

Multimodal Pretraining for Dense Video Captioning

Gabriel Huang^{1*}, Bo Pang², Zhenhai Zhu², Clara Rivera², Radu Soricut²

¹Mila & University of Montreal

²Google Research

gabriel.huang@umontreal.ca

{bopang, zhenhai, rivera, rsoricut}@google.com

Abstract

Learning specific hands-on skills such as cooking, car maintenance, and home repairs increasingly happens via instructional videos. The user experience with such videos is known to be improved by meta-information such as time-stamped annotations for the main steps involved. Generating such annotations automatically is challenging, and we describe here two relevant contributions. First, we construct and release a new dense video captioning dataset, **Video Timeline Tags (ViTT)**, featuring a variety of instructional videos together with time-stamped annotations. Second, we explore several multimodal sequence-to-sequence pretraining strategies that leverage large unsupervised datasets of videos and caption-like texts. We pretrain and subsequently finetune dense video captioning models using both YouCook2 and ViTT. We show that such models generalize well and are robust over a wide variety of instructional videos.

1 Introduction

YouTube recently reported that a billion hours of videos were being watched on the platform every day (YouTubeBlog, 2017). In addition, the amount of time people spent watching online videos was estimated to grow at an average rate of 32% a year between 2013 and 2018, with an average person forecasted to watch 100 minutes of online videos per day in 2021 (ZenithMedia, 2019).

An important reason for this fast-growing video consumption is information-seeking. For instance, people turn to YouTube “hungry for how-to and learning content” (O’Neil-Hart, 2018). Indeed, compared to traditional content format such as text, video carries richer information to satisfy such

*This work was done while Gabriel Huang was an intern at Google Research.



Groundtruth *Varying stitching speeds*
Ø-Pretraining *Showing other parts*
MASS-Pretraining *Explaining how to do a stitch*

Figure 1: Dense video captioning using ViTT-trained models. For the given video scene, we show the ViTT annotation (Groundtruth) and model outputs (no pretraining and MASS-based pretraining).

needs. But as a content media, videos are also inherently more difficult to skim through, making it harder to quickly target the relevant part(s) of a video. Recognizing this difficulty, search engines started showing links to “key moments” within videos in search results, based on timestamps and short descriptions provided by the content creators themselves.¹ This enables users to get a quick sense of what the video covers, and also to jump to a particular time in the video if so desired. This effort echoes prior work in the literature showing how users of instructional videos can benefit from human-curated meta-data, such as a timeline pointing to the successive steps of a tutorial (Kim et al., 2014; Margulieux et al., 2012; Weir et al., 2015). Producing such meta-data in an automatic way would greatly scale up the efforts of providing easier information access to videos. This task is closely related to the dense video captioning task considered in prior work (Zhou et al., 2018a,c; Krishna et al., 2017), where an instructional video is first segmented into its main steps, followed by segment-level caption generation.

To date, the YouCook2 data set (Zhou et al., 2018a) is the largest annotated data set for dense

¹<https://www.blog.google/products/search/key-moments-video-search/>

video captioning. It contains annotations for 2,000 cooking videos covering 89 recipes, with per-recipe training / validation split. Restricting to a small number of recipes is helpful for early exploratory work, but such restrictions impose barriers to model generalization and adoption that are hard to overcome. We directly address this problem by constructing a larger and broader-coverage annotated dataset that covers a wide range of instructional topics (cooking, repairs, maintenance, etc.) We make the results of our annotation efforts publicly available as Video Timeline Tags (ViTT)², consisting of around 8,000 videos annotated with timelines (on average 7.1 segments per video, each segment with a short free-text description).

Using YouCook2 and the new ViTT dataset as benchmarks for testing model performance and generalization, we further focus on the sub-problem of video-segment-level caption generation, assuming segment boundaries are given (Hessel et al., 2019; Sun et al., 2019b; Luo et al., 2020). Motivated by the high cost of collecting human annotations, we investigate pretraining a video segment captioning model using unsupervised signals – ASR (Automatic Speech Recognition) tokens and visual features from instructional videos, and *unpaired* instruction steps extracted from independent sources: Recipe1M (Marin et al., 2019) and WikiHow (Koupae and Wang, 2018). In contrast to prior work that focused on BERT-style pretraining of encoder networks (Sun et al., 2019b,a), our approach entails jointly pretraining both multimodal encoder and text-based decoder models via MASS-style pretraining (Song et al., 2019). Our experiments show that pretraining with either text-only or multi-modal data provides significant gains over no pretraining, on both the established YouCook2 benchmark and the new ViTT benchmark. The results we obtain establish state-of-the-art performance on YouCook2, and present strong performance numbers on the ViTT benchmark. These findings help us conclude that the resulting models generalize well and are quite robust over a wide variety of instructional videos.

2 Related Work

Text-only Pretraining. Language pretraining models based on the Transformer neural net-

²Available at <https://github.com/google-research-datasets/Video-Timeline-Tags-ViTT>

work architecture (Vaswani et al., 2017a) such as BERT (Devlin et al., 2018), GPT (Radford et al., 2018), RoBERTa (Liu et al., 2019), MASS (Song et al., 2019) and ALBERT (Lan et al., 2020) have achieved state-of-the-art results on many NLP tasks. MASS (Song et al., 2019) has been recently proposed as a joint encoder-decoder pretraining strategy. For sequence-to-sequence tasks, this strategy is shown to outperform approaches that separately pretrain the encoder (using a BERT-style objective) and the decoder (using a language modeling objective). UniLM (Dong et al., 2019), BART (Lewis et al., 2019), and T5 (Raffel et al., 2019) propose unified pretraining approaches for both understanding and generation tasks.

Multimodal Pretraining. VideoBERT (Sun et al., 2019b), CBT (Sun et al., 2019a) and ActBERT (Zhu and Yang, 2020) use a BERT-style objective to train both video and ASR text encoders. Alayrac et al. (2016) and Miech et al. (2020) use margin-based loss functions to learn joint representations for video and ASR, and evaluate them on downstream tasks such as video captioning, action segmentation and anticipation, and action localization. An independent and concurrent work (UniViLM) by Luo et al. (2020) is closely related to ours in that we share some similar pretraining objectives, some of the pretraining setup – HowTo100M (Alayrac et al., 2016), and the down-stream video captioning benchmark using YouCook2 (Zhou et al., 2018a). The main difference is that they use BERT-style pretraining for encoder and language-modeling style pretraining for decoder, whereas we use MASS-style pre-training to pretrain encoder and decoder jointly.

Other approaches such as ViLBERT (Lu et al., 2019), LXMERT (Tan and Bansal, 2019), Unicoder-VL (Li et al., 2019), VL-BERT (Su et al., 2019), and UNITER (Chen et al., 2019) focus on pretraining joint representations for text and image, evaluating them on downstream tasks such as visual question answering, image-text retrieval and referring expressions.

Dense Video Captioning. In this paper, we focus on generating captions at the segment-level, which is a sub-task of the so-called dense video captioning task (Krishna et al., 2017), where fine-grained captions are generated for video segments, conditioned on an input video with pre-defined

Name	Type	# segments
<i>Pretraining datasets</i>		
YT8M-cook	ASR+video	186 K
HowTo100M	ASR+video	8.0 M
Recipe1M	CAP-style	10.8 M
WikiHow	CAP-style	1.3 M
<i>Finetuning datasets</i>		
YouCook2	ASR+video+CAP	11.5 K
ViTT-All	ASR+video+CAP	88.5 K

Table 1: Datasets used in this work, along with size of the data measured by the total number of segments.

event segments. This is different from the video captioning models that generate a single summary for the entire video (Wang et al., 2019).

Hessel et al. (2019) make use of ASR and video for segment-level captioning on YouCook2 and show that most of the performance comes from ASR. Shi et al. (2019); Luo et al. (2020) train their dense video captioning models on both video frames and ASR text and demonstrate the benefits of adding ASR as an input to the model. There are also a number of video captioning approaches that do not use ASR directly (Zhou et al., 2018c; Pan et al., 2020; Zheng et al., 2020; Zhang et al., 2020; Lei et al., 2020).

Instructional video captioning data sets. In addition to YouCook2 (Zhou et al., 2018a), there are two other smaller data sets in the instructional video captioning category. Epic Kitchen (Damen et al., 2018) features 55 hours of video consisting of 11.5M frames, which were densely labeled for a total of 39.6K action segments and 454.3K object bounding boxes. How2 (Sanabria et al., 2018) consists of instructional videos with video-level (as opposed to segment-level) descriptions, authored by the video creators themselves.

3 Data

We present the datasets used for pretraining, fine-tuning, and evaluation in Table 1. We also describe in detail the newly introduced dense video captioning dataset, **V**ideo **T**imeline **T**ags (ViTT).

3.1 Dense Video-Captioning Datasets

Our goal is to generate captions (CAP) for video segments. We consider two datasets with segment-level captions for fine-tuning and evaluating ASR+Video→CAP models.

YouCook2. Up to this point, YouCook2 (Zhou et al., 2018a) has been the largest human-annotated dense-captioning dataset of instructional videos publicly available. It originally contained 2,000 cooking videos from YouTube. Starting from 110 recipe types (e.g., “shrimp tempura”), 25 unique videos per recipe type were collected; the recipe types that did not gather enough videos were dropped, resulting in a total of 89 recipe types in the final dataset. In addition, Zhou et al. (2018b) “randomly split the videos belonging to each recipe into 67%:23%:10% as training, validation and test sets³,” which effectively guarantees that videos in the validation and test sets are never about unseen recipes. Annotators were then asked to construct recipe steps for each video — that is, identify the start and end times for each step, and provide a recipe-like description of each step. Overall, they reported an average of 7.7 segments per video, and 8.8 words per description. After removing videos that had been deleted by users, we obtained a total of 11,549 segments.

ViTT. One limitation of the YouCook2 dataset is the artificially imposed (almost) uniform distribution of videos over 89 recipes. While this may help making the task more tractable, it is difficult to judge whether performance on its validation / test sets can be generalized to unseen topics.

The design of our ViTT dataset annotation process is aimed at fixing some of these drawbacks. We started by collecting a large dataset of videos containing a broader variety of topics to better reflect topic distribution in the wild. Specifically, we randomly sampled instructional videos from the YouTube-8M dataset (Abu-El-Haija et al., 2016), a large-scale collection of YouTube videos that also contain topical labels. Since much of prior work in this area revolved around cooking videos, we aimed at sampling a significant proportion of our data from videos with cooking labels (specifically, “Cooking” and “Recipe”). Aside from the intentional bias regarding cooking videos, the rest of the videos were selected by randomly sampling non-cooking videos, including only those that were considered to be instructional videos by our human annotators.

Once candidate videos were identified, timeline annotations and descriptive tags were collected.

³Note that no annotations are provided for the test split; we conducted our own training/dev/test split over available videos.

Our motivation was to enable downstream applications to allow navigating to specific content sections. Therefore, annotators were asked to identify the main steps in a video and mark their start time. They were also asked to produce a descriptive-yet-concise, free-text tag for each step (e.g., “shaping the cookies”, “removing any leftover glass”). A subset of the videos has received more than one complete annotation (main steps plus tags).

The resulting ViTT dataset consists of a total of 8,169 videos, of which 3,381 are cooking-related. A total of 5,840 videos have received only one annotation, and have been designated as the training split. Videos with more than one annotation have been designated as validation / test data. Overall, there are 7.1 segments per video on average (max: 19). Given the dataset design, descriptions are much shorter in length compared to YouCook2: on average there are 2.97 words per tag (max: 16) — 20% of the captions are single-word, 22% are two-words, and 25% are three words. Note that the average caption length is significantly shorter than for YouCook2, which is not surprising given our motivation of providing short and concise timeline tags for video navigation. We standardized the phrases among the top-20 most frequent captions. For instance, {“intro”, “introduction”} → “intro”. Otherwise, we have preserved the original tags as-is, even though additional paraphrasing most definitely exists. Annotators were instructed to start and end the video with an opening and closing segment as possible. As a result, the most frequent tag (post-standardization) in the dataset is “intro”, which accounts for roughly 11% of the 88,455 segments. More details on the data collection process and additional analysis can be found in the Supplementary Material (Section A.1).

Overall, this results in 56,027 unique tags, with a vocabulary size of 12,509 token types over 88,455 segments. In this paper, we consider two variants: the full dataset (ViTT-All), and the cooking subset (ViTT-Cooking).

3.2 Pretraining Datasets: ASR+Video

We consider two large-scale unannotated video datasets for pretraining, as described below. Timestamped ASR tokens were obtained via YouTube Data API,⁴ and split into ASR segments if the timestamps of two consecutive words are more

⁴<https://developers.google.com/youtube/v3/docs/captions>

than 2 seconds apart, or if a segment is longer than a pre-specified max length (in our case, 320 words). They were paired with concurrent video frames in the same segment.

YT8M-cook We extract the cooking subset of YouTube-8M (Abu-El-Haija et al., 2016) by taking, from its training split, videos with “Cooking” or “Recipe” labels, and retain those with English ASR, subject to YouTube policies. After preprocessing, we obtain 186K ASR+video segments with an average length of 64 words (24 seconds) per segment.

HowTo100M. This is based on the 1.2M YouTube instructional videos released by Miech et al. (2019), covering a broad range of topics. After preprocessing, we obtain 7.99M ASR+video segments with an average of 78 words (28.7 seconds) per segment.

3.3 Pretraining Datasets: CAP-style

We also consider two text-only datasets for *pre-training*, containing *unpaired* instruction steps similar in style to the target captions.

Recipe1M is a collection of 1M recipes scraped from a number of popular cooking websites (Marin et al., 2019). We use the sequence of instructions extracted for each recipe in this dataset, and treat each recipe step as a separate example during pretraining. This results in 10,767,594 CAP-style segments, with 12.8 words per segment.

WikiHow is a collection of 230,843 articles extracted from the WikiHow knowledge base (Koupae and Wang, 2018). Each article comes with a title starting with “How to”. Each associated step starts with a step summary (in bold) followed by a detailed explanation. We extract the all step summaries, resulting in 1,360,145 CAP-style segments, with 8.2 words per segment. Again, each instruction step is considered as a separate example during pretraining.

3.4 Differences between Pretraining and Finetuning Datasets

First, note that *video segments* are defined differently for pretraining and finetuning datasets, and may not match exactly. For ASR+Video pretraining datasets, which are unsupervised, the segments are divided following a simple heuristic (e.g., two consecutive words more than 2 seconds apart), whereas for finetuning ASR+Video→CAP datasets, which are supervised, the segments are defined by

human annotators to correspond to instruction steps. Otherwise, the ASR data are relatively similar between pretraining and finetuning datasets, as both come from instructional videos and are in the style of spoken language.

Second, compared to the target captions in finetuning datasets, the CAP-like pretraining datasets are similar in spirit — they all represent summaries of *steps*, but they may differ in length, style and granularity. In particular, the CAP-like pretraining datasets are closer in style to captions in YouCook2, where annotators were instructed to produce a recipe-like description for each step. This is reflected in their similar average length (YouCook2: 8.8 words, Recipe1M: 12.8 words, WikiHow: 8.2 words); whereas captions in ViTT are significantly shorter (2.97 words on average).

Despite these differences — some are inevitable due to the unsupervised nature of pretraining datasets — the pretraining data is very helpful for our task as shown in the experimental results.

4 Method

To model segment-level caption generation, we adopt MASS-style pretraining (Song et al., 2019) with Transformer (Vaswani et al., 2017b) as the backbone architecture. For both pre-training and fine-tuning objectives, we have considered two variants: text-only and multi-modal. They are summarized in Table 2 and more details are given below.

4.1 Separate-Modality Architecture

Both ASR tokens and video segment features are given as input in the multimodal variants. We consider an architecture with a separate transformer for each modality (text or video), see Figure 2 for details. When available, the text and video encoders attend to each other at every layer using cross-modal attention, as in ViLBERT (Lu et al., 2019). The text decoder attends over the final-layer output of both encoders. We discuss in more detail the differences between using a separate-modality architecture vs. a vanilla-Transformer approach for all modalities in Appendix A.2.

The inputs to the text encoder is the sum of three components: text token embeddings, positional embeddings and the corresponding style embeddings,⁵ depending on the style of the text (ASR or Caption-like). The inputs to the video encoder

⁵This is similar to the way language-ID embeddings are used in machine translation.

could be either precomputed frame-level 2D CNN features or 3D CNN features, pretrained on the Kinetics (Carreira and Zisserman, 2017; Kay et al., 2017) data set. The visual features are projected with fully-connected layers to the same dimension as the text embeddings.

The main architecture we consider is a 2-layer encoder (E2), 6-layer decoder (D6) Transformer. We use **E2D6** to refer to the text-only version, and **E2vidD6** to refer to the multimodal version with an active video encoder. We also experiment with E2D2 and E2vidD2 (2-layer decoder).⁶

4.2 Pretraining with Text-only MASS

Text-only pretraining is essentially the unsupervised learning of the style transfer between ASR-style and caption-style texts using *unpaired* data sources: ASR strings from video segments in YT8M-cook or HowTo100M; and CAP-style instruction steps found in Recipe1M or HowTo100M. Just like using MASS for unsupervised machine translation involves pretraining the model on unpaired monolingual datasets, we alternate between ASR→ASR and CAP→CAP MASS steps during our pretraining stage, which does not require the “source” (ASR+Video) and “target” (CAP-style) data to be aligned.

In an ASR→ASR step, we mask a random subsequence of the ASR and feed the masked ASR to the text encoder. The text decoder must reconstruct the hidden subsequence while attending to the encoder output. A CAP→CAP step works similarly by trying to reconstruct a masked sequence of a CAP-style text. The encoder and decoder are trained jointly using teacher-forcing on the decoder. We denote this text-only strategy as **MASS** in the experiments.

4.3 Pretraining with Multimodal MASS

During multimodal pretraining, we alternate between text-only CAP→CAP MASS steps and multimodal MASS steps. During each multimodal MASS step ASR+video→ASR, we feed a masked ASR to the text-encoder and the co-occurring video features to the video-encoder. The text decoder must reconstruct the masked ASR subsequence. We denote this pretraining strategy as **MASSvid** in the experiments. This trains cross-modal attention between the text-encoder and video-encoder

⁶We found in a preliminary study that using 6-layer encoders did not improve performance for our application.

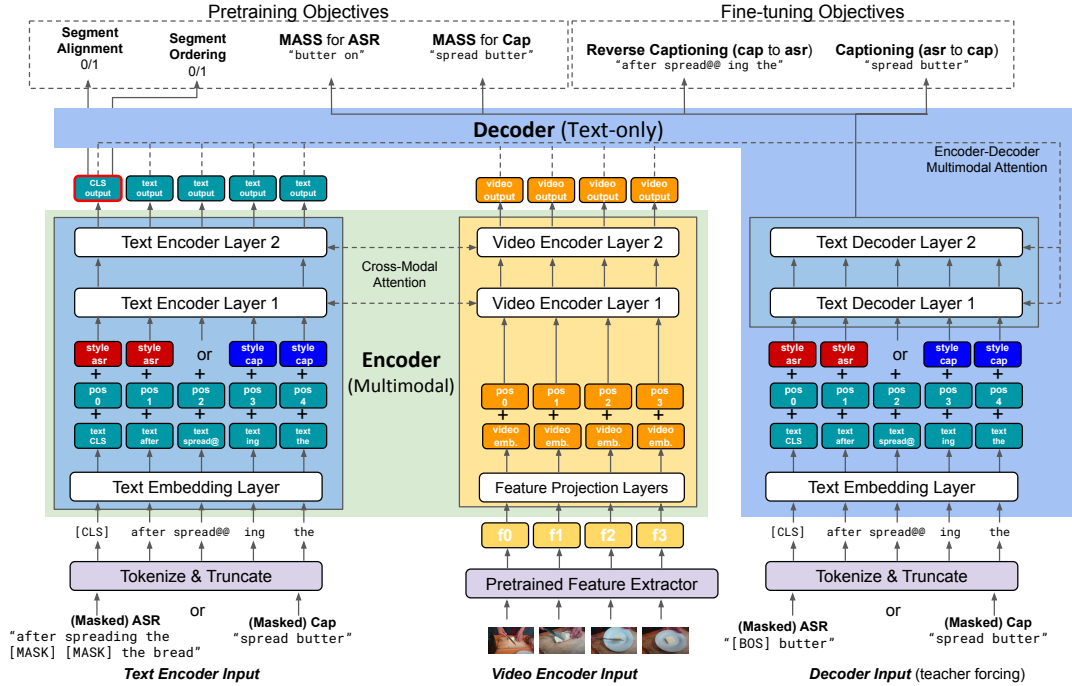


Figure 2: A diagram for the separate-modality architecture. It consists of a two-stream (text and video inputs) encoder with cross-modal attention and a text-only decoder, jointly trained using the MASS objective.

at every layer, jointly with the text decoder that attends to the output layer of both the text and video encoders.⁷

To force more cross-modal attention between encoder and decoder, we also investigate a strategy of hiding the text-encoder output from the decoder for some fraction of training examples. We refer to this strategy as **MASSdrop** in the experiments.

4.4 Pretraining with Alignment and Ordering Tasks

We also explore encoder-only multimodal pretraining strategies. We take the last-layer representation for the CLS (beginning of sentence) token from the encoder, and add a multi-layer perceptron on top of it for binary predictions (Figure 2). Given a pair of ASR and video segment, we train the encoder to predict the following objectives:

- **Segment-Level Alignment.** An (ASR, video) pair is *aligned* if they occur in the same pre-training segment; negative examples are constructed by sampling pairs from the same video but at least 2 segments away.

⁷In preliminary experiments, we had attempted to directly adapt the MASS objective (Song et al., 2019) to video reconstruction — by masking a subsequence of the input video and making the video decoder reconstruct the input using the Noise Contrastive Estimator Loss (Sun et al., 2019a). Due to limited success, we did not further pursue this approach.

- **Segment-Level Ordering.** We sample (ASR, video) pairs that are at least 2 segments away, and train the model to predict whether the ASR occurs before or after the video clip.

During this **MASSalign** pretraining stage, we alternate between two text-only MASS steps (CAP→CAP, ASR→ASR) and the two binary predictions (**Alignment** and **Ordering**) described above.

4.5 Finetuning on Video Captioning

For text-only finetuning, we feed ASR to the text encoder and the decoder has to predict the corresponding CAP (ASR→CAP). For multimodal finetuning, we also feed additional video representations to the video encoder (ASR+video→CAP). When finetuning a multimodal model from text-only pretraining, everything related to video (weights in the video encoder and any cross-modal attention modules) will be initialized randomly. In addition to these *uni-directional* (**UniD**) finetuning, we also experiment with several variants of *bi-directional* (**BiD**) finetuning (Table 2). For instance, adding CAP→ASR (predicting ASR from CAP) to text-only finetuning. In the experiments, we find some variants of bidirectional finetuning beneficial whether training from scratch or finetuning from a pretrained model.

Pretraining Objectives			
Name	T	V	Input→Output
MASS	✓	✗	CAP→CAP, ASR→ASR
MASSvid	✓	✓	CAP→CAP, ASR+video→ASR
MASSdrop	✓	✓	CAP→CAP, ASR+video→ASR
MASSalign	✓	✓	CAP→CAP, ASR→ASR, ASR+video→{0, 1}
Finetuning Objectives			
Name	T	V	Input→Output
UniD	✓	✗	ASR→CAP
BiD	✓	✗	ASR→CAP, CAP→ASR
UniD	✓	✓	ASR+video→CAP
BiD	✓	✓	ASR+video→CAP, CAP→ASR
BiDalt	✓	✓	ASR+video→CAP, CAP+video→ASR

Table 2: Pretraining and Fine-tuning objectives. For each strategy, ✓ indicates whether the text (T) and video (V) encoders are active, followed by a summary of training objectives involved in one training step.

5 Experiments

5.1 Implementation Details

We tokenize ASR and CAP inputs using byte-pair-encoding subwords (Sennrich et al., 2015), and truncate them to 240 subwords. We truncate video sequences to 40 frames (40 seconds of video), compute the 128-dim features proposed by Wang et al. (2014) (which we will refer to as Compact 2D features), and project them to the embedding space using a two-layer perceptron with layer normalization and GeLU activations.

We instantiate the E2xDx models defined in Section 4.1 with 128-dimensional embeddings and 8 heads respectively for self-attention, encoder-decoder, and cross-modal attention modules. We define each epoch to be 3,125 iterations, where each iteration contains one repetition of each training step as represented in Table 2. We pretrain for 200 epochs and finetune for 30 epochs.

For evaluation, we consider BLEU-4 (Papineni et al., 2002), METEOR (Denkowski and Lavie, 2014), ROUGE-L (Lin and Och, 2004) and CIDEr (Vedantam et al., 2015) metrics.

Please refer to Appendix A.3 for full implementation details, hyperparameters and computation cost.

Notes on ViTT evaluation: With the exception of ROUGE-L, all other metrics are sensitive to short groundtruth. 67% of the groundtruth tags in ViTT have less than 4 words, where a perfect prediction will not yield a full score in, say, BLEU-4. Thus, we

focus mainly on ROUGE-L, report BLEU-1 instead of BLEU-4 for ViTT, and provide the other two metrics only as reference points.

We had originally decided to use videos with multiple annotations as validation and test data, so that we could explore evaluation with multiple reference groundtruth captions. But as annotators do not always yield the same set of segment boundaries, this became tricky. Instead, we simply treat each segment as a separate instance with one single reference caption. Note that all segments annotated for the same video will be in either validation or test to ensure no content overlap.

5.2 Main Results

We run several variants of our method on YouCook2, ViTT-All and ViTT-Cooking, using different architectures, modalities, pretraining datasets, pretraining and finetuning strategies.

Comparing with other methods on YouCook2

For YouCook2, we report our method alongside several methods from the literature (Hessel et al., 2019; Sun et al., 2019b; Zhou et al., 2018c; Lei et al., 2020), as well as state-of-the-art concurrent work (Luo et al., 2020). The related work is provided for reference and to give a ballpark estimate of the relative performance of each method, but results are not always strictly and directly comparable. Beyond the usual sources of discrepancy in data processing, tokenization, or even different splits, an additional source of complication comes from the fact that videos are regularly deleted by content creators, causing video datasets to shrink over time. Additionally, when comparing to other work incorporating pretraining, we could differ in (videos available in) pretraining datasets, segmentation strategies, etc. To this end, we perform an extensive ablation study, which at least helps us to understand the effectiveness of different components in our approach.

Effect of pretraining The main experimental results for the three datasets we consider are summarized in Table 3 (YouCook2) and Table 4 (ViTT-All and ViTT-Cooking). Across all three datasets, the best performance is achieved by finetuning a multimodal captioning model under the *Multimodal Pretraining* condition. For instance, on YouCook2, E2vidD6-MASSvid-BiD improves over the no-pretraining model E2vidD6-BiD by 4.37 ROUGE-L, a larger improvement than UniViLM with pretraining (#5) vs without (#2) (Luo et al., 2020). This

Method	Input	Pretraining	BLEU-4	METEOR	ROUGE-L	CIDER
Constant Pred (Hessel et al., 2019)	-	-	2.70	10.30	21.70	0.15
MART (Lei et al., 2020)	Video	-	8.00	15.90	-	0.36
EMT (Zhou et al., 2018c)	Video	-	4.38	11.55	27.44	0.38
CBT (Sun et al., 2019a)	Video	Kinetics + HowTo100M	5.12	12.97	30.44	0.64
AT (Hessel et al., 2019)	ASR	-	8.55	16.93	35.54	1.06
AT+Video (Hessel et al., 2019)	Video + ASR	-	9.01	17.77	36.65	1.12
UniViLM #1 (Luo et al., 2020)	Video	-	6.06	12.47	31.48	0.64
UniViLM #2 (Luo et al., 2020)	Video + ASR	-	8.67	15.38	35.02	1.00
UniViLM #5 (Luo et al., 2020)	Video + ASR	HowTo100M	10.42	16.93	38.02	1.20
<i>∅ Pretraining</i>						
E2D6-BiD	ASR	-	7.90	15.70	34.86	0.93
E2vidD6-BiD	Video + ASR	-	8.01	16.19	34.66	0.91
<i>Text Pretraining</i>						
E2D6-MASS-BiD	ASR	YT8M-cook + Recipe1M	10.60	17.42	38.08	1.20
E2vidD6-MASS-BiD	Video + ASR	YT8M-cook + Recipe1M	11.47	17.70	38.80	1.25
<i>Multimodal Pretraining</i>						
E2vidD6-MASSalign-BiD	Video + ASR	YT8M-cook + Recipe1M	11.53	17.62	39.03	1.22
E2vidD6-MASSvid-BiD	Video + ASR	YT8M-cook + Recipe1M	12.04	18.32	39.03	1.23
E2vidD6-MASSdrop-BiD	Video + ASR	YT8M-cook + Recipe1M	10.45	17.74	38.82	1.22
Human (Hessel et al., 2019)	Video + ASR	-	15.20	25.90	45.10	3.80

Table 3: Segment-level captioning results on YouCook2. We use YT8M-cook and Recipe1M for pretraining. The numbers for the related work (first group) are directly reported from the corresponding papers. The last line is an estimate of human performance as reported by Hessel et al. (2019), and can be taken as a rough upper bound of the best performance achievable.

improvement also holds in ViTT-Cooking (+4.22 in ROUGE-L) and ViTT-All (+2.97 in ROUGE-L). We do not observe consistent and significant trends among the different multimodal pretraining strategies: MASS pretraining with video (**MASSvid**), with video and droptext (**MASSdrop**), or with alignment tasks (**MASSalign**).⁸ Furthermore, we observe that most of the pretraining improvement is achievable via text-only MASS pretraining. Across all three datasets, while *Multimodal Pretraining* (E2vidD6-MASSvid-BiD) is consistently better than *Text Pretraining* (E2vidD6-MASS-BiD), the differences are quite small (under one ROUGE-L point).

It’s worthy noting that for MASSalign, the best validation accuracies for the pretraining tasks are reasonably high: for YT8M-cook, we observed 90% accuracy for the alignment task, and 80% for the ordering task (for HowTo100M: 87% and 71.4%, respectively), where random guess would yield 50%. This suggests that our video features are reasonably strong, and the findings above are not due to weak visual representations.

⁸Limited improvement with MASSalign suggests that such alignment tasks are better suited for retrieval (Luo et al., 2020).

Effect of other modeling choices We experiment with 2-layer decoder (D2) vs 6-layer decoder (D6), combined with either unidirectional fine-tuning (**UniD**) or bidirectional fine-tuning (**BiD**). Table 5 shows ablation results of the four possible combinations when finetuning a multimodal model using text-only pretraining on YouCook2 (a more complete list of results can be found in Appendix A.5, showing similar trends). The D6xBiD combination tends to yield the best performance, with the differences among the four configurations being relatively small (under one ROUGE-L point). For visual features, we also explored using 3D features (Xie et al., 2018) instead of 2D features during finetuning (with no pretraining or text-only pretraining), and do not find much difference in model performance on YouCook2. As a result, we use the simpler 2D features in our multimodal pretraining. We leave more extensive experiments with visual features as future work.

Generalization implications An important motivation for constructing the ViTT dataset and evaluating our models on it has been related to generalization. Since the YouCook2 benchmark is restricted to a small number of cooking recipes, there is little to be understood about how well models

Method	Input	ViTT-All				ViTT-Cooking			
		BLEU-1	METEOR	ROUGE-L	CIDEr	BLEU-1	METEOR	ROUGE-L	CIDEr
Constant baseline (“intro”)	-	1.42	3.32	11.15	0.28	1.16	2.93	10.21	0.25
<i>∅ Pretraining</i>									
E2D6-BiD	ASR	19.60	9.12	27.88	0.68	20.77	10.08	28.63	0.72
E2vidD6-BiD	Video + ASR	19.49	9.23	28.53	0.69	20.45	9.88	28.88	0.69
<i>Text Pretraining</i>									
E2D6-MASS-BiD	ASR	21.93	10.60	30.45	0.79	24.79	12.25	32.40	0.88
E2vidD6-MASS-BiD	Video + ASR	22.44	10.83	31.27	0.81	24.22	12.22	32.60	0.89
<i>Multimodal Pretraining</i>									
E2vidD6-MASSalign-BiD	Video + ASR	22.31	10.66	31.13	0.79	24.92	12.25	33.09	0.90
E2vidD6-MASSvid-BiD	Video + ASR	22.45	10.76	31.49	0.80	24.87	12.43	32.97	0.90
E2vidD6-MASSdrop-BiD	Video + ASR	22.37	11.00	31.40	0.82	24.48	12.22	33.10	0.89
Human	Video + ASR	43.34	33.56	41.88	1.26	41.61	32.50	41.59	1.21

Table 4: Segment-level captioning results on ViTT. For ViTT-All we pretrain on HowTo100M and WikiHow; for ViTT-Cooking we pretrain on YT8M-cook and Recipe1M. We report baseline scores for predicting the most common caption “intro”. We also estimate the human performance as a rough upper bound (details in Supplementary Material A.1; Table 9).

Method	BLEU-4	METEOR	ROUGE-L	CIDEr
D2-UniD	10.84	17.39	38.24	1.16
D6-UniD	11.39	18.00	38.71	1.22
D2-BiD	11.38	18.04	38.67	1.19
D6-BiD	11.47	17.70	38.80	1.25
D6-BiDalt	11.07	17.68	38.43	1.22
D6-BiD (S3D)	11.64	18.04	38.75	1.24

Table 5: Ablation study on YouCook2. We finetune a multimodal captioning model (E2vid) with either 2-layer decoder (D2) or 6-layer decoder (D6) using YT8M-cook /Recipe1M for MASS pretraining, combined with either unidirectional (UniD) or bidirectional (BiD) finetuning. We find no significant difference between using 2D and 3D features (marked as S3D).

trained and evaluated on it generalize. In contrast, the ViTT benchmark has a much wider coverage (for both cooking-related videos and general instructional videos), and no imposed topic overlap between train/dev/test. As such, there are two findings here that are relevant with respect to generalization: (a) the absolute performance of the models on the ViTT benchmark is quite high (ROUGE-L scores above 0.30 are usually indicative of decent performance), and (b) the performance on ViTT vs. YouCook2 is clearly lower (31.5 ROUGE-L vs. 39.0 ROUGE-L, reflecting the increased difficulty of the new benchmark), but it is maximized under similar pretraining and finetuning conditions, which allows us to claim that the resulting models generalize well and are quite robust over a wide variety of instructional videos.

6 Conclusions

Motivated to improve information-seeking capabilities for videos, we have collected and annotated a new dense video captioning dataset, ViTT, which is larger with higher diversity compared to YouCook2. We investigated several multimodal pretraining strategies for segment-level video captioning, and conducted extensive ablation studies. We concluded that MASS-style pretraining is the most decisive factor in improving the performance on all the benchmarks used. Even more to the point, our results indicate that most of the performance can be attributed to leveraging the ASR signal. We achieve new state-of-the-art results on the YouCook2 benchmark, and establish strong performance baselines for the new ViTT benchmark, which can be used as starting points for driving more progress in this direction.

Acknowledgements

We send warm thanks to Ashish Thapliyal for helping the first author debug his code and navigate the computing infrastructure, and to Sebastian Goodman for his technical help (and lightning fast responses!). We also thank the anonymous reviewers for their comments and suggestions.

References

Sami Abu-El-Haija, Nisarg Kothari, Joonseok Lee, Paul Natsev, George Toderici, Balakrishnan Varadarajan, and Sudheendra Vijayanarasimhan.

2016. Youtube-8m: A large-scale video classification benchmark. *arXiv preprint arXiv:1609.08675*.
- Jean-Baptiste Alayrac, Piotr Bojanowski, Nishant Agrawal, Ivan Laptev Josef Sivic, and Simon Lacoste-Julien. 2016. Unsupervised learning from narrated instruction videos. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- João Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4724–4733.
- Yen-Chun Chen, Linjie Li, Licheng Yu, Ahmed El Kholy, Faisal Ahmed, Zhe Gan, Yu Cheng, and Jingjing Liu. 2019. Uniter: Learning universal image-text representations. *arXiv preprint arXiv:1909.11740*.
- Dima Damen, Hazel Doughty, Giovanni Maria Farinella, Sanja Fidler, Antonino Furnari, Evangelos Kazakos, Davide Moltisanti, Jonathan Munro, Toby Perrett, Will Price, and Michael Wray. 2018. Scaling egocentric vision: The EPIC-KITCHENS dataset. In *European Conference on Computer Vision (ECCV)*.
- Michael Denkowski and Alon Lavie. 2014. Meteor universal: Language specific translation evaluation for any target language. In *Proceedings of the ninth workshop on statistical machine translation*, pages 376–380.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2018. Bert: Pre-training of deep bidirectional transformers for language understanding. *arXiv preprint arXiv:1810.04805*.
- Li Dong, Nan Yang, Wenhui Wang, Furu Wei, Xiaodong Liu, Yu Wang, Jianfeng Gao, Ming Zhou, and Hsiao-Wuen Hon. 2019. Unified language model pre-training for natural language understanding and generation. *arXiv preprint arXiv:1905.03197*.
- Jack Hessel, Bo Pang, Zhenhai Zhu, and Radu Soricut. 2019. A case study on combining asr and visual features for generating instructional video captions. In *Proceedings of CoNLL*.
- Will Kay, João Carreira, Karen Simonyan, Brian Zhang, Chloe Hillier, Sudheendra Vijayanarasimhan, Fabio Viola, Tim Green, Trevor Back, Apostol Natsev, Mustafa Suleyman, and Andrew Zisserman. 2017. The kinetics human action video dataset. *ArXiv*, abs/1705.06950.
- Juho Kim, Phu Tran Nguyen, Sarah Weir, Philip J Guo, Robert C Miller, and Krzysztof Z Gajos. 2014. Crowdsourcing step-by-step information extraction to enhance existing how-to videos. In *CHI*.
- Mahnaz Koupaee and William Yang Wang. 2018. Wikihow: A large scale text summarization dataset. *arXiv preprint arXiv:1810.09305*.
- Ranjay Krishna, Kenji Hata, Frederic Ren, Li Fei-Fei, and Juan Carlos Niebles. 2017. Dense-captioning events in videos. In *International Conference on Computer Vision (ICCV)*.
- Zhenzhong Lan, Mingda Chen, Sebastian Goodman, Kevin Gimpel, Piyush Sharma, and Radu Soricut. 2020. Albert: A lite bert for self-supervised learning of language representations. *ArXiv*, abs/1909.11942.
- Jie Lei, Liwei Wang, Yelong Shen, Dong Yu, Tamara L. Berg, and Mohit Bansal. 2020. Mart: Memory-augmented recurrent transformer for coherent video paragraph captioning. In *The 58th Annual Meeting of the Association for Computational Linguistics*.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *ArXiv*, abs/1910.13461.
- Gen Li, Nan Duan, Yuejian Fang, Daxin Jiang, and Ming Zhou. 2019. Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training. *arXiv preprint arXiv:1908.06066*.
- Chin-Yew Lin and Franz Josef Och. 2004. Automatic evaluation of machine translation quality using longest common subsequence and skip-bigram statistics. In *Proceedings of the 42nd Annual Meeting on Association for Computational Linguistics*, page 605. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized BERT pretraining approach. *arXiv preprint arxiv:1907.11692*.
- Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. 2019. Vilt: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks. In *Advances in Neural Information Processing Systems*, pages 13–23.
- Huaishao Luo, Lei Ji, Botian Shi, Haoyang Huang, Nan Duan, Tianrui Li, Xilin Chen, and Ming Zhou. 2020. Univilm: A unified video and language pre-training model for multimodal understanding and generation. *arXiv preprint arXiv:2002.06353*.
- Lauren E Margulieux, Mark Guzdial, and Richard Catrambone. 2012. Subgoal-labeled instructional material improves performance and transfer in learning to develop mobile applications. In *Conference on International Computing Education Research*.
- Javier Marin, Aritro Biswas, Ferda Ofli, Nicholas Hynes, Amaia Salvador, Yusuf Aytar, Ingmar Weber, and Antonio Torralba. 2019. Recipe1m+: A dataset for learning cross-modal embeddings for cooking recipes and food images. *IEEE transactions on pattern analysis and machine intelligence*.

- Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, and Andrew Zisserman Josef Sivic. 2020. End-to-end learning of visual representations from uncurated instructional videos. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- Antoine Miech, Dimitri Zhukov, Jean-Baptiste Alayrac, Makarand Tapaswi, Ivan Laptev, and Josef Sivic. 2019. Howto100m: Learning a text-video embedding by watching hundred million narrated video clips. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2630–2640.
- Celie O’Neil-Hart. 2018. Why you should lean into how-to content in 2018. www.thinkwithgoogle.com/advertising-channels/video/self-directed-learning-youtube/. Accessed: 2019-09-03.
- Boxiao Pan, Haoye Cai, De-An Huang, Kuan-Hui Lee, Adrien Gaidon, Ehsan Adeli, and Juan Carlos Niebles. 2020. Spatio-temporal graph for video captioning with knowledge distillation. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the 40th annual meeting on association for computational linguistics*, pages 311–318. Association for Computational Linguistics.
- Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. 2018. Improving language understanding by generative pre-training. https://s3-us-west-2.amazonaws.com/openai-assets/research-covers/language-unsupervised/language_understanding_paper.pdf.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2019. Exploring the limits of transfer learning with a unified text-to-text transformer. *ArXiv*, abs/1910.10683.
- Ramon Sanabria, Ozan Caglayan, Shruti Palaskar, Desmond Elliott, Loïc Barrault, Lucia Specia, and Florian Metze. 2018. How2: A large-scale dataset for multimodal language understanding. *arXiv preprint arXiv:1811.00347*.
- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2015. Neural machine translation of rare words with subword units. *arXiv preprint arXiv:1508.07909*.
- Botian Shi, Lei Ji, Yaobo Liang, Nan Duan, Peng Chen, Zhendong Niu, and Ming Zhou. 2019. Dense procedure captioning in narrated instructional videos. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6382–6391.
- Kaitao Song, Xu Tan, Tao Qin, Jianfeng Lu, and Tie-Yan Liu. 2019. Mass: Masked sequence to sequence pre-training for language generation. *arXiv preprint arXiv:1905.02450*.
- Weijie Su, Xizhou Zhu, Yue Cao, Bin Li, Lewei Lu, Furu Wei, and Jifeng Dai. 2019. Vi-bert: Pre-training of generic visual-linguistic representations. *arXiv preprint arXiv:1908.08530*.
- Chen Sun, Fabien Baradel, Kevin Murphy, and Cordelia Schmid. 2019a. Contrastive bidirectional transformer for temporal representation learning. *arXiv preprint arXiv:1906.05743*.
- Chen Sun, Austin Myers, Carl Vondrick, Kevin Murphy, and Cordelia Schmid. 2019b. Videobert: A joint model for video and language representation learning. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 7464–7473.
- Hao Tan and Mohit Bansal. 2019. Lxmert: Learning cross-modality encoder representations from transformers. *arXiv preprint arXiv:1908.07490*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017a. Attention is all you need. *ArXiv*, abs/1706.03762.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017b. Attention is all you need. In *Advances in neural information processing systems*, pages 5998–6008.
- Ramakrishna Vedantam, C Lawrence Zitnick, and Devi Parikh. 2015. Cider: Consensus-based image description evaluation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4566–4575.
- Jiang Wang, Yang Song, Thomas Leung, Chuck Rosenberg, Jingbin Wang, James Philbin, Bo Chen, and Ying Wu. 2014. Learning fine-grained image similarity with deep ranking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1386–1393.
- Xin Wang, Jiawei Wu, Junkun Chen, Lei Li, Yuanfang Wang, and William Yang Wang. 2019. VateX: A large-scale, high-quality multilingual dataset for video-and-language research. *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 4580–4590.
- Sarah Weir, Juho Kim, Krzysztof Z Gajos, and Robert C Miller. 2015. Learnersourcing subgoal labels for how-to videos. In *CSCW*.
- Saining Xie, Chen Sun, Jonathan Huang, Zhuowen Tu, and Kevin Murphy. 2018. Rethinking spatiotemporal feature learning: Speed-accuracy trade-offs in video classification. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 305–321.

- YouTubeBlog. 2017. You know what's cool? a billion hours. <https://youtube.googleblog.com/2017/02/you-know-whats-cool-billion-hours.html>. Accessed: 2020-06-23.
- ZenithMedia. 2019. Online video viewing to reach 100 minutes a day in 2021. <https://www.zenithmedia.com/online-video-viewing-to-reach-100-minutes-a-day-in-2021/>. Accessed: 2020-06-23.
- Ziqi Zhang, Yaya Shi, Chunfeng Yuan, Bing Li, Peijin Wang, Weiming Hu, and Zhengjun Zha. 2020. Object relational graph with teacher-recommended learning for video captioning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- Qi Zheng, Chaoyue Wang, and Dacheng Tao. 2020. Syntax-aware action targeting for video captioning. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.
- Luowei Zhou, Chenliang Xu, and Jason J Corso. 2018a. Towards automatic learning of procedures from web instructional videos. In *Thirty-Second AAAI Conference on Artificial Intelligence*.
- Luowei Zhou, Chenliang Xu, and Jason J. Corso. 2018b. YouCookII dataset. http://youcook2.eecs.umich.edu/static/YouCookII/youcookii_readme.pdf. Accessed: 2020-06-23.
- Luowei Zhou, Yingbo Zhou, Jason J. Corso, Richard Socher, and Caiming Xiong. 2018c. End-to-end dense video captioning with masked transformer. *2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 8739–8748.
- Linchao Zhu and Yi Yang. 2020. Actbert: Learning global-local video-text representations. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE.

A Appendix

Supplementary Material for “Multimodal Pretraining for Dense Video Captioning”.

A.1 The ViTT dataset

Sampling video for annotation. The goal of the ViTT dataset design is to mirror topic distribution in the “wild”. Therefore, instead of starting from specific how-to instructions and searching for corresponding videos, we sampled videos from the validation set of the YouTube-8M dataset (Abu-El-Haija et al., 2016), a large-scale collection of YouTube videos with topical labels, subject to YouTube policies.

Exclusion criteria were lack of English ASR and the topic label “Game”. The latter was motivated by the fact that in this type of videos, the visual information predominantly features video games, while the ViTT dataset was intended to contain only videos with real-world human actions. Cooking videos can be easily identified by sampling videos that came with “Cooking” or “Recipe” topic labels. Given the convenience and the fact that much of prior work in this area had focused on cooking videos, approximately half of the dataset was designed to include cooking videos only, while the remaining videos would be randomly sampled non-cooking videos, as long as they were verified as instructional by human annotators.

Annotation process Annotators were presented with a video alongside its timestamped, automatic transcription shown in sentence-length paragraphs. They were asked to watch the video and first judge whether the video was instructional. For the purpose of our dataset, we determine that a video is instructional if it focuses on real-world human actions that are accompanied by procedural language explaining what is happening on screen, in reasonable details. Also for our purposes, instructional videos need to be grounded in real life, with a real person in the video exemplifying the action being verbally described.

For videos judged to be instructional, annotators were then asked to:

- Delimit the main segments of the video.
- Determine their start time if different from the automatically suggested start time (explained below).

- Provide a label summarizing or explaining the segment.

Annotation guidelines Annotators were instructed to identify video segments with two potential purposes:

- Allow viewers to jump straight to the start of a segment for rewatch.
- Present viewers with an index to decide whether to watch the video in full or directly skip to the segment of interest.

Our guidelines suggested a range of five to ten segments as long as the the structure and content of the video permitted. For short videos, the direction was to prioritize quality over quantity and to only define those segments that formed the narrative structure of the video, even if the resulting number of segments was below 5.

To help annotators determine segment start times, transcriptions were shown in “sentences” — we expected that sentence start times might be good candidates for segment start times. We obtained sentence boundaries automatically as follows. Given the stream of timestamped ASR tokens for a video, we first separated them into blocks by breaking two consecutive tokens whenever they were more than 2 seconds apart. We then used a punctuation prediction model to identify sentence boundaries in each resulting block. Each sentence was shown with the timestamp corresponding to its first token. Annotators were advised that transcriptions had been automatically divided into paragraphs that may or may not correspond to a video segment — if they decided that a segment started from a particular sentence, they could choose to use the start time of the sentence as the start time for the segment, or, if needed, they could put in an adjusted start time instead.

Once the start time had been identified, annotators were asked to provide a free-text label to summarize each segment. We instructed the annotators to use nouns or present participles (-ing form of verbs) to write the labels for the video segments, whenever possible. Additionally, we asked that the labels be succinct while descriptive, using as few words as possible to convey as much information as possible.

Data statistics and post-processing The resulting dataset consists of 8,169 instructional videos that received segment-level annotations, of which

3,381 are cooking-related. Overall there are an average of 7.1 segments per video (max: 19). Given our instructions, the descriptions are much shorter in lengths compared to a typical captioning dataset: on average there are 2.97 words per description (max: 16); 20% of the captions are single-word, 22% are two-words, and 25% are three words. We refer to these descriptions as “tags” given how short they are.

When possible, annotators were also asked to start and end the video with an opening and closing segment. As a result, most annotations start with an introduction segment: this accounts for roughly 11% of the 88455 segments in the dataset (“intro”: 8%, “introduction”: 2.3%). Note that while “intro” and “introduction” are clearly paraphrases of each other, an automatic metric will penalize a model predicting “intro” when the groundtruth is “introduction”. Similarly, the ending segment was described in several varieties: “outro”: 3.4%, “closing”: 1%, “closure”, “conclusion”, “ending”, “end of video”: each under 1%. Penalizing paraphrases of the ground truth is an inherent weakness of automatic metrics. To mitigate this, we decided to reduce the chance of this happening for the most frequent tags in the dataset. That is, in our experiments, we identified three groups of tags among the top-20 most frequent tags, and standardized them as follows.

intro	intro, introduction, opening
outro	outro, closing, closure, conclusion, ending, end of video, video closing
result	finished result, final result, results

Table 6: Standardization of top tags

Note that this does not mean we can solve this problem as a classification task like in visual question answering (VQA): overall, there are 56,027 unique tags with a vocabulary size of 12,509 for the 88,455 segments; 51,474 tags appeared only once in the dataset, making it infeasible to reduce the segment-level captioning problem into a pure classification task. Table 7 shows the top 10 most frequent tags after standardization.

Estimate of human performance. A subset of the candidate videos were given to three annotators⁹, to help us understand variations in human annotations. 5,840 videos received dense captioning

⁹A small set were unintentionally given to six annotators.

Tag	% of segments
intro	11.4
outro	6.6
result	0.9
ingredients	0.8
listing ingredients	0.2
supplies	0.2
mixing ingredients	0.2
materials	0.1
what you’ll need	0.1
lining the eyes	0.1

Table 7: 10 most frequent tags after standardization.

from exactly one annotator and were used as training data. Videos with more than one annotation were used as validation / test data. Note that not all the videos with multiple timeline annotations have exactly three sets of them — in fact, 1368 videos received 3-way segment-level annotations. This is because not all annotators agreed on whether a video was instructional. Computing annotator agreement for the annotated timelines is non-trivial. Here we focus on an estimate of tagging agreement when a pair of annotators agreed over the segment start time. Specifically, we go through each video that received multiple segment-level annotations. For each segment where two annotators chose the same ASR sentence as its starting point, we take the tags they produced for this segment and consider one of them as groundtruth, the other as prediction, and add that into our pool of (groundtruth, prediction) pairs. We can then compute standard automatic evaluations metrics over this pool. The results are as follows.

BLEU-1	METEOR	ROUGE-L	CIDEr
43.34	33.56	41.88	1.26

Table 8: Estimate of human performance for the segment-level captioning on ViTT-All (computed over 7528 pairs).

BLEU-1	METEOR	ROUGE-L	CIDEr
41.61	32.50	41.59	1.21

Table 9: Estimate of human performance for the segment-level captioning on ViTT-Cooking (computed over 2511 pairs).

Note that METEOR, and CIDEr scores are both penalized by the lack of n-grams for higher n. That

is, when both groundtruth and prediction are single-word, say, “intro”, this pair will not receive a full score from any of these metrics. But the ROUGE-L score is in the same ballpark as estimate of human performance in prior work (Hessel et al., 2019). One might note that perhaps this pool of label pairs contains a higher share of “intro”, since annotators might be more likely to agree over where an opening segment starts. Indeed, 20% of the time, one of the tags is “intro”. Interestingly, in spite of standardization of top tags, 14% of the time one tag is “intro”, the other tag is *not* “intro”: they can be less frequent paraphrases (e.g., “welcoming”, “greeting”, “opening and welcoming”) or something semantically different (e.g., “using dremel tool”).

A.2 Separated vs. Concatenated-Modality Architecture

Prior work has explored both concatenating different modalities and feeding them into the same multimodal Transformer encoder (Sun et al., 2019b; Hessel et al., 2019), as well as separating them into unimodal transformers (Sun et al., 2019a; Lu et al., 2019). We opt for the separated architecture because it offers more flexibility. First, the concatenated architecture requires embedding the text and video features into the same space. When the video features are projected using a simple network, there is no guarantee that we can meaningfully project them into the text embedding space. VideoBERT (Sun et al., 2019b) gives more flexibility to the video embeddings by quantizing video features and learning an embedding for each code-word. However, the quantization step has subsequently been claimed to be detrimental (Sun et al., 2019a). Moreover, the concatenated architecture uses the same sets of forward and attention weights to process text and video, and performs layer normalization jointly between the two modalities, which is not necessarily meaningful. Finally, the separated architecture makes it easy to switch between variable length text-only, video-only, or text+video modalities, whereas concatenated architectures might rely on separating tokens, modalities embeddings, and using fixed sequence lengths (Luo et al., 2020).

A.3 Additional Implementation Details

We optimize all models on a nVidia v100 GPU using the Adam optimizer with inverse square root schedule, batch size 32, warm-up period of 4,000

iterations, and maximum learning rate of 0.0001, following MASS (Song et al., 2019). The positional embeddings are initialized randomly. We use dropout and attention dropout with probabilities 0.1. With E2vidD6, pretraining takes 3-6 days depending on the objective and bidirectional finetuning takes up to 1.5 days, however those times could be improved by optimizing the data pipeline.

A.4 Example Predictions

We show examples of **good** and **bad** predictions on YouCook2 (Figure 5 and ViTT-All (Figure 4 and 5). The captions are generated by E2vidD6-BiD (no pretraining) and E2vidD6-MASS-BiD (text-only MASS pretraining).

A.5 Full result tables

We present here tables with all the ablation results that we run. There are two main takeaway messages from the results involving the pretraining approach: (a) the accuracy improvements, as measured across all the metrics we use, indicate the value of using a pretraining approach to this problem, specifically one that is capable of leveraging the ASR signals at both pretraining and finetuning stages, and (b) the training speedup achieved from pretraining is impressive, as a pretrained model converges much faster than training from scratch. This is especially visible on ViTT-All where finetuning after MASS pretraining reaches best ROUGE-L score at epoch 2, whereas it takes around 11 epochs to converge when training from scratch.




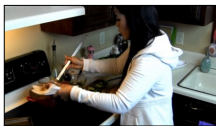

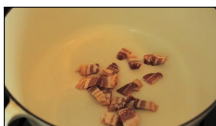


Sample Frame	Ground Truth	Ø-Pretraining	MASS-Pretraining	Comments
	<i>crush and grate the garlic</i>	<i>grate garlic and add to bowl (good)</i>	<i>crush ginger and garlic (good)</i>	ginger is correct despite not appearing in ground truth.
	<i>crimp shut with fork</i>	<i>place the filling on the wrapper (ok)</i>	<i>seal the edges of the wrapper (good)</i>	pretrained model is more specific
	<i>place wings on the baking sheet and cook flipping</i>	<i>bake the pizza in the oven (bad)</i>	<i>cook the wings on the grill (good)</i>	only pretrained model predicted correct food
	<i>add the pork back into the hot oil</i>	<i>add the rice to the pot (bad)</i>	<i>place the meat on the pan (good)</i>	Ø model hallucinates the rice and pot
	<i>add thyme bay leaves onion and clam juice and boil the mixture</i>	<i>add diced tomatoes tomato puree and mix well (bad)</i>	<i>add thyme thyme onion and clam juice to the pot and stir (ok)</i>	Ø hallucinates a lot of nonexistent ingredients
	<i>cook bacon in a pot with oil and pepper</i>	<i>add chopped tomatoes to pan and stir (bad)</i>	<i>add bacon and stir (ok)</i>	both models missed oil and pepper (not mentioned in ASR)
	<i>pour dressing on top of the salad and toss</i>	<i>add dressing to the bowl (good)</i>	<i>serve the soup over the salad (bad)</i>	pretrained model referred to dressing as "soup"
	<i>slice the ginger into pieces</i>	<i>slice a celery (bad)</i>	<i>slice the chicken (bad)</i>	both models had wrong ingredients (ASR segment does not mention what is being sliced)

Figure 3: Example good and bad predictions on YouCook2. The pretrained model is generally but not always better. Note that there are no “intro” or “outro”-like labels on YouCook2 because the dataset was specifically curated to only contain actual recipe steps.

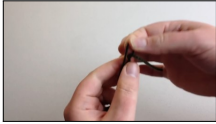
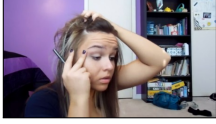
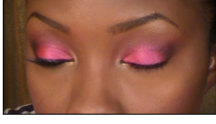


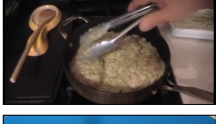
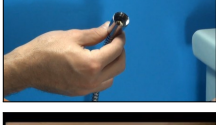

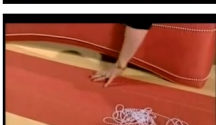
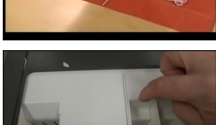

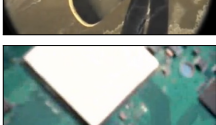
Sample Frame	Ground Truth	Ø-Pretraining	MASS-Pretraining	Comments
	<i>tightening extra loop</i>	<i>tightening the loop (good)</i>	<i>tightening the loop (good)</i>	both models perform well
	<i>adding eyeshadow</i>	<i>blending eye shadow (good)</i>	<i>applying eye shadow (good)</i>	both models perform well
	<i>showcasing the finished look</i>	<i>showing finished look (good)</i>	<i>showing finished look (good)</i>	both models perform well
	<i>rolling and folding the clay</i>	<i>rolling and blending (ok)</i>	<i>rolling and folding the clay (good)</i>	MASS is a bit more specific
	<i>highlighting brow bone</i>	<i>applying eye shadow (ok)</i>	<i>brushing on the brows (good)</i>	MASS is a bit more specific
	<i>covering the chicken and cooking</i>	<i>cooking the bread (bad)</i>	<i>cooking the chicken (good)</i>	only MASS got the right ingredient
	<i>connecting spray hose and sprayer</i>	<i>connecting the new cover (ok)</i>	<i>connecting the valve (good)</i>	spray hose is more specific than valve
	<i>implementing second layer</i>	<i>showing finished product (ok)</i>	<i>showing second layer (good)</i>	MASS is more specific
	<i>making decorative trim</i>	<i>cutting the edges (good)</i>	<i>cutting the fabric (good)</i>	both models yield good predictions
	<i>checking bleach container</i>	<i>outrou (bad)</i>	<i>checking the container (good)</i>	MASS is a bit more specific
	<i>demonstrating the flip</i>	<i>checking the battery (bad)</i>	<i>flipping the board (good)</i>	Ø model got influenced by car mechanics tutorials
	<i>tilting board</i>	<i>setting up the oven (bad)</i>	<i>turning the board (good)</i>	Ø overfitted on cooking videos

Figure 4: Example **good** predictions on ViTT-All (Part 1). The pretrained model is generally but not always better.

Method	Input	Pretraining	BLEU-4	METEOR	ROUGE-L	CIDEr
Constant Pred (Hessel et al., 2019)	-	-	2.70	10.30	21.70	0.15
MART (Lei et al., 2020)	Video	-	8.00	15.90	-	0.36
DPC (Shi et al., 2019)	Video + ASR	-	2.76	18.08	-	-
EMT (Zhou et al., 2018c)	Video	-	4.38	11.55	27.44	0.38
CBT (Sun et al., 2019a)	Video	Kinetics + HowTo100M	5.12	12.97	30.44	0.64
AT (Hessel et al., 2019)	ASR	-	8.55	16.93	35.54	1.06
AT+Video (Hessel et al., 2019)	Video + ASR	-	9.01	17.77	36.65	1.12
UniViLM #1 (Luo et al., 2020)	Video	-	6.06	12.47	31.48	0.64
UniViLM #2 (Luo et al., 2020)	Video + ASR	-	8.67	15.38	35.02	1.00
UniViLM #5 (Luo et al., 2020)	Video + ASR	HowTo100M	10.42	16.93	38.02	1.20
<i>∅ Pretraining</i>						
E2D2-UniD	ASR	-	7.42	15.15	33.26	0.85
E2D6-UniD	ASR	-	7.88	15.29	34.10	0.87
E2D2-BiD	ASR	-	6.85	15.64	34.26	0.91
E2D6-BiD	ASR	-	7.90	15.70	34.86	0.93
E2vidD2-UniD	Video + ASR	-	7.47	15.11	34.77	0.90
E2vidD6-UniD	Video + ASR	-	7.61	15.57	34.28	0.89
E2vidD2-BiD	Video + ASR	-	8.39	15.36	34.54	0.91
E2vidD6-BiD	Video + ASR	-	8.01	16.19	34.66	0.91
E2vidD2-BiDalt	Video + ASR	-	8.12	15.83	34.83	0.93
E2vid,D6-BiDalt	Video + ASR	-	7.70	16.11	34.78	0.91
E2vidD2-BiD (S3D)	Video + ASR	-	8.04	16.17	36.01	0.96
E2vidD6-BiD (S3D)	Video + ASR	-	7.91	16.28	35.23	0.93
<i>Text Pretraining</i>						
E2D2-MASS-UniD	ASR	YT8M-cook + Recipe1M	10.52	17.14	37.39	1.14
E2D6-MASS-UniD	ASR	YT8M-cook + Recipe1M	10.72	17.74	37.85	1.17
E2D2-MASS-BiD	ASR	YT8M-cook + Recipe1M	10.84	17.44	37.20	1.13
E2D6-MASS-BiD	ASR	YT8M-cook + Recipe1M	10.60	17.42	38.08	1.20
E2vidD2-MASS-UniD	Video + ASR	YT8M-cook + Recipe1M	10.84	17.39	38.24	1.16
E2vidD6-MASS-UniD	Video + ASR	YT8M-cook + Recipe1M	11.39	18.00	38.71	1.22
E2vidD2-MASS-BiD	Video + ASR	YT8M-cook + Recipe1M	11.38	18.04	38.67	1.19
E2vidD6-MASS-BiD	Video + ASR	YT8M-cook + Recipe1M	11.47	17.70	38.80	1.25
E2vid,D2-MASS-BiDalt	Video + ASR	YT8M-cook + Recipe1M	11.49	17.85	38.60	1.18
E2vid,D6-MASS-BiDalt	Video + ASR	YT8M-cook + Recipe1M	11.07	17.68	38.43	1.22
E2vidD2-MASS-BiD (S3D)	Video + ASR	YT8M-cook + Recipe1M	11.13	17.71	38.57	1.12
E2vidD6-MASS-BiD (S3D)	Video + ASR	YT8M-cook + Recipe1M	11.64	18.04	38.75	1.24
<i>Multimodal Pretraining</i>						
E2vidD2-MASSalign-BiD	Video + ASR	YT8M-cook + Recipe1M	11.54	17.57	37.70	1.15
E2vidD6-MASSalign-BiD	Video + ASR	YT8M-cook + Recipe1M	11.53	17.62	39.03	1.22
E2vidD2-MASSvid-BiD	Video + ASR	YT8M-cook + Recipe1M	11.17	17.71	38.32	1.17
E2vidD6-MASSvid-BiD	Video + ASR	YT8M-cook + Recipe1M	12.04	18.32	39.03	1.23
E2vidD2-MASSdrop-BiD	Video + ASR	YT8M-cook + Recipe1M	11.21	17.99	38.72	1.23
E2vidD6-MASSdrop-BiD	Video + ASR	YT8M-cook + Recipe1M	10.45	17.74	38.82	1.22
Human (Hessel et al., 2019)	Video + ASR	-	15.20	25.90	45.10	3.80

Table 10: Video Captioning Results on YouCook2. We use YT8M-cook/Recipe1M for pretraining. All video features are Compact 2D (Wang et al., 2014) except when marked as S3D (Xie et al., 2018).

Method	Input	Pretraining	BLEU-1	METEOR	ROUGE-L	CIDEr
Constant baseline (“intro”)	-	-	1.42	3.32	11.15	0.28
<i>∅ Pretraining</i>						
E2D2-UniD	ASR	-	17.94	8.55	27.06	0.64
E2D6-UniD	ASR	-	18.91	8.96	27.80	0.67
E2D2-BiD	ASR	-	18.81	8.82	27.63	0.65
E2D6-BiD	ASR	-	19.60	9.12	27.88	0.68
E2vidD2-UniD	Video + ASR	-	18.94	8.99	28.05	0.67
E2vidD6-UniD	Video + ASR	-	19.29	9.15	27.97	0.69
E2vidD2-BiD	Video + ASR	-	19.37	9.21	28.56	0.69
E2vidD6-BiD	Video + ASR	-	19.49	9.23	28.53	0.69
<i>Text Pretraining</i>						
E2D2-MASS-UniD	ASR	HowTo100M + WikiHow	21.53	10.24	29.95	0.77
E2D6-MASS-UniD	ASR	HowTo100M + WikiHow	22.09	10.58	30.67	0.79
E2D2-MASS-BiD	ASR	HowTo100M + WikiHow	20.73	10.20	30.15	0.76
E2D6-MASS-BiD	ASR	HowTo100M + WikiHow	21.93	10.60	30.45	0.79
E2vidD2-MASS-UniD	Video + ASR	HowTo100M + WikiHow	21.46	10.45	30.56	0.78
E2vidD6-UniD	Video + ASR	HowTo100M + WikiHow	22.21	10.75	30.86	0.81
E2vidD2-MASS-BiD	Video + ASR	HowTo100M + WikiHow	21.78	10.64	30.72	0.79
E2vidD6-MASS-BiD	Video + ASR	HowTo100M + WikiHow	22.44	10.83	31.27	0.81
<i>Multimodal Pretraining</i>						
E2vidD2-MASSalign-BiD	Video + ASR	HowTo100M + WikiHow	22.07	10.33	30.60	0.77
E2vidD6-MASSalign-BiD	Video + ASR	HowTo100M + WikiHow	22.31	10.66	31.13	0.79
E2vidD2-MASSvid-BiD	Video + ASR	HowTo100M + WikiHow	22.15	10.75	31.06	0.80
E2vidD6-MASSvid-BiD	Video + ASR	HowTo100M + WikiHow	22.45	10.76	31.49	0.80
E2vidD2-MASSdrop-BiD	Video + ASR	HowTo100M + WikiHow	21.84	10.55	31.10	0.79
E2vidD6-MASSdrop-BiD	Video + ASR	HowTo100M + WikiHow	22.37	11.00	31.40	0.82
Human estimate	Video + ASR	-	43.34	33.56	41.88	1.26

Table 11: Video captioning results on ViTT-All. We use HowTo100M/WikiHow for pretraining. We also estimate human performance (details in Appendix A.1; Table 9).

Method	Input	Pretraining	BLEU-1	METEOR	ROUGE-L	CIDEr
Constant baseline (“intro”)	-	-	1.16	2.93	10.21	0.25
<i>∅ Pretraining</i>						
E2D2-UniD	ASR	-	19.73	9.43	27.95	0.69
E2D6-UniD	ASR	-	20.24	9.93	28.59	0.71
E2D2-BiD	ASR	-	19.73	9.72	27.92	0.68
E2D6-BiD	ASR	-	20.77	10.08	28.63	0.72
E2vidD2-UniD	Video + ASR	-	19.97	9.75	28.30	0.69
E2vidD6-UniD	Video + ASR	-	20.46	9.93	28.62	0.69
E2vidD2-BiD	Video + ASR	-	20.60	10.08	29.45	0.71
E2vidD6-BiD	Video + ASR	-	20.45	9.88	28.88	0.69
<i>Text Pretraining</i>						
E2D2-MASS-UniD	ASR	YT8M-cook + Recipe1M	22.89	11.53	31.62	0.84
E2D6-MASS-UniD	ASR	YT8M-cook + Recipe1M	24.47	12.22	32.51	0.90
E2D2-MASS-BiD	ASR	YT8M-cook + Recipe1M	22.75	11.63	31.54	0.84
E2D6-MASS-BiD	ASR	YT8M-cook + Recipe1M	24.79	12.25	32.40	0.88
E2vidD2-MASS-UniD	Video + ASR	YT8M-cook + Recipe1M	23.86	11.85	32.32	0.86
E2vidD6-MASS-UniD	Video + ASR	YT8M-cook + Recipe1M	24.32	12.32	32.90	0.90
E2vidD2-MASS-BiD	Video + ASR	YT8M-cook + Recipe1M	22.93	11.68	32.15	0.87
E2vidD6-MASS-BiD	Video + ASR	YT8M-cook + Recipe1M	24.22	12.22	32.60	0.89
<i>Multimodal Pretraining</i>						
E2vidD2-MASSalign-BiD	Video + ASR	YT8M-cook + Recipe1M	24.02	11.91	32.73	0.86
E2vidD6-MASSalign-BiD	Video + ASR	YT8M-cook + Recipe1M	24.92	12.25	33.09	0.90
E2vidD2-MASSvid-BiD	Video + ASR	YT8M-cook + Recipe1M	24.15	12.10	32.96	0.88
E2vidD6-MASSvid-BiD	Video + ASR	YT8M-cook + Recipe1M	24.87	12.43	32.97	0.90
E2vidD2-MASSdrop-BiD	Video + ASR	YT8M-cook + Recipe1M	23.70	12.01	32.71	0.88
E2vidD6-MASSdrop-BiD	Video + ASR	YT8M-cook + Recipe1M	24.48	12.22	33.10	0.89
Human estimate	Video + ASR	-	41.61	32.50	41.59	1.21

Table 12: Video captioning results on ViTT-Cooking. We use YT8M-cook and Recipe1M for optional pretraining.