eval.ai/web/challenges/challenge-page/2082>.

leveraging either (a) fully-annotated datasets where the task shown in the video is given (*video-level labeling*) and manually annotated temporal segments are provided for each step (*segment-level labeling*) [53] or (b) weakly-annotated training sets where the task and the order in which the steps appear in the video is given [70]. However, due to the inherent manual cost involved in collecting step annotations, these works have relied on datasets that are small-scale both in the number of tasks (e.g., at most few hundreds [70]) and in the number of video samples (e.g., 12k videos [53]). These limitations affect both the generality and the complexity of the models that can be trained on these benchmarks. In this paper, we therefore pose the following question: *can we leverage large-scale, unlabeled video datasets to train a model that can ground procedural steps in how-to videos?*

To answer this question, we propose a novel training framework for weakly-supervised step grounding that utilizes two freely available sources of information: (a) instructional articles which define ordered lists of steps for a wide variety of tasks (e.g., from wikiHow) and (b) narrations which provide instance-specific rich commentaries of the execution of the task in the video, e.g., from ASR transcriptions. Our work treats the former as an abstraction of the latter and uses the video-specific narrations to support the grounding of the article steps. Specifically, *during training*, our method leverages narrations as an auxiliary signal to (i) identify the task shown in the video, (ii) temporally ground the article steps that are visually-demonstrated and (iii) filter out steps that are not executed in the given instance. To further motivate this mechanism, let us look at the example in Figure 1. The narrations help disambiguate the task (*make a pumpkin puree*), enabling the automatic retrieval of relevant instructional articles for the video. Furthermore, the narrations can be matched to steps described in the articles to roughly localize the steps that are represented in the video. In this example, the timestamp of “First thing you’ll need to do is cut off the stem” provides a loose temporal prior for the matching step “Cut the pumpkins.” On the other hand, steps that do not have any matching narrations (e.g., “Wash the pumpkins”) are unlikely to be represented in the video and thus can be rejected. Based on this intuition, we propose a procedure that learns to align steps to video by fusing information from two pathways. The first is an *indirect* pathway inferring step-frame alignments by composing step-to-narration assignments with narration-to-frame correspondences. The second is a *direct* pathway that learns associations between step descriptions and frames by leveraging information from all videos having steps in common.

In our experiments we demonstrate that our multi-modality alignment leads to significant performance gains over several baselines, including single-pathway temporal grounding, as well adaptations of prior works to our prob-

lem. *During inference*, the direct pathway can be used by itself to temporally ground steps in absence of transcribed narrations. When narrations are available at test time, our method improves further the accuracy of temporal grounding by fusing the inference outputs of the two pathways.

To summarize, our work makes the following contributions: 1) we learn to align steps to frames in how-to videos, using only weak supervision in the form of noisy ASR narrations and instructional articles; 2) we propose a novel approach for joint dense temporal grounding of instructional steps and video narrations; 3) we introduce a new benchmark for evaluating instructional step grounding which we will make available to the community; 4) we demonstrate state-of-the-art results on multiple benchmarks for both step as well as narration grounding.

## 2. Related Work

**Procedural step recognition.** Prior work on procedural step localization [5, 8, 9, 16, 21, 31, 42, 45, 46, 58, 64, 67, 68, 70] can be roughly divided into two categories, based on the query formulation: the first class approaches the problem in an open-world setting, where the use of text queries transforms it into a temporal grounding task [2, 4, 18]. Such approaches can be further sub-divided into single step grounding, where single steps are queried over the whole video [24, 51] and dense temporal grounding methods [12, 19] where the objective is to jointly ground a sequence of steps or whole article into the video. The second body of works uses fixed taxonomies of steps, often as part of activities [26, 44, 53]. Our work is somehow related to Lin et al. [29] who use semantic similarity between steps and narrations to obtain supervision for learning strong video representations. Although we also associate steps from wikiHow articles to video frames through the use of narrations, the two works differ in several aspects: we align steps to video by a global procedure that takes into account all ordered steps in the article (inspired by dense temporal grounding methods [4]) and temporally grounds them in the whole video, instead of matching individual video clips to an orderless collection of steps; our step grounding uses *video* in addition to steps and narrations while the method proposed in Lin et al. relies purely on text-matching narrations to step descriptions; finally, the works differ in objectives with our aim being step grounding in long how-to videos rather than learning video-clip representations.

Existing methods also vary by the level of supervision used during training. One option is leveraging fully-annotated datasets with known temporal segments for each step [9, 23, 26, 44, 47, 54, 65, 71], using weakly-annotated training sets where the task and the order in which the steps appear in the video are known [6, 7, 10, 14, 17, 25, 41, 70], only the task and potential steps are known [40], or only

loose association between video and instructional articles is given [12, 15]. Video narrations are a commonly used source of weak supervision [1, 17, 34, 43], while instructional steps from knowledge bases have recently been used as supervision: [12, 29]. Chen et al. [12] use video-level instructional step labels for (weak) supervision of a model that grounds instructional articles to videos. This approach attempts to localize steps without using any narration information; we instead show that the task knowledge is not necessary and heavily exploit narrations via multi-task learning and complementary inference pathways [66]: we argue that narrations provide a much richer source of supervision for training step grounding models, while essentially coming for free.

**Video-Text alignment** The availability of large-scale video-text datasets such as HowTo100M has prompted many works on joint video-language embedding training [29, 35]. A form of contrastive loss is often adopted for bringing together the representations of the two modalities [3, 32, 33, 35, 38, 39, 61, 63], while masked objectives are also gaining popularity [13, 20, 28, 33, 50, 51, 52, 55, 69]. Some works perform end-to-end representation learning [35, 36], while others freeze representation and focus on longer-term temporal modelling, which aims to capture context [61]. More recently Han *et al.* investigated directly aligning contextualized narration representations to video frames [19]. We build our method off of this approach – we note however that our objective is complementary: rather than aligning a video’s narrations as an end-goal, we use this functionality to ground a set of independent steps sourced from instructional articles; in that process we show that the synergy that develops while training jointly on the two tasks results in improved performance for both.

### 3. Narration-Aided Step Grounding

We first present our architecture for joint narration and step grounding (Sec. 3.2), followed by learning objectives (Sec. 3.3.1) and pseudo-labeling strategy (Sec. 3.3.3); we discuss inferring the video task in (Section 3.3.2).

#### 3.1. Problem Formulation

Let  $(\mathcal{V}, \mathcal{N})$  be a video-narration pair, consisting of  $T$  video frames and a sequence of  $N$  narrations. Also, let  $\mathcal{S}$  be an ordered list of  $S$  steps from an instructional article for a candidate task  $\tau$ . Our objective is to ground each step of  $\mathcal{S}$  to the video, conditioned on the other steps and the ASR transcript<sup>2</sup>. In particular, the desired output of our model is an alignment matrix  $Y^{SV} \in \{0, 1\}^{S \times T}$ , where  $Y_{st} = 1$  only if frame  $t$  is depicting the  $s$ -th step of task  $\tau$ , and zero otherwise. Note that some steps might not be represented in the video.

<sup>2</sup>ASR transcripts are assumed to be always available for training and optionally during inference.

#### 3.2. Joint Narration and Article Step Grounder

As shown in Figure 2, our proposed VINA model follows the popular paradigm of leveraging Transformers for modeling multimodal interactions [56, 62].

**Unimodal Encoders.** Before feeding the video, narrations and article steps to our model, we preprocess them to extract a sequence of tokens. Given a video-narration pair  $(\mathcal{V}, \mathcal{N})$  we extract visual features,  $X^v \in \mathbb{R}^{T \times D_v}$  and narration features  $X^n \in \mathbb{R}^{N \times D_n}$  using standard backbone networks (e.g., a frozen S3D [60] network for visual features, and pooled Word2Vec [37] embeddings for narration features). Similarly, we encode the sequence of steps in a sequence of features  $X^s \in \mathbb{R}^{S \times D_s}$ . The features of each modality  $m$  are embedded into a common embedding space of dimensionality  $D$  using a Unimodal Encoder that consists of a modality-specific MLP network, and then learnable, modality-specific positional embeddings  $P^m$  are added to them:

$$H^m = MLP(X^m; \theta_m) + P^m, \quad (1)$$

where  $m \in \{v, n, s\}$  denotes the modality.

**Multimodal Encoder.** The outputs of the Unimodal Encoders are concatenated into a sequence of tokens:  $H = [H^v; H^n; H^s] \in \mathbb{R}^{(T+N+S) \times D}$  and fed to the Multimodal Encoder, which is a standard Transformer with multiple layers of multi-head self-attention:

$$Z = \text{Transformer}(H) \in \mathbb{R}^{(T+N+S) \times D}. \quad (2)$$

The contextualized embeddings  $Z = [Z^v; Z^n; Z^s]$  computed by the Multimodal Encoder capture interactions within each modality (e.g., temporal relationships within the video and context among steps of an article) and across modalities. We can then compute cosine similarity matrices between all pairs of modalities: narrations-to-video  $A^{NV} \in \mathbb{R}^{N \times T}$ , steps-to-video  $A^{SV} \in \mathbb{R}^{S \times T}$ , and steps-to-narrations  $A^{SN} \in \mathbb{R}^{S \times N}$ . For example, the narrations-to-video similarity matrix  $A^{NV}$  is obtained by simply computing the cosine similarity between each frame embedding and each narration embedding:  $A_{nt}^{NV} = \mathbf{z}_n^n \top \mathbf{z}_t^v / (\|\mathbf{z}_n^n\| \|\mathbf{z}_t^v\|)$ .

**Narration-aided Step Grounding.** A straightforward inference path for temporally grounding the steps in the video is directly through the  $A^{SV}$  similarity matrix, which captures the similarity of each video frame with each instructional step. However, this alignment does not explicitly take into account the narrations of the video (only implicitly, through the Multimodal Transformer). We observe that an alternative way to ground steps in a video is to first identify narrations in the ASR transcript that are relevant to the step and then exploit the similarity of those narrations with video frames to get a loose prior over the step location.

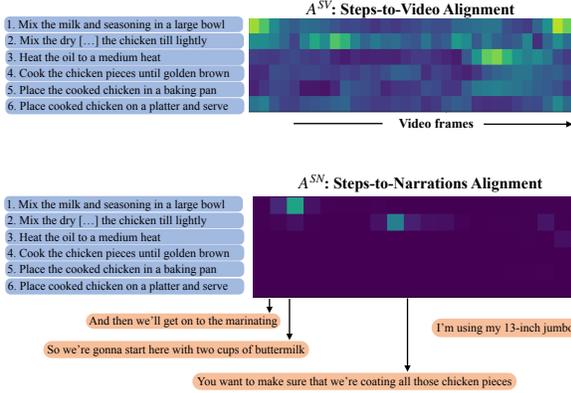
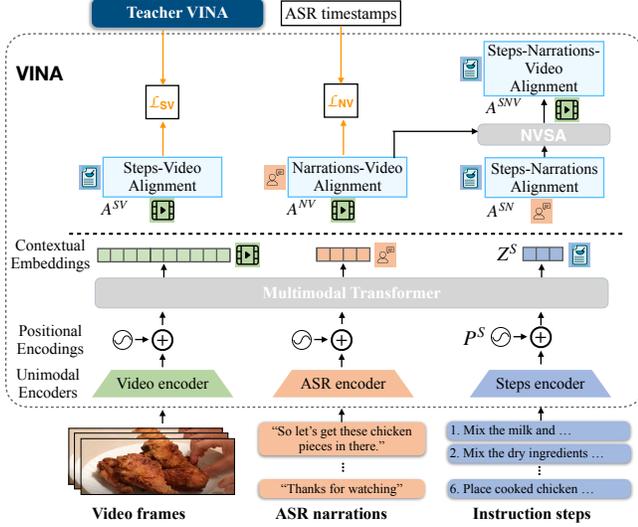


Figure 2: (left) Schematic illustration of our system. First, it extracts representations for each modality (video, ASR narrations, task steps) with three *Unimodal Encoders*. These representations are fed to a Transformer-based *Multimodal Encoder* for capturing interactions among video frames as well as among the textual modalities and the video. The contextualized embeddings for video frames, narrations and steps are used to compute correspondences between *all pairs of modalities*. Step grounding is achieved by fusing the output of two pathways: a *direct* pathway aligning steps to video ( $A^{SV}$ ) and an *indirect* pathway that composes steps-to-narration ( $A^{SN}$ ) with narration-to-video ( $A^{NV}$ ) alignments to produce a second step-to-video alignment ( $A^{SNV}$ ). We train our model using a teacher-student strategy with iteratively refined and filtered step pseudo-labels. (right) Qualitative examples of our learned steps-to-video alignment and steps-to-narrations alignment matrices for a video snippet.

This is computed by combining the information captured in the steps-to-narrations and narrations-to-video alignment matrices  $A^{SN}$  and  $A^{NV}$ :

$$A^{SNV} = \tilde{A}^{SN} A^{NV} \in \mathbb{R}^{S \times T}, \quad (3)$$

where  $\tilde{A}^{SN}$  is the predicted steps-to-narrations alignment matrix  $A^{SN}$  after being normalized with a softmax function with temperature  $\xi$ :  $\tilde{A}_{sn}^{SN} = \frac{\exp(A_{sn}/\xi)}{\sum_{j=1}^N \exp(A_{sj}/\xi)}$ .

The resulting  $A^{SV}$  and  $A^{SNV}$  alignment matrices provide two complementary inference paths to align steps to video frames. The mutual agreement between the direct steps-to-video alignment provided by  $A^{SV}$  and indirect, narration-based steps-to-video alignment provided by  $A^{SNV}$  can be used to better ground steps. Intuitively, if a frame is both very similar to a step in the joint embedding space learned by the Multimodal Transformer, and also very similar to a narration that is relevant to the step, then it is more likely to be indeed relevant to the step. Hence, we fuse the  $A^{SV}$  and  $A^{SNV}$  alignment matrices to a matrix  $A^F = (A^{SV} + A^{SNV})/2$ .

### 3.3. Weakly-Supervised Training from Narrated Instructional Videos

Next, we discuss how to supervise the VINA model in order to learn steps-to-video alignment and narrations-to-

video alignment. We first present the training objective assuming that the ground-truth temporal segments for each narration and step in the video are given. Then we describe our approach for obtaining automatic pseudo-labels for the temporal segments.

#### 3.3.1 Learning on Labeled Data

Let  $\mathcal{B} = \{\mathcal{V}_i, \mathcal{N}_i, \mathcal{S}_i, Y_i^{NV}, Y_i^{SV}\}_{i=1}^B$  denote a set of training tuples, each comprising a video-narration pair, an ordered list of relevant task steps, and the target video-narrations and video-steps alignment matrices, we train the VINA model by optimising the following objective:

$$\mathcal{L} = \frac{1}{B} \left[ \sum_{i=1}^B \lambda_{NV} \mathcal{H}(Y_i^{NV}, A_i^{NV}) + \lambda_{SV} \mathcal{H}(Y_i^{SV}, A_i^{SV}) \right], \quad (4)$$

where  $\mathcal{H}(\cdot, \cdot)$  is the modified InfoNCE loss used by [19] for aligning video with text using noisy ground-truth temporal segments:

$$\mathcal{H}(Y, A) = -\frac{1}{K} \sum_{k=1}^K \log \frac{\sum_t Y_{t,k} \exp(A_{t,k}/\eta)}{\sum_t \exp(A_{t,k}/\eta)}, \quad (5)$$

where  $\eta$  is a temperature constant. Note that although we do not explicitly supervise the steps-narrations alignment

$A^{SN}$ , meaningful alignments emerge during training due to the joint grounding of narrations and steps to the same video samples, as seen in Figure 2. Note that although we do not directly supervise the steps-to-narrations alignment, our model is able to learn meaningful correspondences, which go beyond simple pairwise textual matching.

### 3.3.2 Pairing Videos with Articles

We assume access to a set of instructional articles  $\mathcal{A} = \{\mathcal{S}_j, \tau_j\}_{j=1}^W$ , where  $\tau_j$  denotes the article title and  $\mathcal{S}_j$  the associated set of steps. To assign a set of steps to a given video from our training set  $\mathcal{B}$  we need to associate it with an article from  $\mathcal{A}$ . If our video dataset provides metadata (e.g., a task id for every video), then this can be used to obtain the association – although there is no guarantee that this will result in the best article-match for the video (see discussion in supplementary materials for more analysis). If such metadata is not available, we can predict a task id, using the similarity between the narration and the titles of the available articles. To that end we use an off-the-shelf language model (e.g. MPNet [49]) to compute semantic embeddings of the ASR captions of every video and the title of each article  $\tau_j \in \mathcal{W}$ . For every video  $\mathcal{V}_i$  we then calculate the semantic similarity between all the  $N$  captions in  $\mathcal{N}_i$  and all task titles  $\tau_j \in \mathcal{W}$ , and assign  $N$  votes; the vote of every caption goes to the task that best matches it. Finally the video is assigned the task with the most votes. Alternatively, in order to obtain multiple sets of steps for a video, we rank the tasks by the number of votes.

### 3.3.3 Narration-Aided Pseudo-Labeling

Once a task  $\tau_j$  has been associated with a video, we have access to a list of steps  $\mathcal{S}_j$  from the article of the task. However, whether these steps appear in the video and their temporal location remain unknown. Inspired by self-labeling approaches from the SSL literature [27, 48, 59], we follow a teacher-student approach where a teacher version of our models generates pseudo-labeled temporal segments for training the student. For every step represented by a row in the learned steps-to-video alignment matrix we obtain a pseudo-ground truth segment by finding the maximal activation (peak) and expanding a temporal segment on both sides until the activation falls below an adaptive threshold  $\zeta$  (e.g., 70% of the peak). To avoid training with unreliable pseudo-labels, we filter out pseudo-labels with low confidence: if the peak activation is below a fixed threshold  $\gamma$ , the alignment of that step is treated as unreliable for pseudo-labeling, and is altogether ignored.

**Training curriculum.** For the first  $E_b$  epochs we perform burn-in training of the student model on fixed pseudo-labels generated by feeding the video and the list of steps  $\mathcal{S}_j$  to

TAN [19], an off-the-shelf model pre-trained on the task of video-text alignment. Afterwards, we switch to using pseudo-labels generated from the teacher, where the teacher is initialized by duplicating the burn-in student model and then updated every  $\nu$  epochs. During both stages, we utilize the original ASR timestamps for supervising the video-to-narrations alignment.

## 4. Experiments

### 4.1. Datasets and Metrics

We train our models on narrated videos from the HowTo100M dataset by leveraging the dataset release of wikiHow instructional articles [22], without using any form of manual annotations. In order to evaluate the effectiveness of our method, we evaluate: step grounding on HT-Step (a new benchmark, described below), narration alignment on HTM-Align [19], and zero-shot step localization on CrossTask [70].

**HowTo100M (Training).** The HowTo100M dataset [36] contains instructional videos from YouTube. Following Han *et al.* [19], we use the Food & Entertainment subset containing approximately 370K videos, where each video is complemented by the “sentencified” ASR transcription of its audio narration.

**wikiHow (Training).** We train using 14,541 cooking tasks from the wikiHow-Dataset [22]. For each task, we generate an ordered list of steps by extracting the step headlines.

**CrossTask (Evaluation).** We use this established instructional video benchmark for *zero-shot* grounding, i.e., by directly evaluating on CrossTask our model learned from HowTo100M. Following common practices, we use two evaluation protocols: the first one – *step localization* – aims at predicting a single timestamp for each occurring step in videos from 18 primary tasks [70]. Performance is evaluated by computing the recall (denoted as Avg. R@1) of the most confident prediction for each task and averaging the results over all query steps in a video, where R@1 measures whether the predicted timestamp for a step falls within the ground truth boundaries. We report average results over 20 random sets of 1850 videos [70]. The second task – *article grounding* – requires predicting temporal segments for each step of an instructional article describing the task represented in the video. We use the mapping between CrossTask and *simplified* wikiHow article steps provided in Chen *et al.* [12] and report results on 2407 videos of 15 primary tasks obtained excluding three primary tasks following the protocol of [12] (see supplementary materials for details). Performance for this task is measured with Recall@K at different IoU thresholds [12].

**HT-Step (Evaluation).** To evaluate the effectiveness of our model in grounding steps, we introduce an evaluation

benchmark consisting of 1200 HowTo100M videos spanning a total of 177 unique tasks, with each video manually annotated with temporal segments for each occurring step. For each video, annotators were provided with the task name (e.g., Make Pumpkin Puree) and the recipe steps from the corresponding [wikiHow article](#). We refer the reader to supplementary materials for details about the data annotation. We split the annotated videos into a validation and a test set, each containing 600 videos, with 5 videos per task. We ensure that our validation set does not contain videos from HTM-Align.

**HTM-Align (Evaluation).** This benchmark is used to evaluate our model on narration grounding. It contains 80 videos where the ASR transcriptions have been manually aligned temporally with the video. We report the R@1 metric [19], which evaluates whether the model can correctly localize the narrations that are alignable with the video.

## 4.2. Implementation Details

As video encoder we adopt the S3D [60] backbone pre-trained with the MIL-NCE objective on HowTo100M [35]. Following previous work [19, 61], we keep this module frozen and use it to extract clip-level features (one feature per second for video decoded at 16 fps). For extracting context-aware features for each sentence (step or narration), we follow the Bag-of-word (BoW) approach based on Word2Vec embeddings [37]. Our methods hyperparameters were selected on the HT-Step validation set and are:  $\lambda_{SV} = \lambda_{NV} = 1$ , temperatures  $\eta, \xi = 0.07$ , and threshold  $\gamma = 0.65$ . We train our model for 12 epochs, with 3 epochs burn-in and then we update the teacher every 3 epochs. Pseudo-labels are obtained based on the steps-to-video alignment matrix. To obtain temporal segment detections from the step-to-video alignment output of our model (e.g. for evaluating on the CrossTask article grounding setting) we use a simple 1D blob detector [57]. Unless otherwise specified, we use the fused alignment matrix for step grounding when narrations are available during inference time. More details are included in supplementary materials.

## 4.3. Results

### 4.3.1 Comparison with the State of the Art

**Weakly-Supervised Narration and Step Grounding.** Table 1 compares the step and narration grounding performance of our method with recent state-of-the-art video-text alignment methods trained on HowTo100M using ASR narrations: MIL-NCE [35] and TAN [19]. When using them for narration alignment, we feed them with ASR as input. But we also evaluate them as strong baselines for zero-shot step grounding by feeding them with the sequence of steps as input. Our model achieves 66.5% R@1 on narration alignment on HTM-Align, leading to an absolute

Method	Train. Inp.	HT-Step $\uparrow$ R@1		HTM-Align $\uparrow$ R@1
		w/o nar.	w/ nar.	
CLIP (ViT-B/32) [39]	-	-	-	23.4
MIL-NCE [35]	N	<u>30.7</u>	-	34.2
TAN (Joint+Dual, S2) [19]	N	-	-	<u>49.4</u>
TAN* (Joint, S1, LC) [19]	N	31.2	-	47.1
TAN* (Joint, S1, PE+LC) [19]	N	7.9	-	63.0
Ours	N+S	<b>35.6 <math>\pm</math> 0.4</b>	<b>37.4 <math>\pm</math> 0.4</b>	<b>66.5 <math>\pm</math> 0.9</b>

Table 1: **Comparison with state-of-the-art methods for step and narration grounding.** We report results on the HT-Step and HTM-Align test sets, respectively. TAN\* refers to our improved baselines of [19]. S1 and S2 refer to the training stages followed in [19]. PE denotes the addition of positional encoding to the output of the narration encoder. LC denotes long context, *i.e.*, our improved TAN\* baseline using 1024 seconds of context as opposed to 64 for TAN. Previous best results are shown underlined. Our VINA results are reported after 5 random runs. VINA clearly outperforms all previous work – as well as our improved TAN baselines – by large margins on both narration alignment and step grounding.

improvement of 17.1% over the previously reported state-of-the-art (49.4%). Notably, on HTM-Align our method surpasses TAN\* (Joint, S1, PE, LC) which is a new version of TAN [19] implemented by us and much stronger in video-narration alignment. Our re-implementation uses positional encodings for ASR narrations, is trained on long-form videos (up to 17 minutes) only with original ASR timestamps, while TAN was trained on 1 min video-clips with refined narration timestamps and used a fusion of two models during inference (Joint+Dual). Our method also outperforms all baselines for step grounding on HT-Step even when seeing only steps during inference, while being trained with (video, narrations, steps) triplets. It also outperforms TAN\* (Joint, S1, LC), which is a second re-implementation of TAN designed for maximum performance on the task of step grounding. Additionally, VINA is able to use ASR transcripts of videos during test time, if available, to further boost the performance.

**Step localization on CrossTask.** In Table 2 we compare our model against the state-of-the-art in step localization on the CrossTask benchmark. Our approach sets a new state-of-the-art for zero-shot step localization on this challenging benchmark. Importantly, most approaches are evaluated on this dataset by feeding their predicted steps-to-frames alignment matrix to a dynamic programming algorithm which finds the optimal assignment of each step with exactly one short clip *assuming a canonical, fixed ordering of steps for each task*. In contrast, our method, which is naturally aware of context and ordering by densely grounding steps, can

Method	↑Avg. R@1 (%)
<i>Supervised</i>	
TempCLR [64]	52.5
<i>Zero-Shot</i>	
HT100M [36]	33.6
VideoCLIP [61]	33.9
MCN [11]	35.1
DWSA [45]	35.3
MIL-NCE [35]	40.5
Zhukov [70]	40.5
VT-TWINS* [21]	40.7
UniVL [31]	42.0
Ours w/o nar.	<b>44.1</b>
Ours w/ nar.	<b>44.8</b>

Table 2: **Comparison with state-of-the-art methods for zero-shot action step localization on the CrossTask dataset.** The performance of the state-of-the-art fully-supervised method (TempCLR [64]) is reported as an upper-bound to the zero-shot approaches. \* denotes results reported on different test splits, and hence not directly comparable with the rest. Our model outperforms all previous works by a clear margin (2.1% absolute improvement over the previous best result on the standard split). When providing narrations as additional inputs during inference (only text, not the timings), we obtain a further 0.7% boost.

outperform prior results without imposing any constraints during inference.

Model	↑R@50(IOU)			↑R@100(IOU)		
	0.1	0.3	0.5	0.1	0.3	0.5
MIL-NCE-max [35]	33.5	12.0	4.9	39.7	14.3	5.9
MIL-NCE-avg [35]	42.9	24.3	12.9	56.8	32.1	17.0
WSAG [12]	40.1	23.1	10.1	54.3	31.3	14.0
Ours	<b>87.1</b>	<b>59.0</b>	<b>30.0</b>	<b>90.6</b>	<b>61.1</b>	<b>30.9</b>

Table 3: **Comparison with state-of-the-art approaches for article grounding on the CrossTask dataset.**

**Article grounding on CrossTask.** VINA is robust to the type of language in which task steps are described. It can handle both atomic phrases (as demonstrated by our results on step localization on CrossTask), but also rich, natural language step descriptions, as evidenced by performance on HT-Step. To further demonstrate this, we compare against the state-of-the-art on the article grounding task of CrossTask in Table 3. Our model outperforms all previous works by a large margin. We emphasize the performance improvement we obtain compared to WSAG, which highlights the importance of exploiting the narration information for training.

### 4.3.2 Ablation Studies

We perform ablations to assess the impact of the various design choices in our method by measuring step grounding performance on the HT-Step validation set and video-narration alignment performance on HTM-Align.

Method	Train. Inp.	Iter. Pseudo.	HT-Step ↑R@1		HTM-Align
			w/o nar.	w/ nar.	
<i>Baseline/Initial Step Pseudo-labels</i>					
(1) TAN Joint S1	N		30.7	-	47.1
<i>Single-Task Training</i>					
(2) Ours	N		-	-	63.2
(3) Ours	S		34.0	-	-
(4) Ours	S	✓	35.8	-	-
<i>Multi-Task Training</i>					
(5) Ours	N+S		34.3	36.1	64.8
(6) Ours	N+S	✓	<b>36.9</b>	<b>39.1</b>	<b>67.0</b>

Table 4: **Ablation of main components of our framework.** We study the contribution of (a) multi-task training for narration and step grounding, (b) iterative step pseudo-labeling (*Iter. Pseudo*), and (c) narration-aware step grounding (*w/ nar.*). We report results on the HT-Step val set for STG and HTM-Align for NG. We compare training only with narrations (N), only with wikiHow steps (S), and training with narrations-steps sequence pairs (N+S). We also compare the performance with and without providing narrations during inference.

**Effect of weak supervision from instructional articles.** Row 3 in Table 4 shows the step grounding results obtained from an instance of our model that includes only the direct video-step alignment pathway and that is trained just on wikiHow steps (without narrations) using the fixed step pseudo-labels from TAN\* [19] without any form of iterative pseudo-labeling (row 1). Remarkably, this variant improves by 3.3% over the step-grounding performance of TAN\*. When we let this variant update the step pseudo-labels (row 4), the recall improves further (5.1% over TAN\*). These results provide evidence of the strong benefits of utilizing instructional articles for the learning of step grounding.

**Effect of multimodal training and inference.** Training our model with multi-modal textual inputs (steps and narrations), we observe an improvement of 1.6% in narration grounding (row 5 of Table 4) compared to its single-task counterpart (row 2). However the gain in step grounding is marginal when seeing only video frames during inference (*w/o nar.*, 34% in row 3 vs 34.3% in row 5). Our conjecture is that the missing modality (narrations) leads to some drop in performance. Providing both steps and narrations during inference leads to a stronger step grounding performance, which surpasses the TAN\* baseline by 5.4% (30.7 → 36.1).

**Effect of iterative pseudo-labeling.** By comparing row 5 to row 6 of Table 4 we observe a clear boost in performance on both step grounding and narration alignment. This is a clear indication of the gains produced by iteratively refining the pseudo-labels using our model as a teacher during training.

Alignment	S → V	S → N	N → V	HT-Step ↑R@1
S → V	learned	-	-	34.3
S → N → V	-	learned	learned	30.5
S → N → V	-	learned	ASR	27.9
S → N → V	-	MPNet [49]	ASR	19.0
Fused	learned	learned	learned	<b>36.1</b>

Table 5: **Impact of the alignment matrix used during inference with narrations.** The same model is used for all results (corresponding to row 5 in Table 4).

**Impact of pathways during inference.** In Table 5 we study the effects of using different pathways and alignment information during inference. All results are produced from the same model trained for joint narration and step grounding with fixed pseudo-labels from TAN (row 5 in Table 4). Grounding steps using the indirect steps-to-video alignment only lags by 3.8% behind the direct steps-to-video alignment that directly computes the similarity between steps and video frames (30.5% vs 34.3%). Their fusion outperforms their individual grounding performance. This suggests that they capture complementary information. We also explore substituting our learned steps-to-narrations alignment with an alignment computed with an off-the-shelf language model. This significantly degrades performance (19.0%) showing that our joint steps and narrations grounding model learns relationships between steps and narrations that go beyond textual similarity between pairs of sentences. Similarly, substituting our learned narrations-to-video alignment with an alignment based on the original ASR timestamps reduces performance by 2.6%.

**Iterative pseudo-labeling strategies.** In Table 6 we ablate design choices for the iterative pseudo-labeling stage. We can observe that using aggressive filtering (i.e., high thresholds translating to a high maximum percentage of pseudo-labels that are discarded) is key to observing gains from iterative pseudo-labeling (using either the S → V or Fused alignment matrices) compared to training with fixed pseudo-labels from TAN. Intuitively, a large percentage of steps described in wikiHow articles are not represented in the given instructional video due to task mismatch, variations in recipe execution, and some steps being optional. Therefore, starting with a small subset of reliable pseudo-labels can facilitate step grounding.

**Task selection.** In Table 7 we investigate different strategies to select the wikiHow articles during training. This selection determines the set of steps to be grounded. We

Alignment	$\gamma$	max % step discarded	HT-Step ↑R@1	
			w/o nar.	w/ nar.
S → V	0.40	24	34.3	36.5
S → V	0.65	91	<b>36.9</b>	<b>39.1</b>
Fused	0.40	60	34.1	35.9
Fused	0.55	88	36.2	35.5

Table 6: **Ablation of the type of alignment matrix and filtering threshold used for pseudo-label generation.** Pseudo-label generation with the steps-to-video alignment matrix and the fusion of the direct and indirect pathways perform comparably for step grounding. Aggressive unreliable pseudo-label filtering with high confidence thresholds  $\gamma$  (large maximum step discard ratio) helps in both cases.

Task ID selection	HT-Step ↑R@1
HT100M metadata	34.3
Top-1 prediction	34.7
Random / top-5 pred	34.3

Table 7: **Sensitivity to task id selection.** We assess how the performance of our method changes when using different strategies to associate videos with articles. We experiment with using the task ids available from HT100M, as well as the the two predictive strategies presented in Section 3.3.2. We conclude that our method is robust to the task selection, and the task labels are not necessary for training.

evaluate two strategies for video-task association from narrations and compare them with using the task id provided in the HowTo100M metadata for each video<sup>3</sup>. We see that our automatic selection approaches yield results on par with or even slightly better than those based on metadata.

## 5. Conclusion

We have presented a method for learning how to temporally ground sequences of steps in instructional videos, without any manual supervision. Our proposed method exploits the weak supervision naturally provided in such videos through their narrations, and solves for joint alignment of narrations and steps, while fusing two complementary pathways for step-to-video alignment. We demonstrated strong quantitative performance, surpassing the state-of-the-art on multiple benchmarks for both narration and step grounding.

**Acknowledgements.** We thank Huiyu Wang, Yale Song, Mandy Toh, and Tengda Han for helpful discussions.

<sup>3</sup>During inference, metadata task ids are used in all of our HowTo100M experiments in order to evaluate against the ground-truth step annotations.

## References

- [1] Jean-Baptiste Alayrac, Piotr Bojanowski, Nishant Agrawal, Josef Sivic, Ivan Laptev, and Simon Lacoste-Julien. Unsupervised learning from narrated instruction videos. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4575–4583, 2016. [3](#)
- [2] Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. Localizing moments in video with natural language. In *Proceedings of the IEEE international conference on computer vision*, pages 5803–5812, 2017. [2](#)
- [3] Max Bain, Arsha Nagrani, Gül Varol, and Andrew Zisserman. Frozen in time: A joint video and image encoder for end-to-end retrieval. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 1728–1738, 2021. [3](#)
- [4] Peijun Bao, Qian Zheng, and Yadong Mu. Dense events grounding in video. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 920–928, 2021. [2](#)
- [5] Jing Bi, Jiebo Luo, and Chenliang Xu. Procedure planning in instructional videos via contextual modeling and model-based policy learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15611–15620, 2021. [2](#)
- [6] Piotr Bojanowski, Rémi Lajugie, Francis Bach, Ivan Laptev, Jean Ponce, Cordelia Schmid, and Josef Sivic. Weakly supervised action labeling in videos under ordering constraints. In *Computer Vision—ECCV 2014: 13th European Conference, Zurich, Switzerland, September 6–12, 2014, Proceedings, Part V 13*, pages 628–643. Springer, 2014. [2](#)
- [7] Piotr Bojanowski, Rémi Lajugie, Edouard Grave, Francis Bach, Ivan Laptev, Jean Ponce, and Cordelia Schmid. Weakly-supervised alignment of video with text. In *Proceedings of the IEEE international conference on computer vision*, pages 4462–4470, 2015. [2](#)
- [8] Kaidi Cao, Jingwei Ji, Zhangjie Cao, Chien-Yi Chang, and Juan Carlos Niebles. Few-shot video classification via temporal alignment. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 10615–10624, 2020. [2](#)
- [9] Meng Cao, Tianyu Yang, Junwu Weng, Can Zhang, Jue Wang, and Yuexian Zou. Locvtp: Video-text pre-training for temporal localization. In *European Conference on Computer Vision*, 2022. [2](#)
- [10] C. Chang, De-An Huang, Yanan Sui, Li Fei-Fei, and Juan Carlos Niebles. D3tw: Discriminative differentiable dynamic time warping for weakly supervised action alignment and segmentation. *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3541–3550, 2019. [2](#)
- [11] Brian Chen, Andrew Rouditchenko, Kevin Duarte, Hilde Kuehne, Samuel Thomas, Angie Boggust, Rameswar Panda, Brian Kingsbury, Rogerio Feris, David Harwath, et al. Multimodal clustering networks for self-supervised learning from unlabeled videos. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 8012–8021, 2021. [7](#)
- [12] Long Chen, Yulei Niu, Brian Chen, Xudong Lin, Guangxing Han, Christopher Thomas, Hammad Ayyubi, Heng Ji, and Shih-Fu Chang. Weakly-supervised temporal article grounding. In *Conference on Empirical Methods in Natural Language Processing*, pages 9402–9413, Abu Dhabi, United Arab Emirates, Dec. 2022. Association for Computational Linguistics. [2](#), [3](#), [5](#), [7](#), [14](#)
- [13] Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. Uniter: Universal image-text representation learning. In *ECCV*, 2020. [3](#)
- [14] Nikita Dvornik, Isma Hadji, Konstantinos G. Derpanis, Animesh Garg, and Allan D. Jepson. Drop-dtw: Aligning common signal between sequences while dropping outliers. *ArXiv*, abs/2108.11996, 2021. [1](#), [2](#)
- [15] Nikita Dvornik, Isma Hadji, Hai Pham, Dhaivat Bhatt, Brais Martinez, Afsaneh Fazly, and Allan D Jepson. Flow graph to video grounding for weakly-supervised multi-step localization. In *European Conference on Computer Vision*, pages 319–335. Springer, 2022. [3](#)
- [16] Ehsan Elhamifar and Dat Huynh. Self-supervised multi-task procedure learning from instructional videos. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVII 16*, pages 557–573. Springer, 2020. [2](#)
- [17] Daniel Fried, Jean-Baptiste Alayrac, Phil Blunsom, Chris Dyer, Stephen Clark, and Aida Nematzadeh. Learning to segment actions from observation and narration. In *Annual Meeting of the Association for Computational Linguistics*, pages 2569–2588, Online, July 2020. Association for Computational Linguistics. [2](#), [3](#)
- [18] Simon Ging, Mohammadreza Zolfaghari, Hamed Pirsiavash, and Thomas Brox. Coot: Cooperative hierarchical transformer for video-text representation learning. *Advances in neural information processing systems*, 33:22605–22618, 2020. [2](#)
- [19] Tengda Han, Weidi Xie, and Andrew Zisserman. Temporal alignment networks for long-term video. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2906–2916, June 2022. [1](#), [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [13](#), [14](#), [15](#)
- [20] Haoyang Huang, Yaobo Liang, Nan Duan, Ming Gong, Linjun Shou, Daxin Jiang, and Ming Zhou.

- Unicoder: A universal language encoder by pre-training with multiple cross-lingual tasks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2485–2494, 2019. [3](#)
- [21] Dohwan Ko, Joonmyung Choi, Juyeon Ko, Shinyeong Noh, Kyoung-Woon On, Eun-Sol Kim, and Hyunwoo J Kim. Video-text representation learning via differentiable weak temporal alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5016–5025, 2022. [2](#), [7](#)
- [22] Mahnaz Koupaee and William Yang Wang. Wikihow: A large scale text summarization dataset. *ArXiv*, abs/1810.09305, 2018. [5](#), [13](#), [14](#)
- [23] 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. [2](#)
- [24] Hilde Kuehne, Juergen Gall, and Thomas Serre. An end-to-end generative framework for video segmentation and recognition. In *Proc. IEEE Winter Applications of Computer Vision Conference (WACV 16)*, Lake Placid, Mar 2016. [2](#)
- [25] Hilde Kuehne, Alexander Richard, and Juergen Gall. Weakly supervised learning of actions from transcripts. *Computer Vision and Image Understanding*, 163:78–89, 2017. [2](#)
- [26] Colin Lea, Rene Vidal, Austin Reiter, and Gregory D Hager. Temporal convolutional networks: A unified approach to action segmentation. In *Computer Vision—ECCV 2016 Workshops: Amsterdam, The Netherlands, October 8-10 and 15-16, 2016, Proceedings, Part III 14*, pages 47–54. Springer, 2016. [2](#)
- [27] Dong-Hyun Lee et al. Pseudo-label: The simple and efficient semi-supervised learning method for deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3, page 896, 2013. [5](#)
- [28] Linjie Li, Yen-Chun Chen, Yu Cheng, Zhe Gan, Licheng Yu, and Jingjing Liu. HERO: Hierarchical encoder for Video+Language omni-representation pre-training. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 2046–2065, Online, Nov. 2020. Association for Computational Linguistics. [3](#)
- [29] Xudong Lin, Fabio Petroni, Gedas Bertasius, Marcus Rohrbach, Shih-Fu Chang, and Lorenzo Torresani. Learning to recognize procedural activities with distant supervision. *arXiv preprint arXiv:2201.10990*, 2022. [2](#), [3](#)
- [30] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *International Conference on Learning Representations*, 2019. [14](#)
- [31] 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](#), [7](#)
- [32] Yiwei Ma, Guohai Xu, Xiaoshuai Sun, Ming Yan, Ji Zhang, and Rongrong Ji. X-CLIP:: End-to-end multi-grained contrastive learning for video-text retrieval. *arXiv preprint arXiv:2207.07285*, 2022. [3](#)
- [33] Yue Ma, Tianyu Yang, Yin Shan, and Xiu Li. Simvtp: Simple video text pre-training with masked autoencoders. *arXiv preprint arXiv:2212.03490*, 2022. [3](#)
- [34] Jonathan Malmaud, Jonathan Huang, Vivek Rathod, Nicholas Johnston, Andrew Rabinovich, and Kevin Murphy. What’s cookin’? interpreting cooking videos using text, speech and vision. In *Proceedings of the 2015 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 143–152, Denver, Colorado, May–June 2015. Association for Computational Linguistics. [3](#)
- [35] 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 *CVPR*, 2020. [3](#), [6](#), [7](#), [14](#)
- [36] Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Ivan Laptev, and Josef Sivic. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2019. [3](#), [5](#), [7](#), [13](#)
- [37] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In C.J. Burges, L. Bottou, M. Welling, Z. Ghahramani, and K.Q. Weinberger, editors, *Advances in Neural Information Processing Systems*, volume 26. Curran Associates, Inc., 2013. [3](#), [6](#), [14](#)
- [38] Yookoon Park, Mahmoud Azab, Seungwhan Moon, Bo Xiong, Florian Metze, Gourab Kundu, and Kirmani Ahmed. Normalized contrastive learning for text-video retrieval. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 248–260, Abu Dhabi, United Arab Emirates, Dec. 2022. Association for Computational Linguistics. [3](#)
- [39] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. In *International conference on*

- machine learning*, pages 8748–8763. PMLR, 2021. 3, 6
- [40] Alexander Richard, Hilde Kuehne, and Juergen Gall. Action sets: Weakly supervised action segmentation without ordering constraints. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 5987–5996, 2018. 2
- [41] Alexander Richard, Hilde Kuehne, Ahsan Iqbal, and Juergen Gall. Neuralnetwork-viterbi: A framework for weakly supervised video learning. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 7386–7395, 2018. 2
- [42] Andrew Rouditchenko, Angie Boggust, David Harwath, Brian Chen, Dhiraj Joshi, Samuel Thomas, Kartik Audhkhasi, Hilde Kuehne, Rameswar Panda, Rogerio Feris, et al. Avlnet: Learning audio-visual language representations from instructional videos. *arXiv preprint arXiv:2006.09199*, 2020. 2
- [43] Ozan Sener, Amir R Zamir, Silvio Savarese, and Ashutosh Saxena. Unsupervised semantic parsing of video collections. In *Proceedings of the IEEE International conference on Computer Vision*, pages 4480–4488, 2015. 3
- [44] Yuhan Shen and Ehsan Elhamifar. Semi-weakly-supervised learning of complex actions from instructional task videos. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3344–3354, 2022. 2
- [45] Yuhan Shen, Lu Wang, and Ehsan Elhamifar. Learning to segment actions from visual and language instructions via differentiable weak sequence alignment. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10156–10165, 2021. 2, 7
- [46] Nina Shvetsova, Brian Chen, Andrew Rouditchenko, Samuel Thomas, Brian Kingsbury, Rogerio S Feris, David Harwath, James Glass, and Hilde Kuehne. Everything at once-multi-modal fusion transformer for video retrieval. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 20020–20029, 2022. 2
- [47] Bharat Singh, Tim K Marks, Michael Jones, Oncel Tuzel, and Ming Shao. A multi-stream bi-directional recurrent neural network for fine-grained action detection. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1961–1970, 2016. 2
- [48] Kihyuk Sohn, David Berthelot, Nicholas Carlini, Zizhao Zhang, Han Zhang, Colin A Raffel, Ekin Dogus Cubuk, Alexey Kurakin, and Chun-Liang Li. Fix-match: Simplifying semi-supervised learning with consistency and confidence. *Advances in neural information processing systems*, 33:596–608, 2020. 5
- [49] Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. MpNet: Masked and permuted pre-training for language understanding. *arXiv preprint arXiv:2004.09297*, 2020. 5, 8
- [50] Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. Vi-bert: Pre-training of generic visual-linguistic representations. In *International Conference on Learning Representations*, 2020. 3
- [51] Chen Sun, Austin Myers, Carl Vondrick, Kevin P. Murphy, and Cordelia Schmid. Videobert: A joint model for video and language representation learning. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 7463–7472, 2019. 2, 3
- [52] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing*, 2019. 3
- [53] 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 *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [54] 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*, 43(9):3138–3153, 2020. 2
- [55] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. VideoMAE: Masked autoencoders are data-efficient learners for self-supervised video pre-training. In *Advances in Neural Information Processing Systems*, 2022. 3
- [56] Yao-Hung Hubert Tsai, Shaojie Bai, Paul Pu Liang, J. Zico Kolter, Louis-Philippe Morency, and Ruslan Salakhutdinov. Multimodal transformer for unaligned multimodal language sequences. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6558–6569, Florence, Italy, July 2019. Association for Computational Linguistics. 3
- [57] Huiyu Wang, Mitesh Kumar Singh, and Lorenzo Torresani. Ego-only: Egocentric action detection without exocentric pretraining, 2023. 6, 15
- [58] Hanlin Wang, Yilu Wu, Sheng Guo, and Limin Wang. Pdpp: Projected diffusion for procedure planning in instructional videos. *arXiv preprint arXiv:2303.14676*, 2023. 2
- [59] Yuchao Wang, Haochen Wang, Yujun Shen, Jingjing Fei, Wei Li, Guoqiang Jin, Liwei Wu, Rui Zhao, and Xinyi Le. Semi-supervised semantic segmentation using unreliable pseudo-labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 4248–4257, 2022. 5

- [60] Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *European Conference on Computer Vision*, 2018. 3, 6, 14
- [61] Hu Xu, Gargi Ghosh, Po-Yao Huang, Dmytro Okhonko, Armen Aghajanyan, Florian Metzger, Luke Zettlemoyer, and Christoph Feichtenhofer. Video-CLIP: Contrastive pre-training for zero-shot video-text understanding. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Online, Nov. 2021. Association for Computational Linguistics. 3, 6, 7, 14
- [62] Peng Xu, Xiatian Zhu, and David A. Clifton. Multimodal learning with transformers: A survey, 2022. 3
- [63] Jianwei Yang, Yonatan Bisk, and Jianfeng Gao. Taco: Token-aware cascade contrastive learning for video-text alignment. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 11542–11552, 2021. 3
- [64] Yuncong Yang, Jiawei Ma, Shiyuan Huang, Long Chen, Xudong Lin, Guangxing Han, and Shih-Fu Chang. Temporal alignment representation with contrastive learning. *arXiv preprint arXiv:2212.13738*, 2022. 2, 7
- [65] Fangqiu Yi, Hongyu Wen, and Tingting Jiang. As-former: Transformer for action segmentation. In *The British Machine Vision Conference (BMVC)*, 2021. 2
- [66] Amir R Zamir, Alexander Sax, Nikhil Cheerla, Rohan Suri, Zhangjie Cao, Jitendra Malik, and Leonidas J Guibas. Robust learning through cross-task consistency. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11197–11206, 2020. 3
- [67] He Zhao, Isma Hadji, Nikita Dvornik, Konstantinos G Derpanis, Richard P Wildes, and Allan D Jepson. P3iv: Probabilistic procedure planning from instructional videos with weak supervision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 2938–2948, 2022. 2
- [68] Yiwu Zhong, Licheng Yu, Yang Bai, Shangwen Li, Xueting Yan, and Yin Li. Learning procedure-aware video representation from instructional videos and their narrations. *arXiv preprint arXiv:2303.17839*, 2023. 2
- [69] Linchao Zhu and Yi Yang. Actbert: Learning global-local video-text representations. pages 8743–8752, 06 2020. 3
- [70] Dimitri Zhukov, Jean-Baptiste Alayrac, Ramazan Gokberk Cinbis, David Fouhey, Ivan Laptev, and Josef Sivic. Cross-task weakly supervised learning from instructional videos. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3537–3545, 2019. 2, 5, 7, 13, 14
- [71] Andrea Ziani, Zicong Fan, Muhammed Kocabas, Sammy Christen, and Otmar Hilliges. Tempclr: Reconstructing hands via time-coherent contrastive learning. In *International Conference on 3D Vision (3DV)*, 2022. 2

This Appendix provides: additional details (annotation procedure, statistics) about the HT-Step dataset that we introduced for evaluating models on step grounding (Section A), additional details for the rest of the datasets that were used for training/evaluation (Section B), implementation details (Section C), qualitative results for step grounding on HT-Step (Section D), additional ablation studies (Section E), and additional details about the evaluation of our models on HTM-Align (Section F).

## A. HT-Step Dataset

Dataset	Step Annot.	# Videos	# Activities	# Steps	# Segments
HowTo100M [36]	✗	1.2M	25k	-	-
HTM-Align [19]	✗	80	80	-	-
CrossTask [70]	✓	4.8k (2.8k)	83 (18)	133	20.9k
HT-Step (val)	✓	600	120	1,204	3,441
HT-Step (test)	✓	600	120	1,242	3,631
wikiHow		-	14k	100k	-

Table 8: Summary statistics for the datasets in our work. For CrossTask, the statistics for primary activities only are shown in parentheses.

In this section we provide details about the creation of the HT-Step benchmark that we used for evaluating our models. This benchmark was designed to provide a high-quality set of step-annotated instructional videos for a plethora of tasks, described in rich, structured language instead of atomic phrases.

**Annotation setup.** We used videos from the HowTo100M dataset; each one of those videos contains a task id label that corresponds to a wikiHow article. This association enabled us to obtain a set of potential step descriptions for every video, directly from the corresponding wikiHow article. We note that this association is noisy, e.g. the video might show a variation of a specific recipe, where some of the steps in the article often do not appear at all, appear partially, are executed in different order, or are repeated multiple times.

**Annotation instructions.** For each video, annotators were provided with the task name (e.g., Make Pumpkin Puree) and the recipe steps from the corresponding [wikiHow article](#). The annotators were asked to watch the whole video and first decide whether it is relevant to the given task – i.e. if at least some of the given steps were visually demonstrated and the task’s end goal was the same (e.g. a specific recipe) – or reject it otherwise. When a video was deemed relevant, annotators were asked to mark all instances of the provided steps with a temporal window. We note that WikiHow articles often contain several variations/methods for completing a given task. For tasks where this was the case, the annotators were asked to select the set of steps corresponding to the variation that best fits every video and only

use those steps for annotating the entire video.

**QA process.** To ensure the quality of the annotations, we followed a rigorous multi-stage Quality Assurance (QA) process: In the first stage, the videos were annotated by a single annotator. These initial annotations were then reviewed by more experienced annotators, who either approved all the annotations on a given video (meaning all the marked steps were correct and no steps were missing) or marked it for redoing with specific comments on which annotations needed fixing and in what way. At the last stage of the QA process, the annotations that were marked as incorrect were redone by third, independent annotators.

**Statistics.** We provide per-activity statistics for the annotations in Table 9. The metrics used, *i.e.* number of unique steps, step and video coverage, are given to provide an understanding of how the number of steps varies between different tasks and how the steps of a task may appear partially in the HowTo100M videos.

**Validation and test (val/test) split.** Overall during the full annotation process, approximately 35% of the videos were rejected as irrelevant to the given tasks. We split the remaining, annotated videos into a validation and a test set, each containing 600 videos, with 5 videos per task. We ensured that our validation set does not contain videos from HTM-Align. In total 87 human annotators manually annotated 1200 videos over 177 tasks: 120 in the validation and 120 in the test set, with 5 videos per task, *i.e.* with 63 tasks overlapping between the two sets.

## B. Datasets Details

We provide a statistics summary for the datasets used for training and evaluation in Table 8.

**HowTo100M (Training).** HowTo100M contains over 1M unique instructional videos, spanning over 24k activities including cooking, DIY, arts and crafts, gardening, personal care, fitness and more. Each instructional video is complemented by the ASR transcription of its audio, which usually contains the real time narration/commentary of the instructor during the activity. We use the “senticified” version of the ASR sentences provided by Han *et al.* [19]. Following Han *et al.* [19] we also train only using the Food & Entertainment subset, which includes a subset of approximately 370k videos.

**wikiHow (Training).** We train using 14,541 cooking tasks from the wikiHow-Dataset [22]. For each task, we generate an ordered list of steps by extracting the step headlines. The HowTo100M dataset was curated using a semi-automatic pipeline that involved searching YouTube with queries based on the titles of wikiHow articles. Consequently there is an almost complete overlap in activities between the two corpora, which makes wikiHow a natural choice for mining step-level articles to associate with

Task	# steps	step coverage	video coverage
Make Zucchini Pancakes	4.0	0.83	0.37
Make a Hearty Stew	3.5	0.82	0.12
Make Beef and Broccoli	3.1	0.78	0.24
Make Coconut Popsicles	3.8	0.76	0.28
Make Yorkshire Pudding	5.3	0.76	0.11
Cook Spaghetti alla Carbonara	4.6	0.73	0.39
Make Vegan Pesto	2.2	0.73	0.15
Make Corn Fritters	6.4	0.72	0.28
Make Buttermilk Fried Chicken	4.2	0.70	0.44
Make a Shrimp Po Boy Sandwich	4.2	0.70	0.27
⋮	⋮	⋮	⋮
Cook Prime Rib	2.6	0.19	0.19
Cure Bacon	2.2	0.18	0.11
Make Dim Sum	4.6	0.18	0.15
Make Vegan Ceviche	2.8	0.17	0.08
Make Lobster Bisque	3.6	0.17	0.28
Make Giblet Gravy	2.8	0.16	0.23
Make Pickled Eggs	4.4	0.16	0.19
Pickle Onions	1.6	0.15	0.12
Cook Rib Eye Roast	2.0	0.12	0.28
Make Pap	2.0	0.12	0.21
<b>Average</b>	<b>4.0</b>	<b>0.42</b>	<b>0.24</b>

Table 9: Statistics of the annotations used to create the HT-Step benchmark. The metrics are computed per task (for 177 tasks in total), averaged over all the annotated videos for a given task. **# steps** denotes the average number of unique steps annotated per video, per activity; **step coverage** denotes the fraction of a task’s steps that have been found and annotated in every video; **video coverage** denotes the fraction of the video’s duration that is covered by step annotations; Rows are sorted by step coverage; only the 10 tasks with the highest and lowest step coverage are shown here for brevity.

instructions in HowTo100M videos. In the context of this paper we used the wikiHow-Dataset [22] to collect the articles for 14,541 cooking tasks.

**CrossTask (Evaluation).** We use this established instructional video benchmark for *zero-shot* grounding, i.e., by directly evaluating on CrossTask our model learned from HowTo100M. The Crosstask dataset [70], is an established benchmark for temporal localization of steps in instructional videos. It consists of 4800 videos from 83 activities, which are divided into 18 primary (14 related to cooking and 4 to DIY car repairs and shelf assembly) and 65 related activities. The videos in the primary activities are annotated with step annotations in the form of temporal segments from a predefined taxonomy of 133 steps. Those steps tend to be atomic, e.g. for activity “Make Taco Salad” the available steps are “add onion”, “add taco”, “add lettuce”, “add meat”, “add tomato”, “add cheese”, “stir”, and “add tortilla”. Following common practices, we use two evaluation protocols: the first one – *step localization* – aims at predicting a single timestamp for each occurring step in videos from 18 primary tasks [70]. Performance is evaluated by computing the recall (denoted as Avg. R@1)

of the most confident prediction for each task and averaging the results over all query steps in a video, where R@1 measures whether the predicted timestamp for a step falls within the ground truth boundaries. We report average results over 20 random sets of 1850 videos [70]. The second task – *article grounding* – requires predicting temporal segments for each step of an instructional article describing the task represented in the video. We use the mapping between CrossTask and *simplified* wikiHow article steps provided in Chen et al. [12] and report results on 2407 videos of 15 primary tasks obtained excluding three primary tasks following the protocol of [12]. Performance for this task is measured with Recall@K at different IoU thresholds [12].

**HTM-Align (Evaluation).** This benchmark is used to evaluate our model on narration grounding. It contains 80 videos where the ASR transcriptions have been manually aligned temporally with the video. In the main submission, we report the R@1 metric [19], which evaluates whether the model can correctly localize the narrations that are alignable with the video. In Section F we also evaluate our model in terms of its capability to decide whether a narration is visually groundable in the video or not using the ROC-AUC metric [19]. AUC denotes the area the ROC curve of the alignment task, and measures the ability of the model to correctly predict whether a given step is alignable within a video or not.

### C. Implementation Details

As video encoder we adopt the S3D [60] backbone pre-trained with the MIL-NCE objective on HowTo100M [35]. Following previous work [19, 61], we keep this module frozen and use it to extract clip-level features (one feature per second for video decoded at 16 fps). For extracting context-aware features for each sentence (step or narration), we follow the Bag-of-word (BoW) approach based on Word2Vec embeddings [37]. These embeddings are initialized based on MIL-NCE Word2Vec and are fine-tuned during training.

The hyperparameters of the model compared with state-of-the-art methods in Tables 1,2,3 of the main submission were selected based on R@1 performance on the HT-Step validation set and are:  $\lambda_{SV} = \lambda_{NV} = 1$ , temperatures  $\eta, \xi = 0.07$ , and pseudo-label filtering threshold  $\gamma = 0.65$ . We train our model for 12 epochs, with 3 epochs burn-in training with step pseudo-labels generated by TAN, and then we update the teacher VINA every 3 epochs. We use the AdamW [30] optimizer, having an initial learning rate of  $2e - 4$  decayed with a cosine learning schedule. Our batch size is 32 videos, with maximum length of 1024 seconds.

Pseudo-labels are obtained based on the steps-to-video alignment matrix and are generated (before filtering) as follows: for each step we find the timestep with maximum similarity with the step and then extend a temporal segment

to the left and right of that peak as long as the similarity score does not follow below 0.7 of the peak height. Pseudo-labels whose peak score falls below the filtering threshold  $\gamma$  are not used for training.

The rest of hyperparameters were selected based on TAN [19]. The multimodal encoder is a pre-norm multi-layer transformer which consists of 6 layers of self-attention, with 8 heads and has hidden dimension  $D = 512$ . A learnable positional encoding of size  $D = 512$  is used to inject temporal information to each frame/narration/step token.

To obtain temporal segment detections from the step-to-video alignment output of our model (e.g. for evaluating on the CrossTask article grounding setting or for the qualitative video included in this supplementary) we use a simple 1D blob detector [57]. Unless otherwise specified, we use the fused alignment matrix for step grounding when narrations are available during inference time.

Our model is trained on 8 GPUs (Tesla V100-SXM2-32GB) and training lasts approximately 10-12 hours. All models were implemented in Python using Pytorch and are based on the PySlowFast (<https://github.com/facebookresearch/SlowFast>) and TAN (<https://github.com/TengdaHan/TemporalAlignNet>) open-source codebases. For ablation studies, we choose the best checkpoint for each configuration based on performance on HT-Step validation set and report its test split performance.

## D. Qualitative Results

In this section, we provide qualitative results for the ground-truth steps-to-video alignment and predicted alignments by our improved baseline that serves as the initial teacher model (TAN\*), and our model (using the direct steps-to-video alignment without narrations) or the fusion with the indirect steps-to-video alignment (with narrations). From these qualitative results (Figure 3), we observe that our VINA model can correctly temporally localize visually groundable steps, despite being trained only with noisy pairs of narrated videos and instructional steps. Predicted alignments tend to also be less noisy than TAN\*, showcasing the effectiveness of training a video-language alignment model with distant supervision from WikiHow articles. Our model can also leverage ASR transcripts (without any temporal information regarding when the instructor uttered each narration) to further improve its results (Figure 4).

## E. Extra Ablations

**Architecture ablations.** In Table 10 we study the design of the unimodal encoder used to embed steps before they are fed to our Multimodal Transformer. Overall, using po-

sitional embeddings capturing the ordering of steps in a task, and using modality-specific projection MLPs leads to a slightly better performance in step grounding (w/o narration input). Narration grounding seems to benefit from using a shared text encoder, possibly because this facilitates knowledge transfer from the WikiHow steps.

PE	Sep. MLP	HT-Step $\uparrow$ R@1		HTM-Align
		w/o nar.	w/ nar.	
	✓	33.5	34.0	65.8
		34.0	34.9	65.9
✓		33.8	34.4	<b>67.0</b>
✓	✓	<b>34.3</b>	<b>36.1</b>	64.8

Table 10: **Ablation study on architecture design.** We study the contribution of positional encodings for steps (*PE*) and of specialized text projection layers for wikiHow article steps (*Sep. MLP*). All models are trained for joint narration and step grounding with fixed pseudo-labels from TAN and evaluated on HT-Step val split (last row corresponds to row 5 in Table 4 of the main text).

## F. Experimental Setup on HTM-Align

As explained in the official code repository of TAN [19] ([https://github.com/TengdaHan/TemporalAlignNet/tree/main/htm\\_align](https://github.com/TengdaHan/TemporalAlignNet/tree/main/htm_align)), the results reported for HTM-Align are obtained with a text moving window of 1 minute, i.e., for each 1-minute temporal segment only ASR captions whose original time-stamps fall within a 3-min window centered around this temporal segment are considered for grounding. Instead, for all our reported results (for TAN\* and VINA) we operate in the more challenging setup where an ASR caption can be grounded in any timestep of the original video (there is no knowledge about the original ASR timestamps during inference). Under this more challenging setup, our model outperforms TAN both in narration retrieval, as measured by Recall@1 (66.5% vs 49.4%, as seen in Table 1 of the main submission).

Our model also performs comparably with TAN in step alignability prediction, as measured by ROC-AUC (76% vs 75.1%). Note that our model does not have dedicated alignability head for predicting whether a narration exists or not in the video as TAN [19]. Instead, we simply obtain an alignability score by using the maximum cosine similarity score over time, where cosine similarities of each narration with each video frame are computed based on the outputs of the unimodal encoders.

## G. Limitations and Ethical Concerns

From the qualitative results, we observe that due to the losses used during training, which do not explicitly penalize

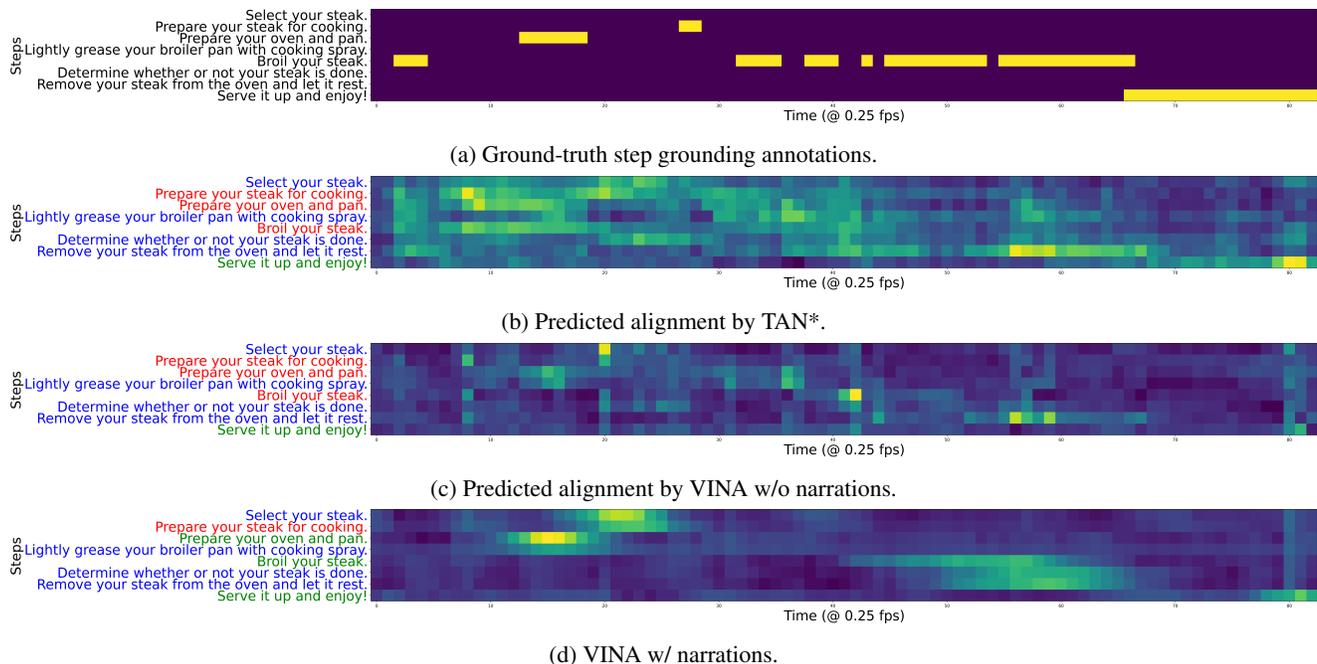
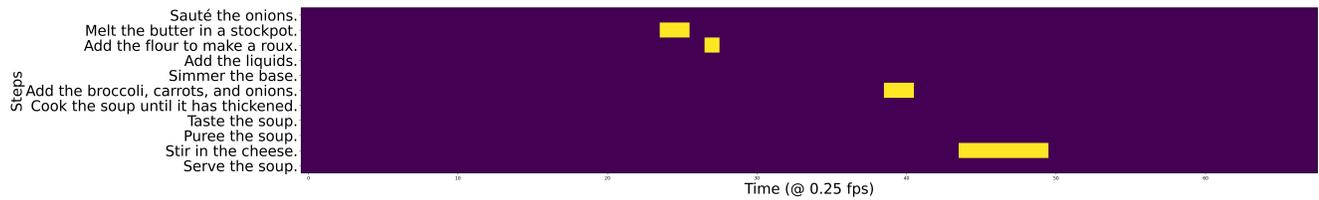
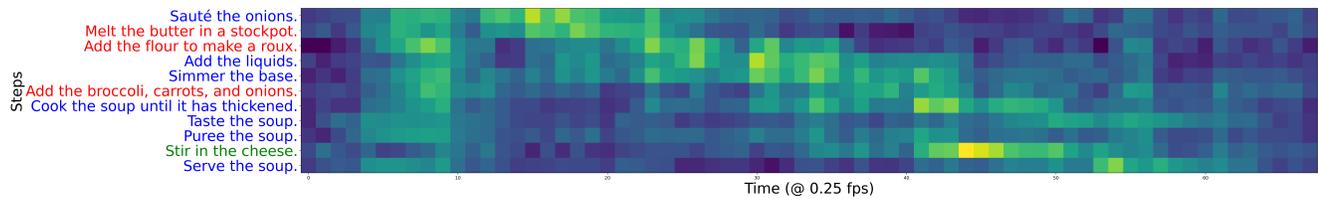


Figure 3: Qualitative results on a validation video of the HT-Step dataset (VIQYQkA3mNU) demonstrating how to *Broil Steak*. Steps that are not visually groundable in the video are highlighted in blue, steps that are correctly retrieved by each model are highlighted in green, while steps that are not retrieved are shown in red. Figure best viewed zoomed in and in color.

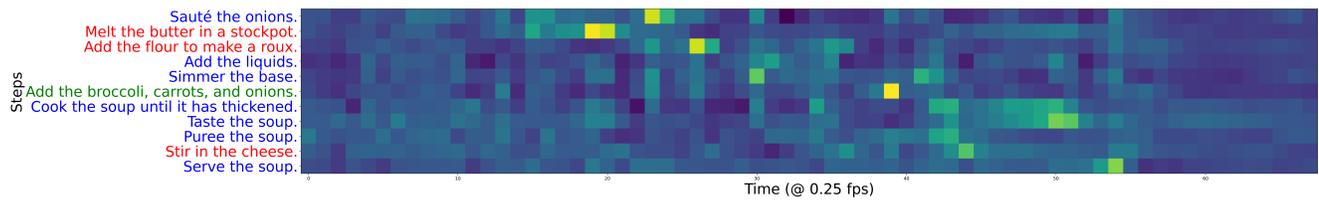
wrong temporal extent (as long as the predicted heatmap has a peak within the target temporal window), grounded temporal segments tend to be short. This is especially prominent when using the direct steps-to-videos alignment that is explicitly supervised (second to last row of the predicted alignment figures). Furthermore, our training objective does not utilize negative examples, e.g. steps that are not visually groundable, to suppress detections. This can lead to confident detections for missing steps. Another limitation of our approach (similar to previous approaches that operate on pre-extracted visual features) is that our performance is limited by the quality of the extracted visual representations. Regarding ethical concerns, public instructional video datasets and public knowledge base datasets may have gender, age, geographical and cultural bias.



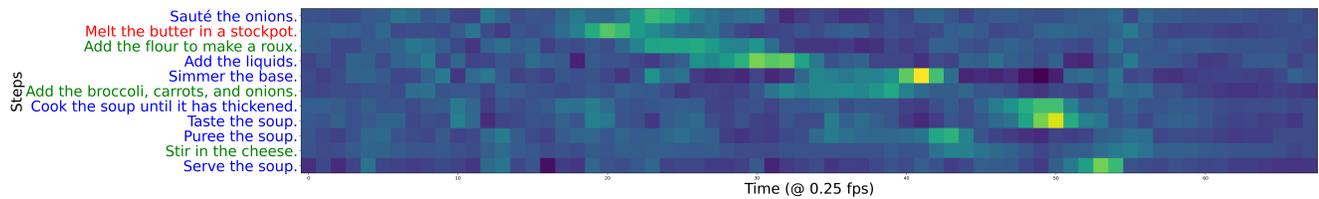
(a) Ground-truth step grounding annotations.



(b) Predicted alignment by TAN\*.



(c) Predicted alignment by VINA w/o narrations.



(d) VINA w/ narrations.

Figure 4: Qualitative results on a validation video of the HTM-Step dataset (0dHofx11qAg) demonstrating how to *Make Broccoli Cheese Soup*. Steps that are not visually groundable in the video are highlighted in blue, steps that are correctly retrieved by each model are highlighted in green, while steps that are not retrieved are shown in red. Figure best viewed zoomed in and in color.