

# Learning a Grammar Inducer from Massive Uncurated Instructional Videos

Songyang Zhang<sup>1\*</sup>, Linfeng Song<sup>2</sup>, Lifeng Jin<sup>2</sup>, Haitao Mi<sup>2</sup>, Kun Xu<sup>2</sup>,  
Dong Yu<sup>2</sup> and Jiebo Luo<sup>1</sup>

<sup>1</sup>University of Rochester, Rochester, NY, USA

szhang83@ur.rochester.edu, jluo@cs.rochester.edu

<sup>2</sup>Tencent AI Lab, Bellevue, WA, USA

{lfsong, lifengjin, haitaomi, kxkunxu, dyu}@tencent.com

## Abstract

Video-aided grammar induction aims to leverage video information for finding more accurate syntactic grammars for accompanying text. While previous work focuses on building systems for inducing grammars on text that are well-aligned with video content, we investigate the scenario, in which text and video are only in loose correspondence. Such data can be found in abundance online, and the weak correspondence is similar to the indeterminacy problem studied in language acquisition. Furthermore, we build a new model that can better learn video-span correlation without manually designed features adopted by previous work. Experiments show that our model trained only on large-scale YouTube data with no text-video alignment reports strong and robust performances across three unseen datasets, despite domain shift and noisy label issues. Furthermore our model yields higher F1 scores than the previous state-of-the-art systems trained on in-domain data.

## 1 Introduction

Grammar induction is a fundamental and long-lasting (Lari and Young, 1990; Clark, 2001; Klein and Manning, 2002) problem in computational linguistics, which aims to find hierarchical syntactic structures from plain sentences. Unlike supervised methods (Charniak, 2000; Collins, 2003; Petrov and Klein, 2007; Zhang and Clark, 2011; Cross and Huang, 2016; Kitaev and Klein, 2018) that require human annotated treebanks, *e.g.*, Penn Treebank (Marcus et al., 1993), grammar inducers do not rely on any human annotations for training. Grammar induction is attractive since annotating syntactic trees by human language experts is expensive and time consuming, while the current treebanks are limited to several major languages and domains.

\* This work was done when Songyang Zhang was an intern at Tencent AI Lab.

Recently, deep learning models have achieved remarkable success across NLP tasks, and neural models have been designed (Shen et al., 2018b,a; Kim et al., 2019a,b; Jin et al., 2018) for grammar induction, which greatly advanced model performance on induction with raw text. Recent efforts have started to consider other useful information from multiple modalities, such as images (Shi et al., 2019; Jin and Schuler, 2020) and videos (Zhang et al., 2021). Specifically, Zhang et al. (2021) show that multi-modal information (*e.g.* motion, sound and objects) from videos can significantly improve the induction accuracy on verb and noun phrases. Such work uses curated multi-modal data publicly available on the web, which all assume that the meaning of a sentence needs to be identical (*e.g.*, being a caption) to the corresponding video or image. This assumption limits usable data to several small-scale benchmarks (Lin et al., 2014; Xu et al., 2016; Hendricks et al., 2017) with expensive human annotations on image/video captions.

The noisy correspondence between form and meaning is one of the main research questions in language acquisition (Akhtar and Montague, 1999; Gentner et al., 2001; Dominey and Dodane, 2004), where different proposals attempt to address this indeterminacy faced by children. There has been computational work incorporating such indeterminacy into their models (Yu and Siskind, 2013; Huang et al., 2021). For modeling empirical grammar learning with multi-modal inputs, two important questions still remain open: 1) *how can a grammar inducer benefit from large-scale multi-media data (e.g., YouTube videos) with noisy text-to-video correspondence?* and 2) *how can a grammar inducer show robust performances across multiple domains and datasets?* By using data with only weak cross-modal correspondence, such as YouTube videos and their automatically generated subtitles, we allow the computational models to face a similar indeterminacy problem, and exam-

ine how indeterminacy interacts with data size to influence learning behavior and performance of the induction models.

In this paper, we conduct the first investigation on both questions. Specifically, we collect 2.4 million video clips and the corresponding subtitles from instructional YouTube videos (HowTo100M [Miech et al. 2019](#)) to train multi-modal grammar inducers, instead of using the training data from a benchmark where text and video are in alignment. We then propose a novel model, named Pre-Trained Compound Probabilistic Context-Free Grammars (PTC-PCFG), that extends previous work ([Shi et al., 2019](#); [Zhang et al., 2021](#)) by incorporating a video-span matching loss term into the Compound PCFG ([Kim et al., 2019a](#)) model. To better capture the video-span correlation, it leverages CLIP ([Miech et al., 2020](#)), a state-of-the-art model pretrained on video subtitle retrieval, as the encoders for both video and text. Compared with previous work ([Zhang et al., 2021](#)) that independently extracts features from each modality before merging them using a simple Transformer ([Vaswani et al., 2017](#)) encoder, the encoders of our model have been pre-trained to merge such multi-modal information, and no human efforts are needed to select useful modalities from the full set.

Experiments on three benchmarks show that our model, which is trained on noisy YouTube video clips and no data from these benchmarks, produces substantial gains over the previous state-of-the-art system ([Zhang et al., 2021](#)) trained on in-domain video clips with human annotated captions. Furthermore, our model demonstrates robust performances across all three datasets. We suggest the limitations of our model and future directions for improvements through analysis and discussions. Code will be released upon paper acceptance.

In summary, the main contributions are:

- We are the first to study training a grammar inducer with massive general-domain noisy video clips instead of benchmark data, introducing the indeterminacy problem to the induction model.
- We propose PTC-PCFG, a novel model for unsupervised grammar induction. It is simpler in design than previous models and can better capture the video-text matching information.
- Trained only on noisy YouTube videos without finetuning on benchmark data, PTC-PCFG reports stronger performances than previous mod-

els trained on benchmark data across three benchmarks.

## 2 Background and Motivation

### 2.1 Compound PCFGs

A PCFG model in Chomsky Normal Form can be defined as a tuple of 6 terms  $(\mathcal{S}, \mathcal{N}, \mathcal{P}, \Sigma, \mathcal{R}, \Pi)$ , where they correspond to the start symbol, the sets of non-terminals, pre-terminals, terminals, production rules and their probabilities. Given pre-defined numbers of non-terminals and pre-terminals, a PCFG induction model tries to estimate the probabilities for all production rules.

The compound PCFG (C-PCFG) model ([Kim et al., 2019a](#)) adopts a mixture of PCFGs. Instead of a corpus-level prior used in previous work ([Kurihara and Sato, 2006](#); [Johnson et al., 2007](#); [Wang and Blunsom, 2013](#); [Jin et al., 2018](#)), C-PCFG imposes a sentence-specific prior on the distribution of possible PCFGs. Specifically in the generative story, the probability  $\pi_r$  for production rule  $r$  is estimated by model  $g$  that assigns a latent variable  $\mathbf{z}$  for each sentence  $\sigma$ , and  $\mathbf{z}$  is drawn from a prior distribution:

$$\pi_r = g(r, \mathbf{z}; \theta), \quad \mathbf{z} \sim p(\mathbf{z}). \quad (1)$$

where  $\theta$  represents the model parameters. The probabilities for all three types of CFG rules are defined as follows:

$$\begin{aligned} \pi_{S \rightarrow A} &= \frac{\exp(\mathbf{u}_A^\top f_s([\mathbf{w}_S; \mathbf{z}]))}{\sum_{A' \in \mathcal{N}} \exp(\mathbf{u}_{A'}^\top f_s([\mathbf{w}_S; \mathbf{z}]))}, \\ \pi_{A \rightarrow BC} &= \frac{\exp(\mathbf{u}_{BC}^\top [\mathbf{w}_A; \mathbf{z}])}{\sum_{B', C' \in \mathcal{N} \cup \mathcal{P}} \exp(\mathbf{u}_{B'C'}^\top [\mathbf{w}_A; \mathbf{z}])}, \\ \pi_{T \rightarrow w} &= \frac{\exp(\mathbf{u}_w^\top f_t([\mathbf{w}_T; \mathbf{z}]))}{\sum_{w' \in \Sigma} \exp(\mathbf{u}_{w'}^\top f_t([\mathbf{w}_T; \mathbf{z}]))}, \end{aligned} \quad (2)$$

where  $A \in \mathcal{N}$ ,  $B$  and  $C \in \mathcal{N} \cup \mathcal{P}$ ,  $T \in \mathcal{P}$ ,  $w \in \Sigma$ . Both  $\mathbf{w}$  and  $\mathbf{u}$  are dense vectors representing words and all types of non-terminals, and  $f_s$  and  $f_t$  are neural encoding functions.

Optimizing the C-PCFG model involves maximizing the marginal likelihood  $p(\sigma)$  of each training sentence  $\sigma$  for all possible  $\mathbf{z}$ :

$$\log p_\theta(\sigma) = \log \int_{\mathbf{z}} \sum_{t \in \mathcal{T}_G(\sigma)} p_\theta(t|\mathbf{z}) p(\mathbf{z}) d\mathbf{z} \quad (3)$$

where  $\mathcal{T}_G(\sigma)$  indicates all possible parsing trees for sentence  $\sigma$ . Since computing the integral over  $\mathbf{z}$

is intractable, this objective is optimized by maximizing its evidence lower bound  $\text{ELBO}(\sigma; \phi, \theta)$ :

$$\text{ELBO}(\sigma; \phi, \theta) = \mathbb{E}_{q_\phi(\mathbf{z}|\sigma)}[\log p_\theta(\sigma|\mathbf{z})] - \text{KL}[q_\phi(\mathbf{z}|\sigma)||p(\mathbf{z})], \quad (4)$$

where  $q_\phi(\mathbf{z}|\sigma)$  is the variational posterior calculated by another neural network with parameters  $\phi$ . Given a sampled  $\mathbf{z}$ , the log-likelihood term  $\log p_\theta(\sigma|\mathbf{z})$  is calculated via the inside algorithm. The KL term can be computed analytically when both the prior  $p(\mathbf{z})$  and the variational posterior  $q_\phi(\mathbf{z}|\sigma)$  are Gaussian (Kingma and Welling, 2014).

## 2.2 Multi-Modal Compound PCFGs

Multi-Modal Compound PCFGs (MMC-PCFG) (Zhang et al., 2021) extends C-PCFG with a model to match a video  $v$  with a span  $c$  in a parse tree  $t$  of a sentence  $\sigma$ . It extracts  $M$  visual and audio features from a video  $v$  and encodes them via a multi-modal transformer (Gabeur et al., 2020), denoted as  $\Psi = \{\psi^i\}_{i=1}^M$ . The word representation  $\mathbf{h}_i$  of the  $i$ th word is computed by BiLSTM. Given a particular span  $c = w_i, \dots, w_j$ , its representation  $\mathbf{c}$  is the weighted sum of all label-specific span representations:

$$\mathbf{c} = \sum_{k=1}^{|\mathcal{N}|} p(k|c, \sigma) f_k \left( \frac{1}{j-i+1} \sum_{l=i}^j \mathbf{h}_l \right), \quad (5)$$

where  $\{p(k|c, \sigma) | 1 \leq k \leq |\mathcal{N}|\}$  are the phrasal label probabilities of span  $c$ . The representation of a span  $\mathbf{c}$  is then correspondingly projected to  $M$  separate embeddings via gated embedding (Miech et al., 2018), denoted as  $\Xi = \{\xi^i\}_{i=1}^M$ . Finally the video-text matching loss is defined as a sum over all video-span matching losses weighted by the marginal probability of a span from the parser:

$$s_{mm}(v, \sigma) = \sum_{c \in \sigma} p(c|\sigma) h_{mm}(\Xi, \Psi), \quad (6)$$

where  $h_{mm}(\Xi, \Psi)$  is a hinge loss measuring the distances from video  $v$  to the matched and unmatched (*i.e.* span from another sentence) span  $c$  and  $c'$  and the distances from span  $c$  to the matched

and unmatched (*i.e.* another video) video  $v$  and  $v'$ :

$$\omega_i(\mathbf{c}) = \frac{\exp(\mathbf{u}_i^\top \mathbf{c})}{\sum_{j=1}^M \exp(\mathbf{u}_j^\top \mathbf{c})}, \quad (7)$$

$$o(\Xi, \Psi) = \sum_{i=1}^M \omega_i(\mathbf{c}) \cos(\xi^i, \psi^i), \quad (8)$$

$$h_{mm}(\Xi, \Psi) = \mathbb{E}_{c'}[o(\Xi', \Psi) - o(\Xi, \Psi) + \epsilon]_+ + \mathbb{E}_{v'}[o(\Xi, \Psi') - o(\Xi, \Psi) + \epsilon]_+, \quad (9)$$

where  $\Xi'$  is a set of unmatched span expert embeddings of  $\Psi$ ,  $\Psi'$  is a set of unmatched video representations of  $\Xi$ ,  $\epsilon$  is a positive margin,  $[\cdot]_+ = \max(0, \cdot)$ ,  $\{\mathbf{u}_i\}_{i=1}^M$  are learned weights, and the expectations are approximated with one sample drawn from the training data. During training, both ELBO and the video-text matching loss are jointly optimized.

## 2.3 Limitation and Motivation

Existing work on multi-modal grammar induction aims at leveraging strict correspondence between image/video and text for information about syntactic categories and structures of the words and spans in the text. However, such datasets are expensive to annotate. Besides, the ambiguous correspondence between language and real-world context, observed in language acquisition, is not really reflected in such training setups.

As a result, we believe that the previous work fails to answer the following important questions: 1) how well a grammar inducer would perform when it is trained only on noisy multi-media data; 2) how the scale of training data would affect the performance and cross-domain robustness?

## 3 Training a Grammar Inducer with Massive YouTube Videos

We make the first investigation into the above questions by leveraging massive video clips from instructional YouTube videos to train our grammar inducer. Different from the benchmark data used by previous work, the YouTube video clips do not contain paired sentences. This section will first introduce the method for generating noisy training instances (video clip and sentence pairs) from YouTube videos (§3.1), before describing a novel grammar induction model (§3.2) with pre-trained text and video encoders.

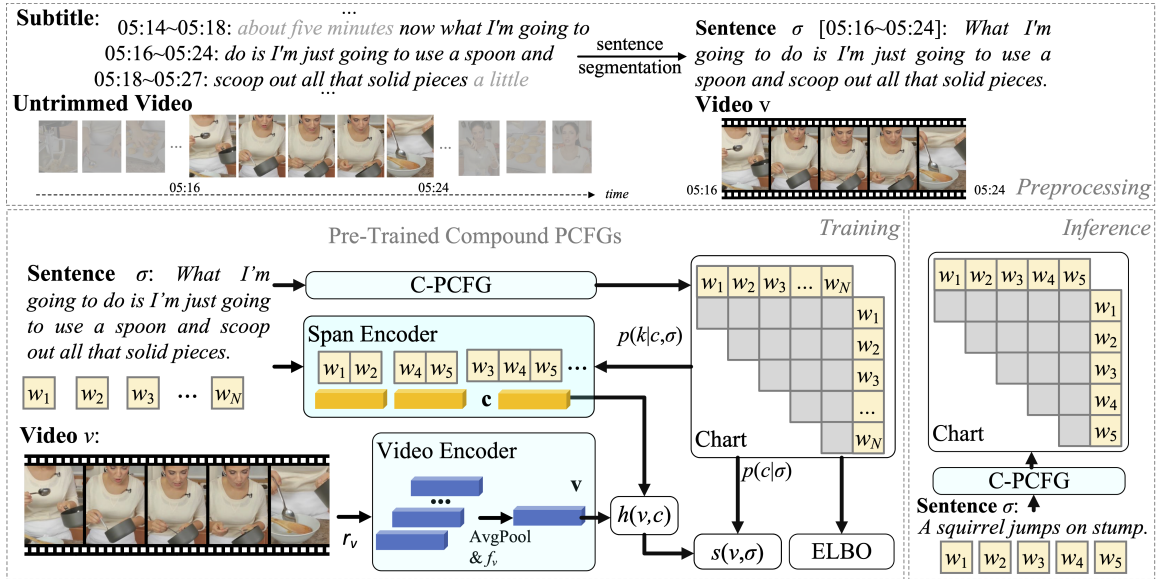


Figure 1: The pipeline of our approach.

### 3.1 Harvesting Training Instances from YouTube Videos

Given a YouTube video, we would like to generate a set of video clip and subtitle pairs  $\Omega = \{(v, \sigma)\}$ , where each subtitle  $\sigma$  is a complete sentence and is aligned in time with its paired video clip  $v$ . To this end, the YouTube API is chosen to obtain all subtitles of the video. But, our observation finds that most obtained subtitles are not complete sentences, and in some cases, a complete sentence can last for several continuous video fragments. Meanwhile, they do not contain any punctuation, which is a key factor for sentence segmentation. As shown in the top part of Figure 1, we design an algorithm that takes the following steps to find each complete sentence and its corresponding video clip.

**Sentence segmentation.** In the first step, we try to find complete sentences from the subtitles. We first concatenate all subtitles from the same video are concatenated into a very long sequence of tokens. Next, a punctuation restoration model<sup>1</sup> (Tilk and Alumäe, 2016) is adopted to insert punctuation into the sequence. Lastly, sentences are segmented based on certain punctuation (e.g., “.”, “?”, “!”).

**Video clip extraction.** In the second step, we trim the corresponding video clips. Each raw subtitle contains its start and an end times. We assume each word within the raw subtitle occupies equal time and record the start and end times for

<sup>1</sup>We manually punctuate subtitles from 10 videos randomly selected from HowTo100M, which contains 461 sentences after annotation. The punctuation restoration model has an overall F1 score of 74.1% with the manual labels.

each word. After that, given a complete sentence  $\sigma = w_1, w_2, \dots, w_N$ , we use the start time of its first word  $w_1$  and the end time of its last word  $w_N$  as the start and end times of  $\sigma$ . Lastly, we segment a complete sentence  $\sigma$ 's corresponding video clip  $v$  based on its start and end times.

### 3.2 Model: Pre-Trained Compound PCFGs

After harvesting large-scale sentence and video pairs, the next step is to build a strong grammar induction model that can benefit from them. In this section, we introduce our Pre-Trained Compound PCFGs (PTC-PCFG) model for unsupervised grammar induction. As shown in the lower part of Figure 1, the PTC-PCFG model composes of a video encoder, a span encoder and a parsing model. Both the video encoder and the span encoder are initialized from the MIL-NCE model (Miech et al., 2020), a pre-trained video-text matching model that takes a simple design and has shown superior zero-shot results on many video understanding tasks, such as video retrieval, video question answering, *etc.* We first introduce the pre-trained video and span encoders, before covering the training and inference details of PTC-PCFG.

**Video encoding.** The first step is to encode a video  $v$  to its representation  $\mathbf{v}$ . To do this, we first segment  $v$  into small video clips, where each video clip  $v_i$  consists of  $T$  frames. Following Zhang et al. (2021), we sample  $L$  video clips with equal interval for efficiency. We use the video encoder from the MIL-NCE model (Miech et al., 2020) as our video encoder and only fine-tune its last fully connected

layer  $f_v$  for efficiency. In more detail, for each sampled video clip, we pre-compute the input of  $f^v$  as its representation, denoted as  $\{\mathbf{h}_i^v\}_{i=1}^L$ . Then we feed them into  $f^v$  and average the output as its representation  $\mathbf{v}$ , denoted as,

$$\mathbf{v} = \text{AvgPool}(\{f^v(\mathbf{h}_i^v)\}_{i=1}^L), \quad (10)$$

where AvgPool indicates average pooling.

**Span encoding.** The next step is to compute a span representation  $\mathbf{c}$  for each particular span  $c = w_i, \dots, w_j$  ( $1 \leq i < j \leq N$ ) in sentence  $\sigma = w_1, w_2, \dots, w_N$ . The pre-trained text encoder of MIL-NCE consists of a word embedding layer and two stacked fully connected layers,  $f_0^c$  and  $f_1^c$ . Motivated by Zhao and Titov (2020); Zhang et al. (2021), we expect to learn  $|\mathcal{N}|$  different span representations, each is specified for one non-terminal node. However, directly applying the pre-trained text encoder is not feasible, since it has only one output layer  $f_1^c$ . Therefore, we duplicate  $f_1^c$  for  $|\mathcal{N}|$  times, denoted as  $\{f_k^c\}_{k=1}^{|\mathcal{N}|}$ , and compose  $|\mathcal{N}|$  label-specific output layers. In more detail, we first encode each word  $w_i$  with the word embedding layer, denoted as  $\mathbf{h}_i^c$ . Then we feed the word embeddings to  $f_0^c$ , ReLU, maximum pooling and each label-specific output layer sequentially. We also compute the probabilities of its phrasal labels  $\{p(k|c, \sigma) | 1 \leq k \leq |\mathcal{N}|\}$ , as illustrated in Section 2.1. Lastly, the span representation  $\mathbf{c}$  is the sum of all label-specific span representations weighted by the probabilities we predicted, denoted as:

$$\begin{aligned} \tau &= \text{MaxPool}(\text{ReLU}(f_0^c(\mathbf{h}_i^c))) \\ \mathbf{c} &= \sum_{k=1}^{|\mathcal{N}|} p(k|c, \sigma) f_k^c(\tau), \end{aligned} \quad (11)$$

where MaxPool is a maximum pooling operation and ReLU is a ReLU activation function.

**Training.** As shown in lower left of Figure 1, we optimize both the video-text matching loss and evidence lower bound during training. We first compute the similarity between a video clip  $v$  and a particular span  $c$  via dot product and then compute a triplet hinge loss as following,

$$\begin{aligned} h(v, c) &= \mathbb{E}_{c'}[\mathbf{c}' \cdot \mathbf{v} - \mathbf{c} \cdot \mathbf{v} + \epsilon]_+ \\ &+ \mathbb{E}_{v'}[\mathbf{c} \cdot \mathbf{v}' - \mathbf{c} \cdot \mathbf{v} + \epsilon]_+, \end{aligned} \quad (12)$$

where  $\epsilon$  is a positive margin,  $[\cdot]_+ = \max(0, \cdot)$ ,  $v'$  is a clip from a different video and  $c'$  is a span from

a different sentence. The video-text matching loss is correspondingly defined as,

$$s(v, \sigma) = \sum_{c \in \sigma} p(c|\sigma) h(v, c), \quad (13)$$

where  $p(c|\sigma)$  is the probability of a particular span  $c$  being a syntactic phrase. Finally, the overall loss function is composed by the ELBO and the video-text matching loss:

$$\mathcal{L}(\phi, \theta) = \sum_{(v, \sigma) \in \Omega} -\text{ELBO}(\sigma; \phi, \theta) + \alpha s(v, \sigma), \quad (14)$$

where  $\alpha$  is a constant balancing these two terms.

**Inference.** During inference, given a sentence  $\sigma$ , we predict the most likely tree  $t^*$  without accessing videos, as shown in the lower right of Figure 1. Since computing the integral over  $\mathbf{z}$  is intractable, we estimate  $t^*$  with the following approximation,

$$\begin{aligned} t^* &= \arg \max_t \int_{\mathbf{z}} p_\theta(t|\mathbf{z}) p_\theta(\mathbf{z}|\sigma) d\mathbf{z} \\ &\approx \arg \max_t p_\theta(t|\sigma, \boldsymbol{\mu}_\phi(\sigma)), \end{aligned} \quad (15)$$

where  $\boldsymbol{\mu}_\phi(\sigma)$  is the mean vector of the variational posterior  $q_\phi(\mathbf{z}|\sigma)$ , and  $t^*$  is obtained by the CYK algo. (Cocke, 1969; Younger, 1967; Kasami, 1966).

## 4 Experiments

### 4.1 Datasets

Following previous work, we evaluate all systems on three benchmarks (i.e., DiDeMo, YouCook2 and MSRVTT). Instead of training on benchmark data, our models are trained on the data harvested from HowTo100M dataset. Below shows more details about these datasets:

**DiDeMo** (Hendricks et al., 2017) contains 10k unedited personal Flickr videos. Each video is associated with roughly 3-5 video-sentence pairs. There are 32 994, 4 180 and 4021 video pairs in the training, validation and testing sets.

**YouCook2** (Zhou et al., 2018) contains 2000 long untrimmed YouTube videos from 89 cooking recipes. The procedure steps for each video are annotated with temporal boundaries and described by imperative English sentences. There are 8 913, 969 and 3 310 video-sentence pairs in the training, validation and testing sets.

**MSRVTT** (Xu et al., 2016) contains 10k generic YouTube videos accompanied by 200k captions annotated by paid human workers. There are 130 260,

9 940 and 59 794 video-sentence pairs in the training, validation and testing sets.

**HowTo100M** (Miech et al., 2019) is a large-scale dataset of 136 million video clips sourced from 1.22M narrated instructional web videos depicting humans performing more than 23k different visual tasks. Noted that there are 404 videos in HowTo100M exists in YouCook2, we exclude these videos during training.

## 4.2 Evaluation

We discard punctuation, lowercase all words, replace numbers with a special token and ignore trivial single-word and sentence-level spans during testing following Kim et al. (2019a). Besides, we follow previous work (Shi et al., 2019; Zhang et al., 2021) by using a state-of-the-art constituency parser (Benepar Kitaev et al. 2019) to obtain the reference trees for evaluation<sup>2</sup>. Following Shi et al. (2020); Zhang et al. (2021), all models are run 5 times for 1 epoch with different random seeds. For each model, we report the averaged sentence-level F1 (S-F1) and corpus-level F1 (C-F1) of its runs on each testing set.

## 4.3 Implementation Details

We use Spacy<sup>3</sup> for tokenization and keep sentences with fewer than 40 words for training due to the limited computational resources. Each video is decoded at 16 fps and  $L = 8$  video clips are sampled in total, where each clip contains  $T = 16$  frames. We train baseline models, C-PCFG and MMC-PCFG with the same hyper-parameters suggested by Kim et al. (2019a) and Zhang et al. (2021). The parsing model of PTC-PCFG has the same hyper-parameter setting as C-PCFG and MMC-PCFG (Please refer their papers for details). The constant  $\alpha$  is set to 1. We select the top 20 000 most common words in HowTo100M as vocabulary for all datasets. All baseline methods and ours are optimized by Adam (Kingma and Ba, 2015) with a learning rate of 0.001,  $\beta_1 = 0.75$  and  $\beta_2 = 0.999$ . All parameters (except the video-text matching model in PTC-PCFG) are initialized with Xavier uniform initializer (Glorot and Bengio, 2010). All our models in experiments are trained for 1 epoch with batch size of 32, without finetuning on the

<sup>2</sup>For each dataset, we randomly select 50 sentences and manually label their constituency parse trees. Benepar has S-F1 scores of 98.1% (DiDeMo), 97.2% (YouCook2) and 98.1% (MSRVTT) with manual labels.

<sup>3</sup><https://spacy.io/>

target dataset.

## 4.4 Main Results

Figure 2-4 compare our proposed PTC-PCFG approach with recently proposed state-of-the-art models: C-PCFG (Kim et al., 2019a) and MMC-PCFG<sup>4</sup> (Zhang et al., 2021). To pinpoint more fine-grained contributions, we also train these models on HowTo100M data.

**The effectiveness of HowTo100M.** We find that C-PCFG achieve better performance when they are trained with more instances from HowTo100M than the original in-domain training sets, where the largest improvements are +18.1%, +21.7% and +1.4% S-F1 scores on DiDeMo, YouCook2 and MSRVTT, respectively. These results indicate that grammar inducers are generally robust against the instances with noisy text-video correspondence. As the results, learning from noisy YouTube videos can benefit model’s overall performance and its generalization ability across multiple domains.

**The effectiveness of PTC-PCFG.** Comparing C-PCFG, MMC-PCFG and PTC-PCFG trained on different amount of HowTo100M data, we found that PTC-PCFG achieves the best performances in all three datasets. It can further improve S-F1 to +6.3% on DiDeMo, +16.7% on YouCook2 and +2.8% on MSRVTT. This demonstrates the effectiveness of the PTC-PCFG model. In particular, utilizing the video and span encoders pre-trained on a relevant tasks (e.g., video retrieval) can benefit unsupervised grammar induction.

**Performance comparison over data scale.** On DiDeMo and MSRVTT, we observe that PTC-PCFG achieves the best performance with 592k HowTo100M training samples, and further increasing the number of training instances does not improve the parsing performance on these two datasets. In contrast, the performance gain of PTC-PCFG on YouCook2 further increases with increasing training data. The reason can be that the domain of HowTo100M is closer to YouCook2 (both are instructional videos) than the other two datasets. Future work includes adding data from other sources to the whole training set more domain generic.

<sup>4</sup>Since audios are removed by HowTo100M authors, we implement MMC-PCFG with video features only, including object features(ResNeXt, SENet), action features (I3D, R2PID, S3DG), scene features, OCR features and face features.

Table 1: Performance comparison across different training set. We use HT to represent HowTo100M dataset for short, where the number in the brackets indicates the number of samples used for training. The values highlighted by **bold** and *italic* fonts indicate the top-2 methods, respectively. All numbers are shown in percentage(%). The remaining tables follow the same notations.

Method	Trainset	DiDeMo		YouCook2		MSRVTT	
		C-F1	S-F1	C-F1	S-F1	C-F1	S-F1
MMC-PCFG	DiDeMo	55.0 $\pm$ 3.7	58.9 $\pm$ 3.4	49.1 $\pm$ 4.4	53.0 $\pm$ 4.9	49.6 $\pm$ 1.4	53.8 $\pm$ 0.9
MMC-PCFG	YouCook2	40.1 $\pm$ 4.4	44.2 $\pm$ 4.4	44.7 $\pm$ 5.2	48.9 $\pm$ 5.7	34.0 $\pm$ 6.4	37.5 $\pm$ 6.8
MMC-PCFG	MSRVTT	<i>59.4</i> $\pm$ 2.9	<i>62.7</i> $\pm$ 3.3	49.6 $\pm$ 3.9	54.2 $\pm$ 4.1	<i>56.0</i> $\pm$ 1.4	60.0 $\pm$ 1.2
MMC-PCFG	HT(592k)	58.5 $\pm$ 7.3	62.4 $\pm$ 7.9	<i>53.9</i> $\pm$ 6.6	<i>58.0</i> $\pm$ 7.1	55.1 $\pm$ 7.0	<i>60.2</i> $\pm$ 8.0
<b>PTC-PCFG</b>	HT(592k)	<b>61.3</b> $\pm$ 3.9	<b>65.2</b> $\pm$ 5.3	<b>58.9</b> $\pm$ 2.5	<b>63.2</b> $\pm$ 2.3	<b>57.4</b> $\pm$ 4.6	<b>62.8</b> $\pm$ 5.7

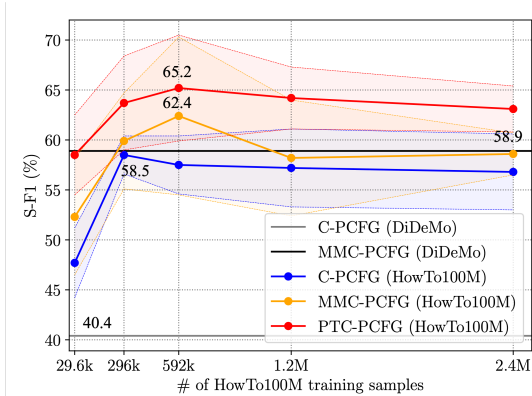


Figure 2: Performance Comparison on DiDeMo. The dotted lines and their enclosed area represent the mean and variance of each model trained on HowTo100M at different scales. We mark the highest average S-F1 achieved by each method with numbers. The remaining figures follow the same notations.

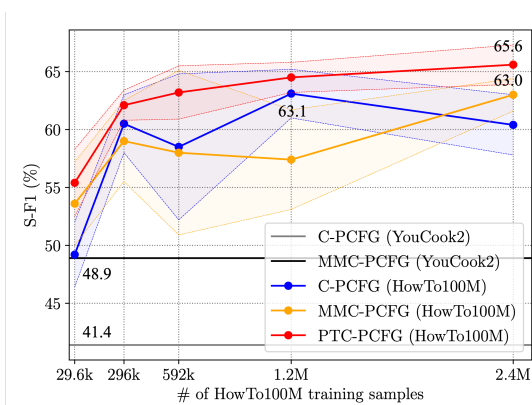


Figure 3: Performance Comparison on YouCook2.

#### 4.5 Cross-dataset Evaluation

We evaluate the robustness of models across different datasets, as shown in Table 1. Comparing MMC-PCFG trained on in-domain datasets (Row 1-3), we can observe that MMC-PCFG trained on MSRVTT achieves the best overall performance,

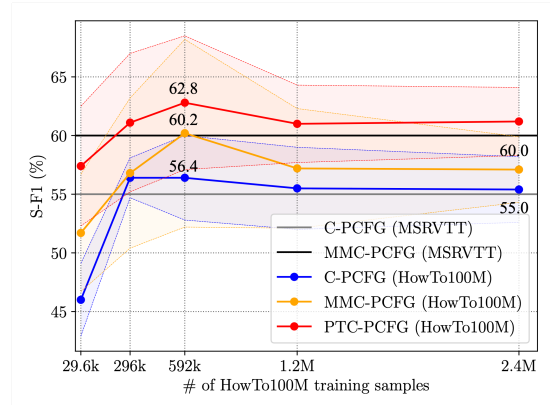


Figure 4: Performance Comparison on MSRVTT.

while MMC-PCFG trained on YouCook2 is the worst. We believe this is due to the different number of training instances<sup>5</sup> and the domain gap between different datasets. Comparing Rows 1-4, we can observe that the MMC-PCFG model trained on HT(592k) (Row 4) is the best or the second place regarding C-F1 and S-F1 compared with its variants trained on in-domain datasets (Rows 1-3). This demonstrates that the our processed video-text training instances are abundant, rich in content and can serve for general purpose. Comparing Rows 4 and 5, PTC-PCFG outperforms MMC-PCFG in both C-F1 and S-F1 in all three datasets and has smaller variance. This demonstrate that our model can leverage pre-trained video-text matching knowledge and learn consistent grammar induction.

#### 4.6 Effectiveness of Pre-Training

In this section, we explore how different pre-trained video and text encoders can affect the parsing performance, and the results are shown in Table 2. In particular, we study different

<sup>5</sup>The number of training instances in YouCook2, DiDeMo and MSRVTT are 8.9K, 32.9K and 130.2K, respectively.

Table 2: Performance comparison across different video and span encoders.

Video-Text Model		Trainset	DiDeMo		YouCook2		MSRVTT	
Video Encoder	Span Encoder		C-F1	S-F1	C-F1	S-F1	C-F1	S-F1
MIL-NCE	LSTM	HT(296k)	52.4 $\pm$ 5.5	54.4 $\pm$ 5.4	51.5 $\pm$ 5.4	56.5 $\pm$ 5.2	49.7 $\pm$ 5.5	53.4 $\pm$ 5.8
MM	LSTM	HT(296k)	53.6 $\pm$ 3.2	55.8 $\pm$ 3.1	53.1 $\pm$ 5.7	57.9 $\pm$ 5.6	48.9 $\pm$ 3.5	52.5 $\pm$ 3.6
MIL-NCE	TinyBERT	HT(296k)	54.8 $\pm$ 5.4	56.4 $\pm$ 6.0	55.7 $\pm$ 4.0	60.2 $\pm$ 3.5	52.3 $\pm$ 4.3	56.0 $\pm$ 5.0
MIL-NCE	MIL-NCE	HT(296k)	<b>59.5</b> $\pm$ 4.3	<b>63.7</b> $\pm$ 4.7	<b>57.1</b> $\pm$ 1.7	<b>62.1</b> $\pm$ 1.3	<b>55.7</b> $\pm$ 5.0	<b>61.1</b> $\pm$ 5.9
CLIP	CLIP	HT(296k)	52.9 $\pm$ 2.3	54.9 $\pm$ 2.6	53.3 $\pm$ 2.2	58.9 $\pm$ 2.1	49.1 $\pm$ 2.6	53.0 $\pm$ 2.9

video encoders<sup>6</sup>, including the S3D-based encoder from MIL-NCE (Miech et al., 2020) (MIL-NCE), the multi-modal video encoder from MMC-PCFG (Zhang et al., 2021) (MM) and the CLIP model for image-text pre-training (Radford et al., 2021) (CLIP). We also investigate various text encoders, including an LSTM encoder with random initialization (Zhang et al., 2021; Zhao and Titov, 2020), a pre-trained TinyBERT (Jiao et al., 2020) model, the text encoder from MIL-NCE (Miech et al., 2020), and the text encoder from CLIP (Radford et al., 2021).

Comparing Rows 1 with 2, we can observe that MM is better than the video encoder of MIL-NCE regarding C-F1 and S-F1 on all three datasets, as MM provides more comprehensive video features. By comparing row 1 with 3, we can also observe that TinyBERT, which is distilled from BERT (Devlin et al., 2019), outperforms the randomly initialized LSTM encoder. However, both MM and TinyBERT are independently trained only on vision or language tasks, where the vision-language correspondences are not considered during pre-training. Therefore, we further investigate the encoders jointly pre-trained on large scale multimedia datasets, including the video-text matching model MIL-NCE (Row 4) and the image-text matching model CLIP (Row 5). We can observe that by leveraging both video and text encoders in MIL-NCE can improve the parsing performance by a large margin on all three datasets. On the other hand, CLIP does not perform well, since it is designed for static images and other multi-modal information (e.g., motion) is ignored.

#### 4.7 Qualitative Analysis

In figure 5, we visualize a parser tree predicted by the best run of C-PCFG trained on MSRVTT, MMC-PCFG trained on MSRVTT, MMC-PCFG trained on HT(296k) and PTC-PCFG trained on HT(296k), as well as its reference tree. We can

<sup>6</sup> We list the video processing details in Appendix A.

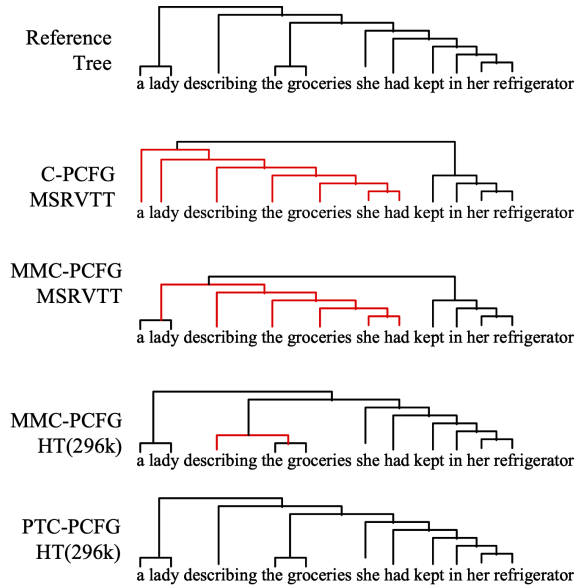


Figure 5: Parse trees predicted by different models for sentence *a lady describing the groceries she had kept in her refrigerator*. The red line shows the difference between the predicted trees and the reference tree.

observe that C-PCFG trained on MSRVTT fails at noun phrase “*a lady*”, while MMC-PCFG trained on MSRVTT succeeds. MMC-PCFG can be further improved by training on HT(296k), however, fails at noun phrase “*the groceries she had kept in her refrigerator*”. Our PTC-PCFG can leverage the pre-trained matching knowledge and make the correct prediction.

## 5 Related Work

**Grammar Induction** has a long and rich history in the computational linguistics. Earlier work (Shen et al., 2018a,b; Drozdov et al., 2019; Kim et al., 2019a; Jin et al., 2019; Yang et al., 2021a,b) on grammar induction with pure unsupervised learning showed promising results. Instead of learning purely from text, recent work improved the parsing performance with paired images (Shi et al., 2019; Zhao and Titov, 2020) or videos (Zhang et al., 2021). However, they are all limited to small



benchmarks and specified for a few domains. In contrast, our work leverages massive noisy video-subtitle pairs from YouTube without any manual annotations.

**Video Retrieval** has been a hot topic in the computer vision field for many years. Earlier approaches focused on model design (Gabeur et al., 2020; Zhang et al., 2019), while more recent approaches (Radford et al., 2021; Miech et al., 2020) focused on the pre-training on a large scale dataset and demonstrated superior zero-shot results on many downstream tasks. These models are simple in design and provide representative features with less human effort in annotations. In this work, we demonstrate that unsupervised grammar induction can also benefit from the pre-trained video-text model.

## 6 Conclusion

In this paper, we have investigated how massive instructional YouTube video and subtitle pairs can improve grammar induction. We have also proposed a new model that leverages the latest advances in multi-modal pre-training to learn better video-span correlation. Experiments on three benchmarks demonstrate superior and robust performances of our model over previous systems. We leave exploring other pre-trained video-text matching models and more publicly available data (e.g., YouTube videos from other domains and TV shows) in future work.

## 7 Limitations

Although our model faces a similar indeterminacy problem like children do, and results show that induction works even with noisy correspondence, there are a few factors which prevent this result from being directly applied to language acquisition. Our models only use instructional video and do not have the capability to interact with the world, both of which are unrealistic for human language learners. The complexity of the PCFG induction algorithm we use is cubic to the number of syntactic categories, therefore potentially limits the usefulness of larger amounts of data, where finer subcategories may be learned. Algorithms such as in Yang et al. (2021b) could be used in conjunction with multimodal inputs to examine this issue.

Following previous work, our experiments are only conducted on English video-text datasets. However, our framework is general for grammar

induction in many languages. Since our training instances are originally collected from Internet and are uploaded by users, the dataset itself might have misinformation. Meanwhile, training a model on a large-scale dataset could have high cost in energy and carbon emission. We list our computational cost of our experiments in Appendix B.

## References

- Nameera Akhtar and Lisa Montague. 1999. Early lexical acquisition: the role of cross-situational learning. *First Language*, 19(57):347–358.
- Joao Carreira and Andrew Zisserman. 2017. Quo vadis, action recognition? a new model and the kinetics dataset. In *CVPR*.
- Eugene Charniak. 2000. A maximum-entropy-inspired parser. In *NAACL*.
- Alexander Clark. 2001. Unsupervised induction of stochastic context-free grammars using distributional clustering. In *ConLL*.
- John Cocke. 1969. *Programming languages and their compilers: Preliminary notes*. New York University.
- Michael Collins. 2003. Head-driven statistical models for natural language parsing. *Computational linguistics*, 29(4):589–637.
- James Cross and Liang Huang. 2016. Span-based constituency parsing with a structure-label system and provably optimal dynamic oracles. In *EMNLP*.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. Bert: Pre-training of deep bidirectional transformers for language understanding. In *NAACL*.
- Peter F Dominey and Christelle Dodane. 2004. Indeterminacy in language acquisition: the role of child directed speech and joint attention. *Journal of Neurolinguistics*, 17(2):121–145.
- Andrew Drodov, Patrick Verga, Mohit Yadav, Mohit Iyyer, and Andrew McCallum. 2019. Unsupervised latent tree induction with deep inside-outside recursive auto-encoders. In *NAACL*.
- Valentin Gabeur, Chen Sun, Karteeek Alahari, and Cordelia Schmid. 2020. Multi-modal transformer for video retrieval. In *ECCV*.
- Dedre Gentner, Lera Boroditsky, Melissa Bowerman, and Stephen Levinson. 2001. Individuation, relativity, and early word. *Language, culture and cognition*, 3:215–256.
- Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *AISTATS*.

- Lisa Anne Hendricks, Oliver Wang, Eli Shechtman, Josef Sivic, Trevor Darrell, and Bryan Russell. 2017. Localizing moments in video with natural language. In *ICCV*.
- Jie Hu, Li Shen, and Gang Sun. 2018. Squeeze-and-excitation networks. In *CVPR*.
- Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely connected convolutional networks. In *CVPR*.
- Zhenyu Huang, Guocheng Niu, Xiao Liu, Wenbiao Ding, Xinyan Xiao, Hua Wu, and Xi Peng. 2021. Learning with noisy correspondence for cross-modal matching. In *NeurIPS*.
- Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. 2020. Tinybert: Distilling bert for natural language understanding. In *EMNLP: Findings*.
- Lifeng Jin, Finale Doshi-Velez, Timothy Miller, William Schuler, and Lane Schwartz. 2018. Unsupervised grammar induction with depth-bounded PCFG. *TACL*.
- Lifeng Jin, Finale Doshi-Velez, Timothy Miller, Lane Schwartz, and William Schuler. 2019. Unsupervised learning of PCFGs with normalizing flow. In *ACL*.
- Lifeng Jin and William Schuler. 2020. Grounded PCFG induction with images. In *AAACL-IJCNLP*.
- Mark Johnson, Thomas L Griffiths, and Sharon Goldwater. 2007. Bayesian inference for PCFGs via markov chain monte carlo. In *NAACL*.
- Tadao Kasami. 1966. An efficient recognition and syntax-analysis algorithm for context-free languages. *Coordinated Science Laboratory Report no. R-257*.
- Yoon Kim, Chris Dyer, and Alexander M Rush. 2019a. Compound probabilistic context-free grammars for grammar induction. In *ACL*.
- Yoon Kim, Alexander M Rush, Lei Yu, Adhiguna Kuncoro, Chris Dyer, and Gábor Melis. 2019b. Unsupervised recurrent neural network grammars. In *NAACL*.
- Diederik P Kingma and Jimmy Ba. 2015. Adam: A method for stochastic optimization. In *ICLR*.
- Diederik P Kingma and Max Welling. 2014. Auto-encoding variational bayes. In *ICLR*.
- Nikita Kitaev, Steven Cao, and Dan Klein. 2019. Multilingual constituency parsing with self-attention and pre-training. In *ACL*.
- Nikita Kitaev and Dan Klein. 2018. Constituency parsing with a self-attentive encoder. In *ACL*.
- Dan Klein and Christopher D Manning. 2002. A generative constituent-context model for improved grammar induction. In *ACL*.
- Kenichi Kurihara and Taisuke Sato. 2006. Variational bayesian grammar induction for natural language. In *ICGI*.
- Karim Lari and Steve J Young. 1990. The estimation of stochastic context-free grammars using the inside-outside algorithm. *Computer speech & language*, 4(1):35–56.
- Tsung-Yi Lin, Michael Maire, Serge J Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *ECCV*.
- Mitch Marcus, Beatrice Santorini, and Mary Ann Marcinkiewicz. 1993. Building a large annotated corpus of english: The penn treebank. *Computational Linguistics*, 19(2):313–330.
- Antoine Miech, Jean-Baptiste Alayrac, Lucas Smaira, Ivan Laptev, Josef Sivic, and Andrew Zisserman. 2020. End-to-end learning of visual representations from uncurated instructional videos. In *CVPR*.
- Antoine Miech, Ivan Laptev, and Josef Sivic. 2018. Learning a text-video embedding from incomplete and heterogeneous data. *arXiv preprint arXiv:1804.02516*.
- 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 *ICCV*.
- Slav Petrov and Dan Klein. 2007. Improved inference for unlexicalized parsing. In *NAACL*.
- Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. *arXiv preprint arXiv:2103.00020*.
- Florian Schroff, Dmitry Kalenichenko, and James Philbin. 2015. Facenet: A unified embedding for face recognition and clustering. In *CVPR*.
- Yikang Shen, Zhouhan Lin, Chin wei Huang, and Aaron Courville. 2018a. Neural language modeling by jointly learning syntax and lexicon. In *ICLR*.
- Yikang Shen, Shawn Tan, Alessandro Sordani, and Aaron Courville. 2018b. Ordered neurons: Integrating tree structures into recurrent neural networks. In *ICLR*.
- Haoyue Shi, Karen Livescu, and Kevin Gimpel. 2020. On the role of supervision in unsupervised constituency parsing. In *EMNLP*.
- Haoyue Shi, Jiayuan Mao, Kevin Gimpel, and Karen Livescu. 2019. Visually grounded neural syntax acquisition. In *ACL*.

- Ottokar Tilk and Tanel Alumäe. 2016. Bidirectional recurrent neural network with attention mechanism for punctuation restoration. In *Interspeech*.
- Du Tran, Heng Wang, Lorenzo Torresani, Jamie Ray, Yann LeCun, and Manohar Paluri. 2018. A closer look at spatiotemporal convolutions for action recognition. In *CVPR*.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *NIPS*.
- Pengyu Wang and Phil Blunsom. 2013. Collapsed variational bayesian inference for PCFGs. In *CoNLL*.
- Wenhai Wang, Enze Xie, Xiaoge Song, Yuhang Zang, Wenjia Wang, Tong Lu, Gang Yu, and Chunhua Shen. 2019. Efficient and accurate arbitrary-shaped text detection with pixel aggregation network. In *ICCV*.
- Saining Xie, Ross Girshick, Piotr Dollár, Zhuowen Tu, and Kaiming He. 2017. Aggregated residual transformations for deep neural networks. In *CVPR*.
- Jun Xu, Tao Mei, Ting Yao, and Yong Rui. 2016. MSR-VTT: A large video description dataset for bridging video and language. In *CVPR*.
- Songlin Yang, Yanpeng Zhao, and Kewei Tu. 2021a. Neural bi-lexicalized PCFG induction. In *ACL*.
- Songlin Yang, Yanpeng Zhao, and Kewei Tu. 2021b. PCFGs can do better: Inducing probabilistic context-free grammars with many symbols. In *NAACL*.
- Daniel H Younger. 1967. Recognition and parsing of context-free languages in time  $n^3$ . *Information and control*, 10(2):189–208.
- Haonan Yu and Jeffrey Mark Siskind. 2013. Grounded language learning from video described with sentences. In *ACL*.
- Kaipeng Zhang, Zhanpeng Zhang, Hao Wang, Zhifeng Li, Yu Qiao, and Wei Liu. 2017. Detecting faces using inside cascaded contextual cnn. In *ICCV*.
- Songyang Zhang, Linfeng Song, Lifeng Jin, Kun Xu, Dong Yu, and Jiebo Luo. 2021. Video-aided unsupervised grammar induction. In *NAACL*.
- Songyang Zhang, Jinsong Su, and Jiebo Luo. 2019. Exploiting temporal relationships in video moment localization with natural language. In *ACMMM*.
- Yue Zhang and Stephen Clark. 2011. Syntactic processing using the generalized perceptron and beam search. *Computational linguistics*, 37(1):105–151.
- Yanpeng Zhao and Ivan Titov. 2020. Visually grounded compound PCFGs. In *EMNLP*.
- Luowei Zhou, Chenliang Xu, and Jason J Corso. 2018. Towards automatic learning of procedures from web instructional videos. In *AAAI*.

## A Video Processing Details

**MIL-NCE.** Following the implementation<sup>7</sup> of MIL-NCE, we extract 1 feature per second from video encoder’s last fully connected layer. All videos are decoded at 16 frames per second (fps).

**MM.** We list the details of object, action, scene, OCR and face feature extraction as below:

- ResNeXt (Xie et al., 2017). We use the *ResNeXt101* version implemented by torchvision<sup>8</sup>. Videos are decoded at 1 fps and we extract 1 feature per second from the last fully connected layer.
- SENet (Hu et al., 2018). We use the *SENet154* version implemented by Cadene<sup>9</sup>. Videos are decoded at 1 fps and we extract 1 feature per second from the last fully connected layer.
- I3D (Carreira and Zisserman, 2017). We use the *I3D\_8x8\_R50* version implemented by SlowFast<sup>10</sup>. We decode videos at 4 fps and extract 1 feature per 2 seconds from the last fully connected layer.
- S3DG (Miech et al., 2020). We use the implementation from HERO<sup>11</sup>. Videos are decoded at 30 fps and we extract 1 feature per 1.5 seconds from the global averaged pooling layer.
- R2P1D (Tran et al., 2018). We use the *r2plus1d\_34* version implemented by torchvision. we decode videos at 16 fps and extract 1 feature per 2 seconds.
- Scene. We use *densenet161* (Huang et al., 2017) implemented by CSAILVision<sup>12</sup>. Videos are decoded at 1 fps and we extract 1 feature per second from the last fully connected layer.
- OCR. We use the text detector PANet (Wang et al., 2019) and the text recognizer *seg\_r31* implemented by MMOCR<sup>13</sup>. we decode videos at 0.5 fps and extract 1 feature per 2 seconds.
- Face. We use face detector MTCNN (Zhang et al., 2017) and face recognizer (Schroff et al., 2015) implemented by FaceNet<sup>14</sup>. We decode videos at 1 fps and extract 1 feature per second.

**CLIP.** Following the implementation<sup>15</sup> of CLIP, we extract 1 video feature per second from *ViT-B/32*’s last fully connected layer. All videos are decoded at 1 fps.

## B Computational Cost

All our models are trained on 2 32GB V100 GPUs. The approximate time cost for each run of different model is listed in Table 1. For each model, we run 5 times with different random seeds in parallel. During training, the video encoder, the span encoder and C-PCFG are involved, which contains 76.6M parameters in total. During inference, since only C-PCFG is involved, there are 23.0M parameters in total.

Table 1: The approximate training time (hours) of different model on a single run.

Model	HT(29.6k)	HT(296k)	HT(592k)	HT(1.2M)	HT(2.4M)
C-PCFG	0.07	0.7	1.5	2.9	5.8
MMC-PCFG	0.75	7.5	15	30	60
PTC-PCFG	0.50	5.0	10	20	40

<sup>7</sup>[https://github.com/antoine77340/S3D\\_HowTo100M](https://github.com/antoine77340/S3D_HowTo100M)

<sup>8</sup><https://github.com/pytorch/vision>

<sup>9</sup><https://github.com/Cadene/pretrained-models.pytorch>

<sup>10</sup><https://github.com/facebookresearch/SlowFast>

<sup>11</sup>[https://github.com/linjieli222/HERO\\_Video\\_Feature\\_Extractor](https://github.com/linjieli222/HERO_Video_Feature_Extractor)

<sup>12</sup><https://github.com/CSAILVision/places365>

<sup>13</sup><https://github.com/open-mmlab/mocr>

<sup>14</sup><https://github.com/timesler/facenet-pytorch>

<sup>15</sup><https://github.com/openai/CLIP>

## C Performance Comparison - Full Tables

We compare the performances of different models trained on different datasets. The full experiment results are demonstrated in Table 2-4. LBranch, RBranch and Random represent left branching tree, right branching tree and random tree, respectively. In addition to C-F1 and S-F1, we also evaluate the recall of each model on different phrase types, including NP, VP, PP, SBAR, ADJP and ADVP. All numbers are shown in percentage(%).

Table 2: Performance comparison on DiDeMo.

Method	Trainset	NP	VP	PP	SBAR	ADJP	ADVP	C-F1	S-F1
LBranch	None	41.7	0.1	0.1	0.7	7.2	0.0	16.2	18.5
RBranch	None	32.8	<b>91.5</b>	66.5	<b>88.2</b>	36.9	<b>63.6</b>	53.6	57.5
Random	None	36.5 $\pm$ 0.6	30.5 $\pm$ 0.5	30.1 $\pm$ 0.5	25.7 $\pm$ 2.8	29.5 $\pm$ 2.3	28.5 $\pm$ 4.8	29.4 $\pm$ 0.3	32.7 $\pm$ 0.5
C-PCFG	DiDeMo	72.9 $\pm$ 5.5	16.5 $\pm$ 6.2	23.4 $\pm$ 16.9	26.6 $\pm$ 15.9	25.0 $\pm$ 11.6	14.7 $\pm$ 12.8	38.2 $\pm$ 5.0	40.4 $\pm$ 4.1
ResNeXt	DiDeMo	64.4 $\pm$ 21.4	25.7 $\pm$ 17.7	34.6 $\pm$ 25.0	40.5 $\pm$ 26.3	16.7 $\pm$ 9.5	28.4 $\pm$ 21.3	40.0 $\pm$ 13.7	41.8 $\pm$ 14.0
SENet	DiDeMo	70.5 $\pm$ 15.3	25.7 $\pm$ 15.9	36.5 $\pm$ 24.6	36.8 $\pm$ 25.9	21.2 $\pm$ 12.5	23.6 $\pm$ 16.8	42.6 $\pm$ 10.4	44.0 $\pm$ 10.4
I3D	DiDeMo	57.9 $\pm$ 13.5	45.7 $\pm$ 14.1	45.8 $\pm$ 17.2	38.2 $\pm$ 14.8	28.4 $\pm$ 9.2	22.0 $\pm$ 9.3	45.1 $\pm$ 6.0	49.2 $\pm$ 6.0
R2P1D	DiDeMo	61.2 $\pm$ 8.5	38.1 $\pm$ 5.4	62.1 $\pm$ 4.1	61.5 $\pm$ 5.1	21.4 $\pm$ 11.4	40.8 $\pm$ 7.3	48.1 $\pm$ 4.4	50.7 $\pm$ 4.2
S3DG	DiDeMo	61.3 $\pm$ 13.4	31.7 $\pm$ 16.7	51.8 $\pm$ 8.0	50.3 $\pm$ 6.5	18.0 $\pm$ 4.5	35.2 $\pm$ 11.4	44.0 $\pm$ 2.7	46.5 $\pm$ 5.1
Scene	DiDeMo	62.2 $\pm$ 9.6	30.6 $\pm$ 12.3	41.1 $\pm$ 24.8	35.2 $\pm$ 21.9	21.4 $\pm$ 14.0	27.6 $\pm$ 17.1	41.7 $\pm$ 6.5	44.9 $\pm$ 7.4
Audio	DiDeMo	64.2 $\pm$ 18.6	21.3 $\pm$ 26.5	34.7 $\pm$ 11.0	37.3 $\pm$ 19.6	26.1 $\pm$ 4.9	18.2 $\pm$ 11.6	38.7 $\pm$ 3.7	39.5 $\pm$ 5.2
OCR	DiDeMo	64.4 $\pm$ 15.0	27.4 $\pm$ 19.5	42.8 $\pm$ 31.2	35.9 $\pm$ 20.7	14.6 $\pm$ 1.7	23.2 $\pm$ 24.0	41.9 $\pm$ 16.9	44.6 $\pm$ 17.5
Face	DiDeMo	60.8 $\pm$ 16.0	31.5 $\pm$ 17.0	52.8 $\pm$ 9.8	49.3 $\pm$ 5.6	12.6 $\pm$ 3.3	32.9 $\pm$ 14.6	43.9 $\pm$ 4.5	46.3 $\pm$ 5.5
Speech	DiDeMo	61.8 $\pm$ 12.8	26.6 $\pm$ 17.6	43.8 $\pm$ 34.5	34.2 $\pm$ 20.6	14.4 $\pm$ 4.8	12.9 $\pm$ 9.6	40.9 $\pm$ 16.0	43.1 $\pm$ 16.1
Concat	DiDeMo	68.6 $\pm$ 8.6	24.9 $\pm$ 19.9	39.7 $\pm$ 19.5	39.3 $\pm$ 19.8	10.8 $\pm$ 2.8	18.3 $\pm$ 18.1	42.2 $\pm$ 12.3	43.2 $\pm$ 14.2
MMC-PCFG	DiDeMo	67.9 $\pm$ 9.8	52.3 $\pm$ 9.0	63.5 $\pm$ 8.6	60.7 $\pm$ 10.8	34.7 $\pm$ 17.0	50.4 $\pm$ 8.3	55.0 $\pm$ 3.7	58.9 $\pm$ 3.4
MMC-PCFG	YouCook2	47.9 $\pm$ 10.4	34.6 $\pm$ 2.7	58.2 $\pm$ 12.9	19.9 $\pm$ 2.6	11.0 $\pm$ 4.0	25.5 $\pm$ 4.8	40.1 $\pm$ 4.4	44.2 $\pm$ 4.4
MMC-PCFG	MSRVTT	56.5 $\pm$ 6.8	70.8 $\pm$ 4.4	<b>82.6</b> $\pm$ 3.3	62.5 $\pm$ 6.0	42.6 $\pm$ 2.6	52.9 $\pm$ 7.6	59.4 $\pm$ 2.9	62.7 $\pm$ 3.3
C-PCFG	HT(29.6k)	74.8 $\pm$ 1.1	27.3 $\pm$ 5.7	43.6 $\pm$ 13.2	32.9 $\pm$ 7.4	32.4 $\pm$ 4.6	44.5 $\pm$ 7.8	45.7 $\pm$ 3.9	47.7 $\pm$ 3.5
MMC-PCFG	HT(29.6k)	75.4 $\pm$ 2.1	33.2 $\pm$ 12.5	54.9 $\pm$ 7.8	33.7 $\pm$ 9.8	39.2 $\pm$ 4.7	43.8 $\pm$ 7.1	49.8 $\pm$ 4.7	52.3 $\pm$ 5.8
<b>PTC-PCFG</b>	HT(29.6k)	66.0 $\pm$ 9.4	53.7 $\pm$ 13.9	68.2 $\pm$ 4.5	50.4 $\pm$ 5.0	35.2 $\pm$ 4.2	52.7 $\pm$ 8.9	54.9 $\pm$ 3.4	58.5 $\pm$ 4.0
C-PCFG	HT(296k)	81.4 $\pm$ 1.6	36.4 $\pm$ 6.9	67.0 $\pm$ 1.9	45.9 $\pm$ 3.5	46.5 $\pm$ 4.8	49.9 $\pm$ 8.3	55.6 $\pm$ 1.4	58.5 $\pm$ 1.9
MMC-PCFG	HT(296k)	81.9 $\pm$ 2.1	42.7 $\pm$ 15.0	65.3 $\pm$ 6.4	39.1 $\pm$ 8.5	48.0 $\pm$ 8.1	43.7 $\pm$ 6.7	57.1 $\pm$ 4.2	59.9 $\pm$ 4.8
<b>PTC-PCFG</b>	HT(296k)	76.0 $\pm$ 4.9	55.3 $\pm$ 11.9	70.7 $\pm$ 6.4	53.7 $\pm$ 9.5	43.4 $\pm$ 4.8	47.2 $\pm$ 13.3	59.5 $\pm$ 4.3	63.7 $\pm$ 4.7
C-PCFG	HT(592k)	82.4 $\pm$ 1.9	37.6 $\pm$ 9.2	63.1 $\pm$ 5.8	37.4 $\pm$ 7.9	45.8 $\pm$ 5.5	53.8 $\pm$ 9.9	55.5 $\pm$ 2.7	57.5 $\pm$ 2.9
MMC-PCFG	HT(592k)	79.0 $\pm$ 5.7	52.5 $\pm$ 19.1	64.3 $\pm$ 6.1	45.5 $\pm$ 9.9	44.1 $\pm$ 4.4	44.5 $\pm$ 11.7	58.5 $\pm$ 7.3	62.4 $\pm$ 7.9
<b>PTC-PCFG</b>	HT(592k)	79.1 $\pm$ 2.9	55.3 $\pm$ 18.4	73.7 $\pm$ 5.1	50.6 $\pm$ 8.2	38.3 $\pm$ 5.6	47.2 $\pm$ 5.5	<b>61.3</b> $\pm$ 3.9	<b>65.2</b> $\pm$ 5.3
C-PCFG	HT(1.2M)	82.9 $\pm$ 1.5	32.2 $\pm$ 12.5	67.5 $\pm$ 2.1	41.3 $\pm$ 8.5	45.6 $\pm$ 4.8	48.8 $\pm$ 8.0	55.1 $\pm$ 3.3	57.2 $\pm$ 3.9
MMC-PCFG	HT(1.2M)	<b>83.4</b> $\pm$ 2.6	43.2 $\pm$ 15.9	51.3 $\pm$ 19.1	36.2 $\pm$ 7.0	<b>49.4</b> $\pm$ 5.5	56.2 $\pm$ 6.9	55.1 $\pm$ 5.0	58.2 $\pm$ 5.8
<b>PTC-PCFG</b>	HT(1.2M)	77.9 $\pm$ 1.9	49.0 $\pm$ 7.6	78.3 $\pm$ 3.9	47.2 $\pm$ 6.5	35.8 $\pm$ 8.7	49.8 $\pm$ 11.6	60.0 $\pm$ 2.3	64.2 $\pm$ 3.1
C-PCFG	HT(2.4M)	<b>83.4</b> $\pm$ 2.2	31.0 $\pm$ 4.7	68.1 $\pm$ 8.4	48.5 $\pm$ 7.7	46.9 $\pm$ 7.3	45.9 $\pm$ 8.2	55.2 $\pm$ 3.0	56.8 $\pm$ 3.8
MMC-PCFG	HT(2.4M)	81.7 $\pm$ 2.8	37.6 $\pm$ 4.8	71.1 $\pm$ 6.9	47.2 $\pm$ 6.1	41.7 $\pm$ 8.1	51.2 $\pm$ 7.2	57.0 $\pm$ 2.1	58.6 $\pm$ 2.1
<b>PTC-PCFG</b>	HT(2.4M)	78.9 $\pm$ 1.6	49.4 $\pm$ 6.9	75.0 $\pm$ 5.2	48.7 $\pm$ 6.0	36.9 $\pm$ 5.4	51.8 $\pm$ 13.8	60.0 $\pm$ 2.5	63.1 $\pm$ 2.3

Table 3: Performance comparison on YouCook2.

Method	Trainset	NP	VP	PP	SBAR	ADJP	ADVP	C-F1	S-F1
LBranch	None	1.7	42.8	0.4	8.1	1.5	0.0	6.8	5.9
RBranch	None	35.6	47.5	67.0	<b>88.9</b>	33.9	65.0	35.0	41.6
Random	None	27.2 $\pm$ 0.3	27.1 $\pm$ 1.4	29.9 $\pm$ 0.5	31.3 $\pm$ 5.2	26.9 $\pm$ 7.7	26.2 $\pm$ 11.9	21.2 $\pm$ 0.2	24.0 $\pm$ 0.2
C-PCFG	YouCook2	47.4 $\pm$ 18.4	49.4 $\pm$ 11.9	58.0 $\pm$ 22.6	45.7 $\pm$ 6.0	27.7 $\pm$ 15.1	36.2 $\pm$ 7.4	37.8 $\pm$ 6.7	41.4 $\pm$ 6.6
ResNeXt	YouCook2	46.5 $\pm$ 13.7	40.8 $\pm$ 9.8	67.9 $\pm$ 12.7	50.5 $\pm$ 13.3	22.3 $\pm$ 6.7	38.8 $\pm$ 21.3	38.2 $\pm$ 8.3	42.8 $\pm$ 8.4
SENet	YouCook2	48.3 $\pm$ 14.4	40.7 $\pm$ 9.2	73.6 $\pm$ 11.2	45.5 $\pm$ 17.0	26.9 $\pm$ 13.6	41.2 $\pm$ 17.5	39.9 $\pm$ 8.7	44.9 $\pm$ 8.3
I3D	YouCook2	48.1 $\pm$ 10.7	39.0 $\pm$ 8.0	79.4 $\pm$ 8.4	50.0 $\pm$ 14.9	18.5 $\pm$ 7.0	41.2 $\pm$ 4.1	40.6 $\pm$ 3.6	45.7 $\pm$ 3.2
R2P1D	YouCook2	52.4 $\pm$ 10.9	33.7 $\pm$ 16.4	66.7 $\pm$ 10.7	49.5 $\pm$ 13.8	25.8 $\pm$ 10.6	33.8 $\pm$ 12.4	39.4 $\pm$ 8.1	44.4 $\pm$ 8.3
S3DG	YouCook2	50.4 $\pm$ 13.1	32.6 $\pm$ 16.3	71.7 $\pm$ 7.5	33.3 $\pm$ 5.9	30.8 $\pm$ 17.5	40.0 $\pm$ 7.1	39.3 $\pm$ 6.5	44.1 $\pm$ 6.6
Audio	YouCook2	51.2 $\pm$ 3.1	42.0 $\pm$ 7.2	61.5 $\pm$ 18.0	51.0 $\pm$ 14.8	23.5 $\pm$ 16.8	48.8 $\pm$ 8.2	39.2 $\pm$ 4.7	43.3 $\pm$ 4.9
OCR	YouCook2	48.6 $\pm$ 8.1	41.5 $\pm$ 4.1	65.5 $\pm$ 17.4	39.9 $\pm$ 4.4	18.5 $\pm$ 6.6	53.8 $\pm$ 14.7	38.6 $\pm$ 5.5	43.2 $\pm$ 5.6
Concat	YouCook2	50.3 $\pm$ 10.3	42.3 $\pm$ 2.9	81.6 $\pm$ 8.7	40.1 $\pm$ 3.9	17.7 $\pm$ 8.2	52.5 $\pm$ 5.6	42.3 $\pm$ 5.7	47.0 $\pm$ 5.6
MMC-PCFG	YouCook2	62.7 $\pm$ 9.8	45.3 $\pm$ 2.8	63.4 $\pm$ 17.7	43.9 $\pm$ 4.8	26.2 $\pm$ 7.5	35.0 $\pm$ 3.5	44.7 $\pm$ 5.2	48.9 $\pm$ 5.7
MMC-PCFG	DiDeMo	63.8 $\pm$ 4.5	62.1 $\pm$ 7.4	70.7 $\pm$ 9.0	56.8 $\pm$ 9.2	35.4 $\pm$ 7.2	51.2 $\pm$ 4.1	49.1 $\pm$ 4.4	53.0 $\pm$ 4.9
MMC-PCFG	MSRVTT	63.1 $\pm$ 9.2	51.5 $\pm$ 7.3	82.7 $\pm$ 1.5	64.9 $\pm$ 10.8	30.4 $\pm$ 3.0	40.0 $\pm$ 6.1	49.6 $\pm$ 3.9	54.2 $\pm$ 4.1
C-PCFG	HT(29.6k)	68.2 $\pm$ 3.3	43.5 $\pm$ 4.9	48.8 $\pm$ 20.2	34.5 $\pm$ 4.2	28.5 $\pm$ 3.2	56.7 $\pm$ 6.2	44.3 $\pm$ 2.4	49.2 $\pm$ 2.8
MMC-PCFG	HT(29.6k)	68.8 $\pm$ 3.0	51.3 $\pm$ 11.5	65.5 $\pm$ 11.2	33.7 $\pm$ 6.3	32.0 $\pm$ 4.3	55.0 $\pm$ 11.3	48.8 $\pm$ 3.5	53.6 $\pm$ 3.6
<b>PTC-PCFG</b>	HT(29.6k)	66.3 $\pm$ 4.2	55.5 $\pm$ 8.3	76.5 $\pm$ 5.7	46.2 $\pm$ 9.0	33.4 $\pm$ 7.4	40.0 $\pm$ 9.7	50.2 $\pm$ 3.4	55.4 $\pm$ 2.9
C-PCFG	HT(296k)	77.5 $\pm$ 3.1	56.6 $\pm$ 4.5	74.4 $\pm$ 5.5	51.1 $\pm$ 9.0	40.5 $\pm$ 6.4	63.3 $\pm$ 8.5	55.0 $\pm$ 2.7	60.5 $\pm$ 2.5
MMC-PCFG	HT(296k)	76.5 $\pm$ 4.1	61.8 $\pm$ 9.7	67.2 $\pm$ 8.7	34.0 $\pm$ 16.2	41.0 $\pm$ 7.9	68.3 $\pm$ 8.2	53.8 $\pm$ 3.4	59.0 $\pm$ 3.5
<b>PTC-PCFG</b>	HT(296k)	77.5 $\pm$ 2.5	65.8 $\pm$ 6.1	78.9 $\pm$ 7.0	61.6 $\pm$ 3.9	42.4 $\pm$ 6.3	60.0 $\pm$ 3.3	57.1 $\pm$ 1.7	62.1 $\pm$ 1.3
C-PCFG	HT(592k)	76.7 $\pm$ 5.6	64.7 $\pm$ 12.3	65.4 $\pm$ 22.5	45.2 $\pm$ 13.9	45.1 $\pm$ 9.4	68.3 $\pm$ 6.2	54.2 $\pm$ 6.0	58.5 $\pm$ 6.3
MMC-PCFG	HT(592k)	72.5 $\pm$ 12.7	68.9 $\pm$ 7.0	68.3 $\pm$ 18.2	54.3 $\pm$ 5.7	<b>50.2</b> $\pm$ 2.1	75.0 $\pm$ 9.1	53.9 $\pm$ 6.6	58.0 $\pm$ 7.1
<b>PTC-PCFG</b>	HT(592k)	78.7 $\pm$ 5.3	69.9 $\pm$ 3.6	80.5 $\pm$ 2.8	58.9 $\pm$ 12.3	43.2 $\pm$ 4.0	65.0 $\pm$ 6.2	58.9 $\pm$ 2.5	63.2 $\pm$ 2.3
C-PCFG	HT(1.2M)	80.1 $\pm$ 2.5	63.5 $\pm$ 11.9	78.5 $\pm$ 2.5	52.5 $\pm$ 13.7	42.9 $\pm$ 6.5	66.7 $\pm$ 7.5	58.1 $\pm$ 2.4	63.1 $\pm$ 2.1
MMC-PCFG	HT(1.2M)	75.7 $\pm$ 2.2	64.8 $\pm$ 9.7	57.7 $\pm$ 21.4	51.2 $\pm$ 9.4	49.3 $\pm$ 4.9	75.0 $\pm$ 7.5	52.5 $\pm$ 4.0	57.4 $\pm$ 4.3
<b>PTC-PCFG</b>	HT(1.2M)	78.1 $\pm$ 2.6	<b>72.2</b> $\pm$ 4.1	<b>85.8</b> $\pm$ 5.1	69.3 $\pm$ 6.6	41.5 $\pm$ 8.7	66.7 $\pm$ 10.5	60.1 $\pm$ 1.4	64.5 $\pm$ 1.3
C-PCFG	HT(2.4M)	75.9 $\pm$ 3.9	61.5 $\pm$ 8.0	78.2 $\pm$ 4.1	59.7 $\pm$ 13.4	45.4 $\pm$ 5.6	75.0 $\pm$ 5.3	55.9 $\pm$ 2.6	60.4 $\pm$ 2.6
MMC-PCFG	HT(2.4M)	78.0 $\pm$ 2.8	69.8 $\pm$ 4.1	79.2 $\pm$ 5.4	50.6 $\pm$ 14.2	41.5 $\pm$ 5.1	71.7 $\pm$ 10.0	58.3 $\pm$ 1.8	63.0 $\pm$ 1.4
<b>PTC-PCFG</b>	HT(2.4M)	<b>81.9</b> $\pm$ 3.1	71.5 $\pm$ 4.5	83.0 $\pm$ 4.1	59.0 $\pm$ 15.5	40.7 $\pm$ 3.4	<b>81.7</b> $\pm$ 6.2	<b>61.1</b> $\pm$ 2.0	<b>65.6</b> $\pm$ 1.7

Table 4: Performance comparison on MSRVTT.

Method	Trainset	NP	VP	PP	SBAR	ADJP	ADVP	C-F1	S-F1
LBranch	None	34.6	0.1	0.9	0.2	3.8	0.3	14.4	16.8
RBranch	None	34.6	<b>90.9</b>	67.5	<b>94.8</b>	25.4	54.8	54.2	58.6
Random	None	34.6 $\pm$ 0.1	26.8 $\pm$ 0.1	28.1 $\pm$ 0.2	24.6 $\pm$ 0.3	24.8 $\pm$ 1.0	28.1 $\pm$ 1.4	27.2 $\pm$ 0.1	30.5 $\pm$ 0.1
C-PCFG	MSRVTT	46.6 $\pm$ 3.2	61.1 $\pm$ 3.3	72.5 $\pm$ 8.3	63.7 $\pm$ 4.0	33.1 $\pm$ 7.1	67.1 $\pm$ 4.7	50.7 $\pm$ 3.2	55.0 $\pm$ 3.2
ResNeXt	MSRVTT	48.6 $\pm$ 3.0	59.0 $\pm$ 6.0	72.0 $\pm$ 3.6	62.1 $\pm$ 5.2	32.6 $\pm$ 2.5	70.4 $\pm$ 6.4	50.7 $\pm$ 1.7	54.9 $\pm$ 2.2
SENet	MSRVTT	49.0 $\pm$ 4.4	63.5 $\pm$ 6.4	71.7 $\pm$ 4.8	60.9 $\pm$ 10.6	34.0 $\pm$ 6.4	<i>74.1</i> $\pm$ 1.9	52.2 $\pm$ 1.2	56.0 $\pm$ 1.6
I3D	MSRVTT	53.9 $\pm$ 10.5	63.2 $\pm$ 9.1	73.7 $\pm$ 2.9	65.3 $\pm$ 9.1	35.0 $\pm$ 6.8	73.8 $\pm$ 4.1	54.5 $\pm$ 1.6	59.1 $\pm$ 1.7
R2P1D	MSRVTT	52.8 $\pm$ 3.6	63.3 $\pm$ 4.6	73.1 $\pm$ 10.1	66.9 $\pm$ 2.0	34.0 $\pm$ 2.2	72.5 $\pm$ 4.2	54.0 $\pm$ 2.5	58.0 $\pm$ 2.3
S3DG	MSRVTT	48.2 $\pm$ 4.4	60.4 $\pm$ 3.9	71.4 $\pm$ 6.4	58.1 $\pm$ 8.2	25.3 $\pm$ 2.2	61.8 $\pm$ 8.4	50.7 $\pm$ 3.2	54.7 $\pm$ 2.9
Scene	MSRVTT	50.7 $\pm$ 1.6	65.0 $\pm$ 4.7	78.6 $\pm$ 3.6	<i>67.3</i> $\pm$ 3.9	34.5 $\pm$ 4.6	71.7 $\pm$ 1.8	54.6 $\pm$ 1.5	58.4 $\pm$ 1.3
Audio	MSRVTT	50.0 $\pm$ 1.1	63.7 $\pm$ 6.1	72.7 $\pm$ 3.0	61.9 $\pm$ 6.5	34.5 $\pm$ 2.3	68.0 $\pm$ 5.9	52.8 $\pm$ 1.3	56.7 $\pm$ 1.4
OCR	MSRVTT	48.3 $\pm$ 8.3	57.1 $\pm$ 4.6	<i>76.9</i> $\pm$ 0.6	60.7 $\pm$ 4.9	33.9 $\pm$ 8.3	72.1 $\pm$ 4.4	51.0 $\pm$ 3.0	55.5 $\pm$ 3.0
Face	MSRVTT	46.5 $\pm$ 6.8	61.3 $\pm$ 3.6	71.5 $\pm$ 7.1	60.8 $\pm$ 11.0	30.9 $\pm$ 3.4	68.4 $\pm$ 6.0	50.5 $\pm$ 2.6	54.5 $\pm$ 2.6
Speech	MSRVTT	48.5 $\pm$ 7.6	60.7 $\pm$ 3.5	74.5 $\pm$ 5.7	62.6 $\pm$ 6.2	27.3 $\pm$ 1.8	74.0 $\pm$ 3.1	51.7 $\pm$ 2.6	56.2 $\pm$ 2.5
Concat	MSRVTT	43.6 $\pm$ 6.0	64.7 $\pm$ 3.0	68.5 $\pm$ 8.0	63.8 $\pm$ 3.8	32.0 $\pm$ 5.5	70.4 $\pm$ 5.9	49.8 $\pm$ 4.1	54.2 $\pm$ 4.0
MMC-PCFG	MSRVTT	52.3 $\pm$ 5.1	<i>68.1</i> $\pm$ 2.9	<b>78.2</b> $\pm$ 1.9	65.8 $\pm$ 2.4	32.0 $\pm$ 2.0	<b>74.7</b> $\pm$ 2.3	56.0 $\pm$ 1.4	60.0 $\pm$ 1.2
MMC-PCFG	DiDeMo	61.8 $\pm$ 7.7	41.5 $\pm$ 11.8	64.6 $\pm$ 5.2	47.1 $\pm$ 11.1	30.5 $\pm$ 7.1	62.2 $\pm$ 5.1	49.6 $\pm$ 1.4	53.8 $\pm$ 0.9
MMC-PCFG	YouCook2	40.7 $\pm$ 14.9	23.9 $\pm$ 3.4	59.9 $\pm$ 10.2	16.2 $\pm$ 2.6	14.5 $\pm$ 4.0	23.7 $\pm$ 3.9	34.0 $\pm$ 6.4	37.5 $\pm$ 6.8
C-PCFG	HT(29.6k)	68.6 $\pm$ 2.1	25.1 $\pm$ 4.8	37.5 $\pm$ 12.1	33.4 $\pm$ 3.6	27.7 $\pm$ 2.9	41.9 $\pm$ 5.0	42.3 $\pm$ 3.3	46.0 $\pm$ 3.1
MMC-PCFG	HT(29.6k)	70.8 $\pm$ 2.7	32.2 $\pm$ 13.2	48.6 $\pm$ 6.3	36.0 $\pm$ 4.6	32.0 $\pm$ 2.1	43.1 $\pm$ 5.5	47.2 $\pm$ 3.8	51.7 $\pm$ 5.0
<b>PTC-PCFG</b>	HT(29.6k)	62.2 $\pm$ 7.7	54.0 $\pm$ 13.0	60.0 $\pm$ 4.9	53.0 $\pm$ 3.6	32.4 $\pm$ 4.0	45.9 $\pm$ 5.1	52.2 $\pm$ 3.9	57.4 $\pm$ 5.1
C-PCFG	HT(296k)	75.5 $\pm$ 1.4	34.8 $\pm$ 4.3	58.6 $\pm$ 1.6	46.9 $\pm$ 2.8	40.0 $\pm$ 3.1	55.4 $\pm$ 7.1	52.0 $\pm$ 1.3	56.4 $\pm$ 1.7
MMC-PCFG	HT(296k)	75.1 $\pm$ 2.5	39.4 $\pm$ 15.5	55.2 $\pm$ 7.2	40.0 $\pm$ 4.5	40.2 $\pm$ 5.8	51.0 $\pm$ 6.4	52.4 $\pm$ 5.5	56.8 $\pm$ 6.4
<b>PTC-PCFG</b>	HT(296k)	70.2 $\pm$ 5.8	51.7 $\pm$ 12.1	64.5 $\pm$ 6.2	54.0 $\pm$ 6.5	39.2 $\pm$ 3.2	54.8 $\pm$ 9.5	55.7 $\pm$ 5.0	61.1 $\pm$ 5.9
C-PCFG	HT(592k)	76.9 $\pm$ 2.6	35.4 $\pm$ 8.3	57.9 $\pm$ 5.4	44.5 $\pm$ 10.2	<b>41.4</b> $\pm$ 3.9	57.8 $\pm$ 5.3	52.5 $\pm$ 3.4	56.4 $\pm$ 3.6
MMC-PCFG	HT(592k)	76.1 $\pm$ 3.4	46.3 $\pm$ 20.0	57.0 $\pm$ 5.8	50.1 $\pm$ 10.1	37.9 $\pm$ 3.1	52.8 $\pm$ 4.6	55.1 $\pm$ 7.0	60.2 $\pm$ 8.0
<b>PTC-PCFG</b>	HT(592k)	74.0 $\pm$ 3.6	50.2 $\pm$ 18.9	67.0 $\pm$ 4.1	54.5 $\pm$ 9.1	34.7 $\pm$ 2.4	55.4 $\pm$ 7.3	<b>57.4</b> $\pm$ 4.6	<b>62.8</b> $\pm$ 5.7
C-PCFG	HT(1.2M)	77.0 $\pm$ 2.0	30.5 $\pm$ 10.4	60.1 $\pm$ 3.4	41.5 $\pm$ 11.1	38.5 $\pm$ 4.4	52.2 $\pm$ 4.8	51.6 $\pm$ 3.1	55.5 $\pm$ 3.5
MMC-PCFG	HT(1.2M)	<b>77.9</b> $\pm$ 2.2	44.2 $\pm$ 13.1	40.6 $\pm$ 22.2	45.6 $\pm$ 5.5	<i>40.5</i> $\pm$ 4.6	56.3 $\pm$ 6.8	52.4 $\pm$ 4.5	57.2 $\pm$ 5.1
<b>PTC-PCFG</b>	HT(1.2M)	72.3 $\pm$ 1.6	44.0 $\pm$ 7.7	70.2 $\pm$ 4.6	53.4 $\pm$ 7.9	34.4 $\pm$ 4.6	57.4 $\pm$ 6.9	55.6 $\pm$ 2.5	61.0 $\pm$ 3.3
C-PCFG	HT(2.4M)	77.2 $\pm$ 2.1	31.1 $\pm$ 4.6	59.8 $\pm$ 8.0	43.3 $\pm$ 9.7	39.1 $\pm$ 3.5	54.5 $\pm$ 6.4	51.9 $\pm$ 2.3	55.4 $\pm$ 2.8
MMC-PCFG	HT(2.4M)	76.2 $\pm$ 2.2	36.1 $\pm$ 6.3	62.5 $\pm$ 7.5	47.8 $\pm$ 8.1	40.0 $\pm$ 4.5	55.8 $\pm$ 5.3	53.5 $\pm$ 2.4	57.1 $\pm$ 2.8
<b>PTC-PCFG</b>	HT(2.4M)	74.2 $\pm$ 2.8	46.0 $\pm$ 7.1	67.5 $\pm$ 4.1	52.7 $\pm$ 8.6	40.2 $\pm$ 5.6	58.3 $\pm$ 10.7	<i>56.6</i> $\pm$ 2.5	<i>61.2</i> $\pm$ 2.9