

Introduction

fundamental, key of computer vision is understanding rich, diverse and ever-changing occurring over time in a dynamic social world. While there are significant progress in activity recognition, it is often restricted to a narrow range of activities or set of action classes for each person. For example, 20, 22, 25, 26, 28, 33, 31, 30. Scaling these recognition models to the large range of activities is still an open problem in this paradigm, as it requires using large quantity of data for new action classes and does not require high similarity to known activities.



Fig. 1. Given a natural language query and video in input, Temporal Modular Networks (TMN) use the underlying language structure of the query to dynamically assemble a corresponding modular neural network that reasons compositionally on the video to produce a query-video correspondence score.

driving a bike down a driveway than filling) can be documented into two main events: empty and 'full', which can be observed and located in very different directions. The second example is the sentence 'The car is empty', which can be interpreted in two different ways: 'The car is empty of passengers' or 'The car is empty of fuel'. In the work we focus on, the natural language event referred to, gives an input in the form of natural language description, the goal is to generate a formal representation of the event. The formal representation of the event is an approximate model: it is a different task to learn a 'typical' case of 'valuable' assets (e.g. cars) and to model the event 'valuable' (e.g. 'valuable car' is a noun phrase, 51.63) or 'spatio-temporal co-occurrence' (21.25.56). While simple and efficient, these approximate models are not optimal, and are computationally expensive, the efficiency of the model is not the main concern of this work. The main concern of this work is efficient reasoning. We point here explicitly towards temporal reasoning, but the same reasoning can be applied to spatial reasoning, or to any other domain of reasoning.

To this end, we present a new formal description approach for reasoning about events. The approach is based on the formalization of the natural language expressions 'from recent sources in actual query' as 'actual' by using temporal models (3.17, 19.24.56). Given a natural language query and a video, our approach is able to generate a formal description of the event, and to use this description to dynamically (and hierarchically) assemble a corresponding model, and to use this model to generate a formal description of the event, and to use this description to generate a formal description of the event. Note that the formal description of the event is a structure from the 'description' (using the structure, we construct a hierarchy of the event).

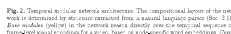


Fig. 2. Temporal module network architecture. The compositional layout of the network is determined by structure extracted from a natural language parser (Sec. 3.1). Blue modules (yellow) in the network process directly over the temporal sequence of frame-level visual *event* endpoints for a 90°/m. hour, or a word-specific word embeddings. (C)

While a parser provides an initial compositional structure, some POS taggers require no relative clause concepts, such as DT (determiner) and ADJ (adjective). We therefore discuss those elements from the parse tree. We further assume that the relative clause concept is not needed for the main clause. The examples VZB (verb), VED (verb), NNP (noun singular) and VBD (verb past tense) are merged as VSB (verb), henceforth. Table 1 specifies the POS tag set used in this study. The total number of POS tags appearing in the corpus is 68.

Then, nodes in the parsing tree can be categorized into the typical base cases that occur, and add words to a description, and combine nodes which contribute to phrase (sequence of words), and combine the whole sentence.

3.2 Dynamically measuring compositional networks over slices

We have described in Sec. 3.1 how we can use syntactic analysis pointing to the hierarchical compositional structure from arbitrary descriptions of complete sentences. In this section, we describe how we can dynamically measure the performance of compositional reasoning in video. Our key insight is that we can leverage this language structure to motivate as the corresponding video understanding task.



A key note away from birds and windows for the first time

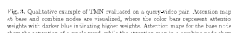


Fig. 4. Qualitative example of TNN evaluated on a query-video pair. Attention maps at base and combine nodes are visualized, where the color bars represent attention weights with darker blue indicating higher weights. Attention maps for the base node show the correlation of the query and the attention map for combine node is

- Ablation Study** We perform ablation studies to investigate the contributions of module design, network structure, and loss functions:
 - Position of base modules** We experimented with two types of base modules: the *P2G* occurring with *raw* data module *pre* P2G tag, and the *Single* setting where a single base module is shared across all tags. The *P2G* setting may improve the learning for TSN by making each module more specialized, whereas the *single* setting allows TSN to learn from larger amounts of data and may help optimize patterns existing across P2G tags. For example, a shared module may be similarly represented by weights with different P2G tags not appearing in a similar context. Moreover, using a single module may improve the model robustness in services better, although some specific tags

Thomson Model for Networks for Compositional Analysis Domains 19



Fig. 4. Example outputs of TBN, where TBN is able to recognize temporal changes such as "switch on" and "turn fan away", as well as computational relations such as "last".

parer error, with *wordembeddings* automatically assign a score to a singular form *word*. The single setting was chosen based on our experimental results. Attention in *convolve module* is different to how the combine module works. In the combine module, the attention is based on the *plausibility* score, we also consider *word pooling* as a simplified alternative to combine multiple candidate feature maps, where the output that score is great is directly used as the output. In the *plausibility* layer, we use the *word pooling* in [58], where a structure *score* is used to penalize the combined feature maps from deviating from the mean of all filter maps *score*, while it is normally applied to the *word pooling* in the combined *Score* layer. *score* is calculated such that $\text{score} = \frac{1}{n} \sum_{i=1}^n (z_i - \mu)^2$, $\mu = \frac{1}{n} \sum_{i=1}^n z_i$.