# Cooking in the kitchen: Recognizing and Segmenting Human Activities in Videos

Hilde Kuehne · Juergen Gall · Thomas Serre

**Abstract** As the field of action recognition matures, research is rapidly moving away from simpler problems such as action recognition in short hand segmented video segments to more complex real-world problems such as the continuous monitoring and analysis of daily human activities.

We propose an end-to-end generative approach for the segmentation and parsing of complex human activities. In this approach, a visual representation based on reduced Fisher Vectors is combined with a structured generative temporal model for recognition. To overcome one of the major limitations of generative models, i.e., their need for large amount of training data, we recorded a large scale activity dataset featuring 52 participants preparing 10 distinct dishes in their own kitchen. We annotated the resulting video at both a coarse and fine level of action granularity. The dataset was used to evaluate the proposed approach on various tasks ranging from basic action unit classification to activity recognition as well as the segmentation and parsing of video sequences. Our results demonstrate the ability of structured temporal generative approaches to cope with the complexity of daily-life activities.

Hilde Kuehne
University of Bonn
E-mail: kuehne@iai.uni-bonn.de

Juergen Gall
University of Bonn
E-mail: gall@iai.uni-bonn.de

Thomas Serre
Brown University
E-mail: thomas_serre@brown.edu

## 1 Introduction

Several real-world applications of computer vision including smart homes, surveillance and assisted living require continuous video monitoring. However, to date, most of the work on action recognition still focuses on the somewhat simpler problem of assigning class labels to short pre-segmented video clips. In comparison, methods for the automated analysis of temporal structures, including methods for parsing and segmentation are still in their early stages of development. Progress in action recognition has been largely spurred by the increasing availability of large realistic video datasets that allow the benchmarking of different visual representations and classification methodologies. Unfortunately, there is currently a dire need for similar large-scale datasets comprising long video sequences of complex activities recorded "in the wild".

To fill in this void, we describe a novel human-activity video dataset that we named the Breakfast dataset. This dataset includes almost 70 hours of hand-annotated videos corresponding to 52 unique participants preparing 10 distinct breakfast dishes in 18 different home kitchens. The dataset provides annotations at two different temporal scales, a finer scale, corresponding to low-level task-oriented motion sequences such as "open drawer" - "reach knife" - "carry knife", and a coarser scale corresponding to higher-level goal-oriented action sequences, such as "take plate" - "cut fruit" etc. The proposed dataset is currently the largest available dataset for human activity recognition and will allow for benchmarking of various activity recognition approaches as well as the parsing, segmentation and detection of activities into finer action units. In addition to supporting general research in computer vision, we hope that our dataset will contribute the test-

ing of brain theories of event perception. A body of the cognitive psychology literature is devoted to understanding how people perceive human actions over time and, in particular, what brain mechanisms support the perception of human movements. In this context, it has been shown that activities are not perceived as a continuous input stream but rather as discrete action units within a behavioral sequence [1]. The segmentation of action streams and the detection of discrete boundaries between action units happens at different levels of granularity such that action segments get combined over time to form a holistic interpretation of perceived actions [47]. This abstraction has also been shown to be a condition for people to predict and react to others' intentions (see, e.g., [46]), and can thus be seen as a fundamental ability towards our understanding of going activities. We thus hope that the release of a large activity dataset together with consistent behavioral annotations at both fine and coarse levels of granularity combined with visual representations derived from state of the art computer vision systems will help further our understanding of the brain mechanisms underlying event perception.

We further describe a novel generative framework for the analysis of temporal structures. In this framework, Fisher Vectors (FVs) are used to represent individual video frames and action units are modeled by Hidden Markov Models (HMMs). Action units are combined through an activity grammar that is learned from data. We extensively evaluate the resulting approach on the proposed Breakfast dataset as well as four other benchmarks. Our evaluation, which includes activity recognition, action unit classification and segmentation, demonstrates that the approach performs on par or better than the state-of-the-art.

A preliminary version of this work appeared in [12, 13]. The present work extends our original release of the Breakfast dataset [12] with additional fine-grained annotations. We also provide a more comprehensive overview of the generative framework used to analyze temporal structures first presented in [13]. The present evaluation has also been extended compared to that in [13] to include action unit classification and segmentation at different levels of granularity, activity recognition as well as a complete run-time analysis of the system and its dependency on the amount of training data available.

## 2 Related Work

Before discussing structured temporal models and activity recognition datasets in Sections 2.2 and 2.3, we

first briefly review the state-of-the-art in video-clip classification 2.1.

### 2.1 Unstructured approaches to video-clip classification

The best current approaches to human action recognition in video clips rely on dense trajectory features [37] that are then quantized using the Fisher Vector (FV) method. The combination of Fisher vector encoding and dense trajectories was first described in [38, 19] and shown to achieve state-of-the-art classification accuracy on several action datasets. The approach was further improved in [21] using stacked FVs. In addition to FVs, it was shown that the accuracy of this approach could be improved by modelling the context of an action. This was done in [10] via the detection of objects in the scene using a convolutional neural network. More generally, deep learning networks have also been described to learn temporal features. For instance, CNNs were used for the training and classification of 1 million YouTube videos [11]. In addition, the combination of learned features derived from a CNN and hand-crafted features was shown to be promising [39].

### 2.2 Structured temporal models in activity recognition

Because of an increased interest in continuous monitoring and analysis of human activities, several structured temporal models have been recently proposed for the recognition of complex activity sequences.

Early approaches were based on motion-captured data [9, 30, 15] or hand-labeled trajectories [23]. One of the first attempts to generalize these methods to raw video data was presented in [18]. In their approach, the authors proposed to classify human activities by aggregating information from motion segments based on visual features and temporal compositions. Video sequences were thus decomposed into temporal segments of variable length and matched against motion segment classifiers. The idea of representing a video by snippets was later adapted in diverse forms including Actoms [7], action spectograms [5], middle-level components [45] or clusters of Tracklets [8].

The development of approaches for the recognition of complex events has also gained in popularity. Early approaches modeled temporal structures using velocity history models [16], Bayes networks [27] or hybrid HMMs [5]. However, the datasets used for evaluating these approaches tended to be relatively simple. With the increased complexity of available datasets, the visual representations used by modern approaches has

also become increasingly complex. The temporal dynamics of video sequences was modeled in [2] using vector time series represented by the principal projections of an eigenvector decomposition of their block Hankel Matrix and harmonic signatures. The resulting mid-level representations were successfully applied to the recognition of complex events using the TRECVID dataset. In [6], the authors used a sequence memorizer [41], i.e., a hierarchical nonparametric Bayesian model that captures long-term dependencies in sequence data. A higher level representation based on a stochastic context-free grammar was proposed in [36]. Another approach for constructing an activity grammar automatically to capture hierarchical temporal structures was proposed in [22]. In this approach, parsing was based on a latent structural SVM which learns sub-actions automatically. A similarly unsupervised method for learning action units was described in [42]. Here, a causal topic model was used to learn the co-occurrence and temporal relation between action units in videos. Another system was described in [32] with the goal to detect missing actions within an activity sequence. The system combines coupled HMMs with a higher-level graph to model the overall structure of an activity. Based on the detection of omitted nodes/action units within the graph, the system produces notifications for missing action units.

## 2.3 Datasets

With the development of novel structured temporal models, various datasets have been described to evaluate the temporal parsing of activities. As there are many datasets for (clip-based) action classification (for an overview see e.g. [4]), we here only focus on video datasets that provide temporal segment annotations with at least one level of granularity. Some of the most recent activity recognition datasets that provide such annotations include the CMU MMAC dataset [33], the MPII Cooking dataset [24] and the 50 Salads dataset [34].

Most of these datasets are mid-size, they comprise only few subjects, and the number of clips per class tend to be small. Much like those available for action classification datasets such as HMDB [14] or UCF101 [31], which comprise 50–100 classes with 100 or more clips per class, there is a need for large-scale datasets for the recognition of complex everyday activity sequences. While this void is partially filled by the MPII Composite dataset [26], the dataset remains limited in that it was recorded in a lab environment with a fixed camera setup and constant lighting conditions. In contrast to the present Breakfast dataset which was recorded

"in the wild" using 18 real-world kitchens with uncontrolled camera positions and uncontrolled light conditions. Additionally, the behavior of staged participants in a lab setting may differ from their everyday behavior at home. Indeed, recent action classification benchmarks have demonstrated the importance of unconstrained settings: Most existing systems achieve near ceiling accuracy on staged datasets such as KTH [28] or Weizmann [3] but, in comparison, very few exhibit high accuracy on datasets collected in the wild such as HMDB or UCF101.

A second limitation of existing datasets is that the level of granularity of the action annotations vary from datasets to datasets. Most datasets include annotations based on human object interactions such as "open brownie box" (CMU MMAC) or "screw open" (MPII Cooking). But the labeling can also be based on more than one level granularity. For instance, in 50 Salads, where coarser elements such as "preparing salad" are made up from a set of finer action units such as "cut tomato" or "peel cucumber" and the finer actions are again partitioned into a preparation, core and post phase. Some datasets further provide labels for complex activity classes. Usually, one video represents one activity class. Examples for such long term activities can be found in the CMU MMAC or the MPII Cooking dataset for example.

Our dataset provides labels for long-term activities and segmentation at two levels of granularity. At the finest level, a segment is about 6 seconds long on average. At the coarser level, segment duration is about 26 seconds on average. See Table 1 for a complete comparison of existing datasets with the proposed Breakfast dataset.

## 3 Breakfast Dataset

### 3.1 Breakfast data collection

We describe the Breakfast dataset for the evaluation of structured temporal recognition models. The dataset features 52 unique participants, each engaged in 10 distinct cooking activities captured in 18 different kitchens[1]. Overall, the dataset includes about 200 clips for each cooking activity including the preparation of coffee (n = 200 samples), orange juice (n = 187), chocolate milk (n = 224), tea (n = 223), a bowl of cereals (n = 214), fried eggs (n = 198), pancakes (n = 173), a fruit salad (n = 185), a sandwich (n = 197) and scrambled eggs (n = 188).

---

[1] `http://serre-lab.clps.brown.edu/resource/breakfast-actions-dataset`

| | Activities | Units | Segments | Clips | Duration | Setting | Persons | Mean length |
|---|---|---|---|---|---|---|---|---|
| Cha LAP | - | 11 | 166 | 7 | 6min | staged | 8 | 2.1sec |
| YouCook | - | 7 | 1422 | 88(46) | 0.7h | youtube | - | 1.9sec |
| Toy assembly | 3 | 40 | 479 | 29 | 1.06h | lab | 2 | 6.4sec |
| CMU MMAC (*) | 1 | 37 | 2238 | 37 | 4.4h | lab | 16 | 6.8sec |
| 50 Salads | 2 | 51 | 2603 | 50 | 5.3h | lab | 25 | 7.0sec |
| MPII Cooking | 14 | 65 | 5609 (1861 BG) | 44 | 8.1h | lab | 12 | 8.0sec |
| MPII Composite | 55 | 78 | 12642 | 256 | 18.2h | lab | 30 | 5.1sec |
| Breakfast | 10 | 48 | 11441 (2992 BG) | 1712 | 66.7h | wild | 52 | 26.0sec |
| Breakfast Fine | 10 | 178 | 31325 (2154 BG) | 804 | 49.6h | wild | 52 | 5.6sec |

Table 1: Overview of existing datasets available for video segmentation evaluation and comparison with the proposed Breakfast dataset. Note that we only consider videos with action segment annotations, which, in some cases (e.g., YouCook or CMU MMAC) correspond to only a subset of the entire dataset. If background class labels are included in the annotations, the corresponding number of segments is included in parenthesis.

All activities were recorded with three to five cameras that were placed at various positions in the kitchens, so that the same activity is recorded from different, varying, viewpoints (Figure 1). For the recording, webcams, standard industry cameras (Prosilica GE680C) as well as a stereo camera (BumbleBee®, Pointgrey, Inc) were used. All videos were normalized to a resolution of $320 \times 240$ pixels with a frame rate of 15 fps. The video streams of each camera was manually synchronized. Overall the dataset provides about 66 hours of video and about 3.5 millions frames. For evaluation purpose, we organized the 52 participants in four groups, and permuted each of these four groups as splits for training and test.

The recording setup is "in the wild" as opposed to a single controlled lab environment [24, 33] in order to closely reflect real-world conditions as it pertains to the monitoring and analysis of daily activities. The actor performance was completely unscripted, unrehearsed and undirected. The actors were only handed a recipe and were instructed to prepare the corresponding food item. The resulting activities are thus highly variable both in terms of the choice of individual action units executed by the actors and their relative ordering. Since the sequences were recorded in various kitchens, the participants used the tools and packages that were locally available. Examples of the various settings and viewpoints are shown in Figure 1.

The set of cooking activities was chosen to include many similar elements (e.g., fried egg vs. scrambled egg preparation, or tea vs. coffee) resulting in shared action units (e.g., crack egg or take cup) across activities. This should yield a low inter-class variance for activities combined with a high intra-class variance because of different recoding locations, view-points and kitchens used. This challenging dataset thus allows for a thorough evaluation of structured temporal approaches.
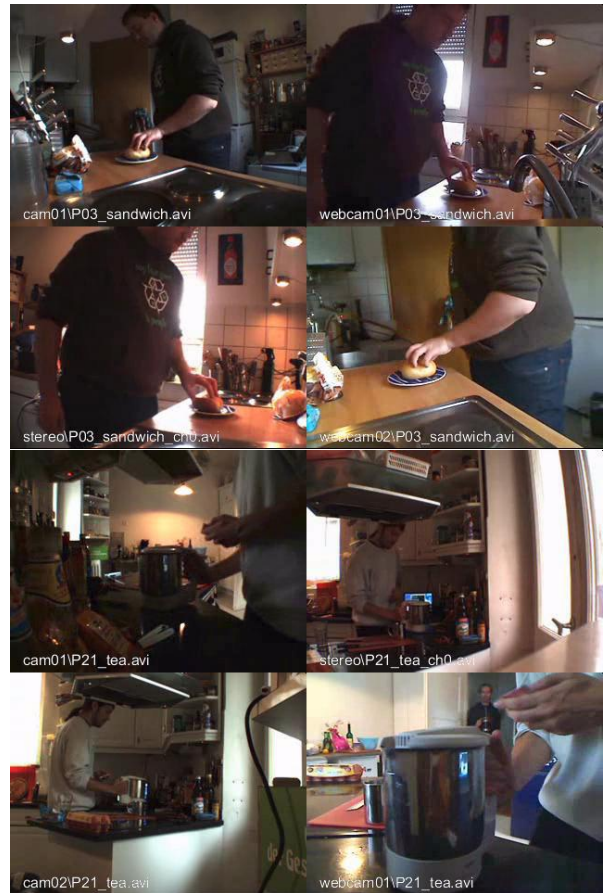


Fig. 1: Sample images from the Breakfast dataset

## 3.2 Data annotation

We asked two sets of annotators to manually label videos at two different levels of granularity: One group consisting of three annotators was asked to annotate action units at a coarse level (e.g., 'pour milk' or 'take plate'). The start and endpoints of a segment at the coarse level was typically based on the usage of a certain tool. For

| Coffee | take cup - pour coffee - pour milk - pour sugar - spoon sugar - stir coffee |
|---|---|
| (Chocolate) Milk | take cup - spoon powder - pour milk - stir milk |
| Juice | take squeezer - take glass - take plate - take knife - cut orange - squeeze orange - pour juice |
| Tea | take cup - add teabag - pour water - spoon sugar - pour sugar - stir tea |
| Cereals | take bowl - pour cereals - pour milk - stir cereals |
| Fried Egg | pour oil - butter pan - take egg - crack egg - fry egg - take plate - add salt and pepper - put egg onto plate |
| Pancakes | take bowl - crack egg - spoon flour - pour flour - pour milk - stir dough - pour oil - butter pan - pour dough into pan - fry pancake - take plate - put pancake onto plate |
| (Fruit) Salad | take plate - take knife - peel fruit - cut fruit - take bowl - put fruit to bowl - stir fruit |
| Sandwich | take plate - take knife - cut bun - take butter - smear butter - take topping - add topping - put bun together |
| Scrambled Egg | pour oil - butter pan - take bowl - crack egg - stir egg - pour egg into pan - stir fry egg - add salt and pepper - take plate - put egg onto plate |

Table 2: Coarse action units for individual activities.

instance, the coarse unit 'pour milk' starts when the milk package is reached and ends when the package is released again. It comprises all action units related to this task such as the opening or closing of the package and the pouring of the milk itself. Overall we identified 48 different coarse action units with about 11,000 samples in total including about 3,000 'silence' samples. Table 2 lists the coarse action units corresponding to individual activities.

To address the question of how granularity influences the overall activity recognition, we asked another group of fifteen annotators to provide annotations at a finer temporal scale. At the fine level, units usually correspond to body-part movements. For instance, a coarse unit such as 'pour milk' is decomposed into finer chunks such as 'grab milk' → 'twist cap' → 'open cap' etc. The mean length of the fine grained labels is 49 frames. As the task of low-level annotation is very time consuming, we only annotated a subset of 802 clips (about 28.4h of video) at the fine grained level. We refer to this subset as "Breakfast Fine". Overall, we collected about 31,000 fine grained units from 178 fine grained unit classes, including about 2,000 'silence' samples. Class labels are usually a composition of verb and object such as 'reach knife', 'pour water' or 'cut bread'. The compositions are made up of 38 unique verbs and 62 unique objects. Note that from the 178 classes about 100 are based on the verbs 'reach' or 'carry' such as 'reach milk' or 'reach spoon'.

Although the annotators for coarse and fine grained action units worked independently, we found a high correlation between breakpoints (corresponding to transitions between action units) derived from fine- and coarse-level annotations. Examples for the segmentation of coarse and fine grained units are shown in Figure 2. One can see that the boundaries of the coarse units are very close to boundaries of fine unit segments.

To assess systematic labeling errors, we computed the frame difference of each coarse breakpoint to the nearest fine grained breakpoint. The mean distance between fine and coarse breakpoints was 12 frames and around 60% of all coarse breakpoints were within 5 frames or less to their closest fine breakpoint. On average each coarse unit comprised 6 fine-grained action units.

## 4 Representation of Activities

### 4.1 Action unit model

To model the temporal extent of an ongoing movement, we present an approach borrowed from speech processing. By analogy to phonemes in speech that are building blocks for words or sentences, we interpret a complex activity sequence as a concatenation of shorter action units. Given this analogy, we model action units using HMMs much like phonemes in speech processing.

In order to model the temporal dynamics of action units, we assume that a video segment may be encoded as a sequence of feature vectors that represent the ongoing motion in each frame. The task of recognizing an action unit is therefore defined as that of finding the action unit $u_i \in \{u_1, u_2, \ldots, u_I\}$ that matches an input sequence $\mathbf{x} = (x_1, x_2, \ldots, x_T)$ best, with $x_t$ representing the feature vector at frame $t$. This can be formulated as maximizing the probability of an action unit $u_i$ given the input sequence $\mathbf{x}$:

$$\underset{i \in 1,\ldots,N}{\operatorname{argmax}} P(u_i|\mathbf{x}) = \underset{i \in 1,\ldots,N}{\operatorname{argmax}} \frac{P(\mathbf{x}|u_i)P(u_i)}{P(\mathbf{x})} \ . \qquad (1)$$

As the observation probability $P(\mathbf{x})$ of the current sequence $\mathbf{x}$ is the same for all units, it is usually omitted. The unit probability $P(u_i)$ is in our case proportional to $\frac{1}{N(u_i)}$ where $N(u_i)$ is the number of class samples in the training data.
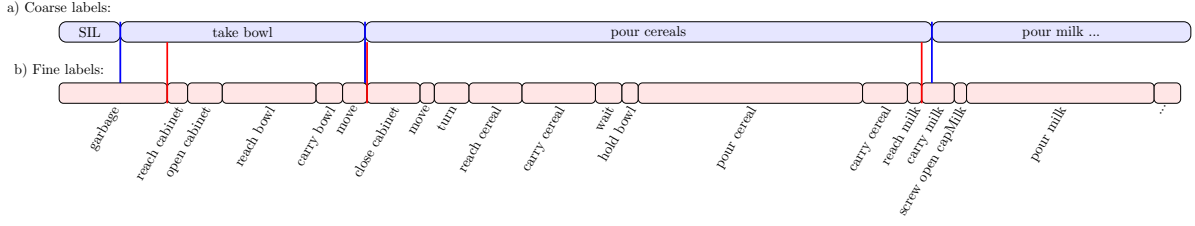
a) Coarse labels:



b) Fine labels:

Fig. 2: Example of coarse- and fine-level annotations for one video with the activity label "preapre cereals". The boundaries of the coarse units (blue) are very close to boundaries of fine unit segments (red), except for the special case of SIL/garbage.
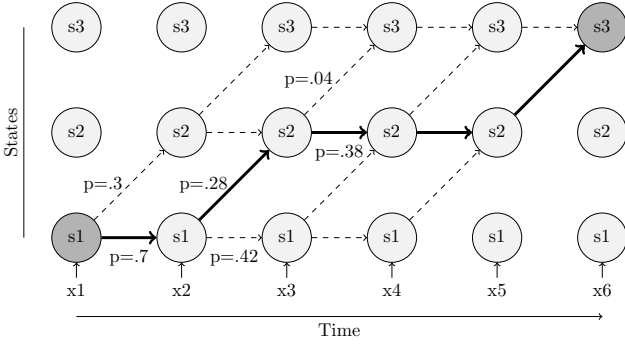


Fig. 3: Inference with a left-to-right feed forward topology. The set of states is given by $S = \{s_1, s_2, s_3\}$ and the input sequence is $\mathbf{x} = (x_1, x_2, \ldots, x_6)$, each $x_t$ corresponds to a feature vector sampled at frame $t$. The dashed lines show all possible solutions. Note that the sequence has to start with $s_1$ and end with $s_3$. Furthermore, only transitions to the next state are allowed. The path with the highest probability (bold) is obtained using the Viterbi algorithm.

To model $P(\mathbf{x}|u_i)$, we represent $u_i$ by a parametric Hidden Markov Model (HMM) $M_{u_i}$ defined by the set of states $S_{u_i} = \{s_1, s_2, s_3, \ldots, s_n\}$, the set of observations $X_{u_i} \subset \mathbb{R}^m$ with $m$ as the dimension of the input sequence, the state transition probability matrix $A_{u_i} \in \mathbb{R}^{n \times n}$ and the observation probability matrix $B_{u_i} \in \mathbb{R}^{n \times m}$. In our approach, the HMMs are defined by a strict left-to-right feed forward topology, thus, only self-transition and transitions to the next state are allowed as shown in Figure 3.

Sampling from a Markov model $M_{u_i}$, produces a sequence of states $\mathcal{S} = (\mathbf{S}(x_t))_{t=1,\ldots,T}$ with $\mathbf{S}(x_t) \in S_{u_i}$. The joint probability that the input sequence $\mathbf{x}$ and the sequence $\mathcal{S}$ generated by the Markov Model $M_{u_i}$ can be calculated as the product of transition probabilities $A_{u_i}$ and observation probabilities $B_{u_i}$:

$$P(\mathbf{x}, \mathcal{S}|M_{u_i}) = b_{s_n}(x_T) \prod_{t=1}^{T-1} a_{(\mathbf{S}(x_t), \mathbf{S}(x_{t+1}))} b_{(\mathbf{S}(x_t))}(x_t) ,$$

$$\tag{2}$$

where the transition probability from state $s_t$ to state $s_{t+1}$ is defined by $a_{(\mathbf{S}(x_t), \mathbf{S}(x_{t+1}))} \in A_{u_i}$ and the observation probability of a state $\mathbf{S}(x_t)$ is defined by a Gaussian mixture model:

$$b_s(x_t) = \sum_{k=1}^{K} \lambda_{ks} N(x_t; \mu_{ks}, \Sigma_{ks}) \tag{3}$$

with

$$N(x; \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^m |\Sigma|}} \exp\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)\right) ,$$

$$\tag{4}$$

where $m$ is the dimension of the input sequence $\mathbf{x}$, $\mu$ the $m$-dimensional mean vector, $\Sigma$ the $m \times m$ covariance matrix and $|\Sigma|$ the determinant of $\Sigma$.

We assume that $P(\mathbf{x}|M_{u_i})$ corresponds to $P(\mathbf{x}, \mathcal{S}|M_{u_i})$ by choosing the sequence $\hat{\mathcal{S}}$ that maximizes $P(\mathbf{x}, \mathcal{S}|M_{u_i})$, i.e.,

$$\hat{\mathcal{S}} = \underset{\mathcal{S}}{\operatorname{argmax}} \left( \prod_{t=1}^{T-1} a_{(\mathbf{S}(x_t), \mathbf{S}(x_{t+1}))} b_{(\mathbf{S}(x_t))}(x_t) \right) , \tag{5}$$

and the probability follows by

$$P(\mathbf{x}|M_{u_i}) = P(\mathbf{x}, \hat{\mathcal{S}}|M_{u_i}) . \tag{6}$$

This leads back to the idea that the model $M_{u_i}$ is a representation of the given unit $u_i$ and that the best path through $M_{u_i}$ corresponds to the probability $P(\mathbf{x}|u_i)$ of the observation of a unit $u_i$ given an input sequence $\mathbf{x}$

$$P(\mathbf{x}|u_i) = P(\mathbf{x}|M_{u_i}) . \tag{7}$$

The parameters $A$ and $B$ of the HMM are optimized using Baum-Welch re-estimation. For the decoding of HMMs, the Viterbi algorithm is used.
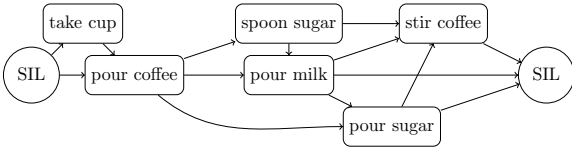
Fig. 4: Sample grammar used for the activity "prepare coffee". Each box represents an action unit (and thus an inidvidual HMM). "SIL" refers to the background (silence) class in our dataset, which is mandatory at the beginning and end of each sequence.

## 4.2 Sequence model

The recognition of individual action units may be thought of as a first step towards the analysis of ongoing event but it is unusual and rather artificial to assume that everyday tasks may consist in only one single action unit. Rather everyday activities consist in meaningful sequences of action units. Modeling activities as sequences of action units exhibit several advantages compared to treating it as a single entity. First, breaking down a complex activity into smaller action units allows not only for the recognition of the activity as a whole, but also for the parsing of the underlying sequence into action units. Second, the bottom-up construction of activities by composition of generic action units allows for a richer representation for efficiently learning novel activities composed of action units previously learned.

We use a grammar notation based on the extended Backus-Naur form to model activities as a combination of action units. The grammar is automatically generated from the segmentation transcripts of the training data. An example is given in Figure 4. The recognition of sequences is based on the token passing concept for connected speech recognition [43], augmenting the partial log probability with unit link records describing the transition from one unit to the next. To compute the most probable sequence, the Viterbi algorithm is used. At any frame $t$, the link records can be traced back to get the current most probable path, i.e., the most probable combination of units, and the position of the unit boundaries, i.e., the segmentation of the sequences until the current frame.

## 5 Evaluation

### 5.1 System description

As features, we use dense trajectories [37]. The dimensionality of the feature descriptors is first reduced from 426 dimensions to 64 dimensions by PCA, following the procedure described in [19]. To compute the Fisher
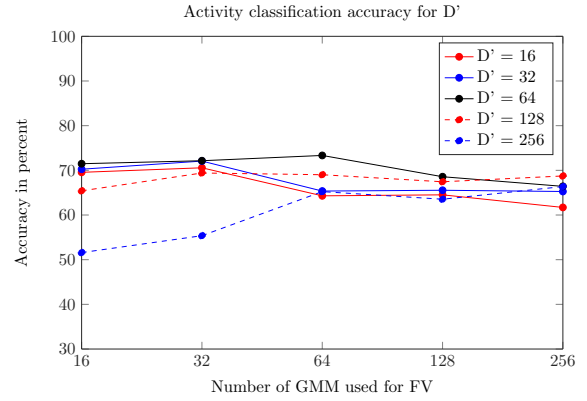


Fig. 6: Results for activity recognition using the first D' = [16, 32, 64, 128, 256] principal components of the FV representation after PCA.

Vectors (FVs), we sample 200,000 random features to learn Gaussian mixture models. The FV representation is computed for each frame over a sliding window of size 20 frames using *vlfeat* [35]. The dimensionality of the resulting vector is then reduced to 64 dimensions again using PCA. Thus, each frame is then represented by a 64-dimensional reduced FV. We further apply a L2-normalization to each feature dimension separately for each video clip.

In our implementation, we use the open source Hidden Markov Toolkit HTK [44]. For the training, inidvidual units are extracted and one HMM is trained for each unit. The number of states of the corresponding HMM is determined relative to the mean length $\hat{T}$ of the training samples, i.e. $n = \frac{\hat{T}}{10}$. We initialize the HMM by splitting all sequences evenly over time and assigning each sub-sequence to individual states. This initialization is possible as the HMMs are built in a left-to-right order. The HMM transition probabilities are initialized such that $a_{j,j} = 0.9$ and $a_{j,j+1} = 0.1$.

The Viterbi algorithm is applied to find the most likely state sequence for each training sequence and the HMM parameters are updated according to the newly estimated state sequence. The process is repeated until no further increase in likelihood is gained. After initialization, the states of individual action units are re-estimated by the Forward-Backward algorithm, optimizing the joint probability of states and frame inputs.

As generative models are prone to overfitting when given imbalanced training data, we set a lower bound of a 50 samples and an upper bound of 80 samples per action unit for training. When fewer training samples are available , we generate artificial samples by minority oversampling. Random down-selection is used when more than 80 samples are available.

$\mathcal{S} = \{SIL\_s_1, SIL\_s_1, SIL\_s_2, SIL\_s_2, SIL\_s_3, TakeCup\_s_1, TakeCup\_s_2, TakeCup\_s_3, TakeCup\_s_3, TakeCup\_s_3, ...\}$

$\mathcal{S} = \{SIL\_s_1, SIL\_s_3, SIL\_s_3, TakeCup\_s_1, TakeCup\_s_2, TakeCup\_s_2, TakeCup\_s_3, TakeCup\_s_3, PourCoffee\_s_1, PourCoffee\_s_1, ...\}$

$\mathcal{S} = \{SIL\_s_1, SIL\_s_1, SIL\_s_2, SIL\_s_2, SIL\_s_3, SIL\_s_3, PourCoffee\_s_1, PourCoffee\_s_2, PourCoffee\_s_2, PourCoffee\_s_2, ...\}$

$\mathcal{S} = \{SIL\_s_1, SIL\_s_2, SIL\_s_2, SIL\_s_3, PourCoffee\_s_1, PourCoffee\_s_2, PourCoffee\_s_2, PourCoffee\_s_2, PourCoffee\_s_3, PourCoffee\_s_3, ...\}$

$\mathcal{S} = \{SIL\_s_1, SIL\_s_2, SIL\_s_3, TakeCup\_s_1, TakeCup\_s_2, TakeCup\_s_3, PourCoffee\_s_1, PourCoffee\_s_2, PourCoffee\_s_2, PourCoffee\_s_2, ...\}$

$\mathcal{S} = \{SIL\_s_1, SIL\_s_1, SIL\_s_1, SIL\_s_2, SIL\_s_2, SIL\_s_2, SIL\_s_3, TakeCup\_s_1, TakeCup\_s_1, ...\}$

$\mathcal{S} = \{SIL\_s_1, SIL\_s_1, SIL\_s_1, SIL\_s_2, SIL\_s_2, SIL\_s_3, SIL\_s_3, SIL\_s_3, SIL\_s_3, PourCoffee\_s_1, ...\}$

$\mathcal{S} = \{SIL\_s_1, SIL\_s_2, SIL\_s_3, TakeCup\_s_1, TakeCup\_s_1, TakeCup\_s_1, TakeCup\_s_2, TakeCup\_s_3, TakeCup\_s_3, TakeCup\_s_3, ...\}$

$\mathcal{S} = \{SIL\_s_1, SIL\_s_2, SIL\_s_3, TakeCup\_s_1, TakeCup\_s_2, TakeCup\_s_3, PourCoffee\_s_1, PourCoffee\_s_1, PourCoffee\_s_2, PourCoffee\_s_2, ...\}$

$\mathcal{S} = \{...\}$

Fig. 5: The nine most likely paths after processing 10 frames of an input video. Inference corresponds to finding the path with the highest probability. This includes a path through individual action units and a path through the HMM of each selected action unit. For the example shown here, it is assumed for simplicity that each action unit is represented by an HMM with three states.

| Action unit classification - Breakfast - Coarse labels | | | | | | |
|---|---|---|---|---|---|---|
| | GMMs = | 16 | 32 | 64 | 128 | 256 |
| SVM w/o PCA | | 19.3 | 20.9 | 21.8 | 22.2 | 23.0 |
| SVM w PCA | $D' = 64$ | 15.2 | 16.1 | 16.6 | 16.6 | 17.8 |
| HMM w PCA | $D' = 64$ | 29.5 | 30.0 | 30.6 | 26.3 | 25.5 |

Table 3: Action unit classification on the Breakfast dataset based on coarse labels with 48 classes (1712 clips).

| Action unit classification - Breakfast Fine - Fine labels | | | | | | |
|---|---|---|---|---|---|---|
| | GMMs = | 16 | 32 | 64 | 128 | 256 |
| SVM w/o PCA | | 3.4 | 4.5 | - | - | - |
| SVM w PCA | $D' = 64$ | 3.0 | 3.2 | - | - | - |
| HMM w PCA | $D' = 64$ | 7.8 | 7.8 | 7.9 | 7.1 | 7.0 |

Table 4: Action unit classification for the Breakfast Fine dataset based on fine labels with 178 classes (804 clips).

## 5.2 Action unit classification

We first evaluate the performance of the model for the classification of individual action units, i.e., the classification of pre-segmented videos into 48 coarse or 178 fine scale action unit classes. The task is analogous to action clip-based action classification. Each video segment is classified independently using (1). We compare the classification accuracy of the HMMs with a linear SVM using the same feature representation with the exception that the FV representation is computed for the entire segment instead of each frame. In contrast to HMMs, which model the temporal relation of the observations, the SVM approach aggregates the observations independently of their temporal order.

Table 3 and 4 provide a comparison of the two approaches based on FV representations using GMMS with 16, 32, 64, 128 and 256 components. After PCA, the dimensionality of the FV representation is reduced to $D' = 64$. We found this value to work best on the proposed dataset (see Figure 6). HMMs work better with lower-dimensional features (64 components) while SVMs work better with higher-dimensional features.

| Activity classification - Breakfast | | | | | | |
|---|---|---|---|---|---|---|
| | GMMs = | 16 | 32 | 64 | 128 | 256 |
| SVM w/o PCA | | 52.0 | 52.6 | 48.7 | 39.6 | 23.2 |
| SVM w PCA | $D' = 64$ | 42.0 | 42.5 | 42.8 | 40.3 | 41.2 |
| Grammar | $D' = 64$ | **71.5** | **72.2** | **73.3** | **68.6** | **66.4** |

Table 5: Activity classification for the Breakfast dataset. For the grammar, action units are modeled using HMMs and PCA.

For coarse-level action unit classification, HMMs outperform linear SVMs. When comparing the best setting for HMM and SVM, the difference in accuracy is about 7%. On Breakfast Fine, HMMs with 64 components achieve an accuracy of 24.7% for the classification of coarse units (chance: 2.1%), thus a drop of about 5% compared to Breakfast. The drop can be explained by the reduced training set of Breakfast Fine compared to Breakfast cf. Table 1.

The results for fine unit classification on Breakfast Fine are reported in Table 4. Although the overall accuracy is lower for both approaches, the results are similar to the coarse action units, i.e., the HMMs outperform the linear SVMs.

## 5.3 Activity classification

We evaluate the approach for activity classification on both the Breakfast and Breakfast Fine dataset. For activity classification, a complete sequence needs to be classified into one of the 10 activity classes listed in Table 2. In our approach, the action units are modelled using HMMs (Section 4.1) and the activity sequence by a grammar (Section 4.2). For the SVM baseline, we encode the entire sequence with a single FV representation similar to the action unit classification task described above.

We report the activity recognition accuracy for both approaches on the Breakfast and Breakfast Fine dataset

| Activity classification - Breakfast Fine | | | | | | |
|---|---|---|---|---|---|---|
| | GMMs = | 16 | 32 | 64 | 128 | 256 |
| SVM w/o PCA | | 43.8 | 43.2 | 47.0 | 48.8 | 48.9 |
| SVM w PCA | $D' = 64$ | 37.8 | 32.4 | 36.1 | 34.0 | 41.0 |
| Grammar (coarse) | $D' = 64$ | **61.1** | **63.1** | **64.5** | **57.8** | **56.6** |
| Grammar (fine) | $D' = 64$ | **67.3** | **68.8** | **70.1** | **61.5** | **63.9** |

Table 6: Activity classification for the Breakfast Fine dataset. For the grammar (coarse), action units are defined using coarse labels. For the grammar (fine), action units are defined using fine labels.

in Table 5 and Table 6, respectively. On both datasets, the grammar with HMMs outperforms the SVM by about 20%. For the Breakfast Fine dataset, we also assess how the level of granularity used for training the system affect accuracy. The comparison shown in Table 6 shows that using 178 action units based on the fine labels instead of 48 action units based on the coarse labels improves the accuracy for activity classification. We will analyze the impact of the granularity further in Section 5.3. Compared to Breakfast, the overall recognition accuracy for the grammar based on coarse units decreases by about 10%. The same effect was observed above for action unit classification.

Figure 7 and 8 show the confusion matrices for both datasets. The grammar with HMMs, in contrast to the SVM, tends to confuse semantically-similar activities such as preparing coffee, tea or chocolate milk compared to semantically-dissimilar ones such as preparing a sandwich. Overall, the confusion matrix shows a clear grouping of activities related to the preparation of drinks (top left) vs. food (bottom right). The only exception is the preparation of cereals, which tends to get more confused with the preparation of drinks rather than food. When comparing the unit list in Table 2, however, one can see that the preparation of cereals shares more elements with the preparation of coffee, e.g., pouring and stirring, than with the preparation of scrambled egg or other food related activities.

### 5.3.1 Granularity in activity classification

In the previous experiments, we have observed that fine-grained action units result in a higher activity recognition rate compared to their coarse-grained counterparts. As can be seen from Table 1, while there are many more classes for the fine-level action units, their mean unit length tends to be shorter (5 seconds vs. 26 seconds for the coarse-level action units).

To assess the influence of the mean unit length on the overall activity recognition performance, we used the coarse labels of the Breakfast Fine dataset and split them evenly into 3, 5, 10, 15, and 20 parts resulting in a multiple for the original 48 classes, and thus reducing
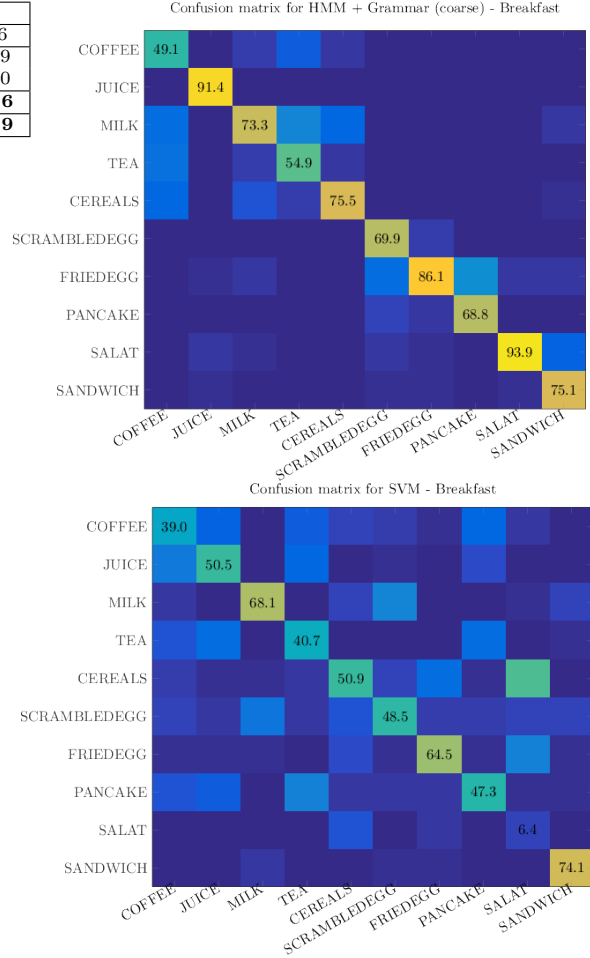


Fig. 7: Confusion matrix for the grammar (left) and the SVM (right) for Breakfast. The grammar gets mainly confused by semantically similar activities. For instance, the preparation of various drinks (coffee, milk, tea) are confused among each other. The confusion matrix for the SVM approach does not show a clear pattern.

the mean length of the units. For 5 splits, the number of action units is artificially increased to 240 and the mean length of each unit is reduced to 5.2 seconds. Table 7 compares the subdivided coarse units with the original coarse units and the fine units. The overall activity recognition increases with the level of granularity and almost reaches the performance of the recognition based on the fine-grained unit annotation. This implies that the approach works better with finer units.

Confusion matrix for HMM + Grammar (fine) - Breakfast Fine
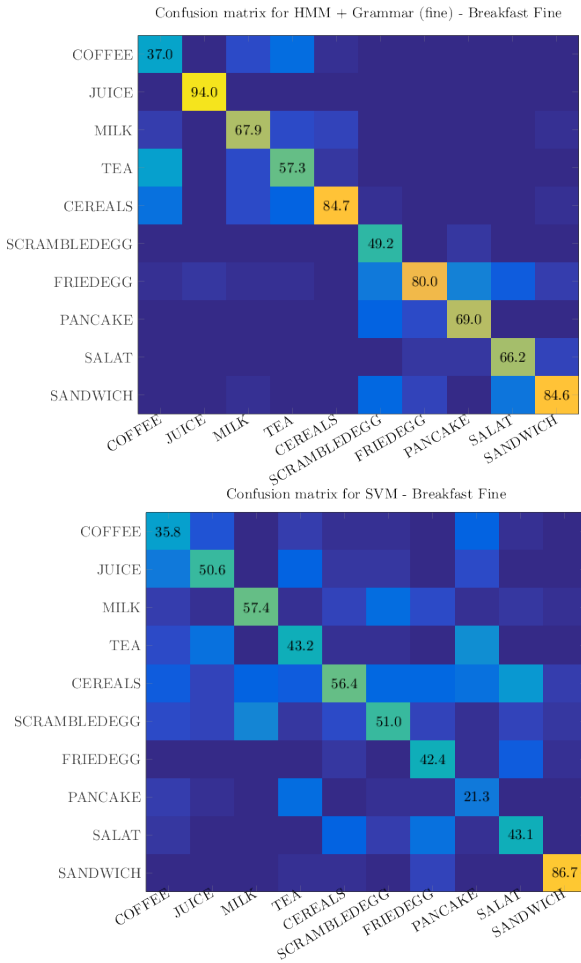


Confusion matrix for SVM - Breakfast Fine

Fig. 8: Confusion matrix for the grammar with fine action units (left) and the SVM (right) for Breakfast Fine with 804 clips. As in the large dataset HMM based recognition groups semantically similar activities.

| Granularity for activity recognition | | | | | | | |
|---|---|---|---|---|---|---|---|
| Splits: | 3 | 5 | 10 | 15 | 20 | coarse | fine |
| | 67.3 | 68.7 | 68.3 | 69.9 | 68.9 | 64.5 | 70.1 |

Table 7: Activity classification for Breakfast Fine. The coarse units are artificially split into smaller units.

## 5.4 Segmentation

The third task is the segmentation of long sequences, i.e., the detection of action units as they appear in an unknown sequence including the start and end frames of each segment corresponding to one unit. The proposed approach based on a grammar and HMMs directly predicts the action unit and the state of an action unit for each frame as illustrated in Figure 4. We evaluate the segmentation for both sets by looking at the over-

| Segmentation - Breakfast | | | | | |
|---|---|---|---|---|---|
| GMMs | 16 | 32 | 64 | 128 | 256 |
| Grammar (MoC) | **36.2** | **36.9** | **38.1** | **34.0** | **32.7** |
| HMMs (MoC) | 18.7 | 19.2 | 19.8 | 16.5 | 15.9 |
| Grammar (MoF) | **54.2** | **54.4** | **56.3** | **51.9** | **50.7** |
| HMMs (MoF) | 24.2 | 24.9 | 26.5 | 20.8 | 20.5 |

Table 8: Segmentation results for Breakfast with 48 action units. For HMMs, the grammar is replaced by a transition graph that allows transitions to and from any action unit. MoC denotes mean over class, MoF denotes mean over frames.

| Segmentation - Breakfast Fine | | | | | |
|---|---|---|---|---|---|
| GMMs | 16 | 32 | 64 | 128 | 256 |
| Grammar (MoC) | **11.2** | **11.6** | **12.2** | **10.4** | **10.5** |
| Grammar (MoF) | **28.7** | **28.6** | **31.3** | **24.2** | **26.4** |

Table 9: Segmentation results for Breakfast Fine with 178 action units. MoC denotes mean over class, MoF denotes mean over frames.

all number of frames that were correctly classified in terms of action units. Following the evaluation protocol of [13], we report mean over class (MoC) and mean over frames (MoF).

In Table 8 and 9, we report the segmentation accuracy for Breakfast and Breakfast Fine, respectively. On Breakfast, the coarse units are very well segmented with 56.3% of all frames correctly classified. The evaluation includes the sequences that were wrongly classified, which are about 25% of the sequences, cf. Table 5. When only the correctly classified sequences are considered, the overall amount of correctly classified frames increases to 70.5%. To asses the segmentation performance without grammar, we replace the grammar by a transition graph that allows transitions to and from any action unit without constraints. Without the grammar, the segmentation accuracy drops to 26.5% correctly classified frames.

We also evaluate the segmentation accuracy for the Breakfast Fine dataset using the fine-grained action units. While 31.3% of the frames are correctly classified, the mean over class is only 12.2%. This shows the difficulty of segmenting fine-grained action units.

The problem becomes visible when looking at the related examples for the fine segmentation in Figure 7. Even if the overall sequence is correctly recognized, the alignment, especially in case of short units, can be off and so, decreasing the overall frame segmentation performance.
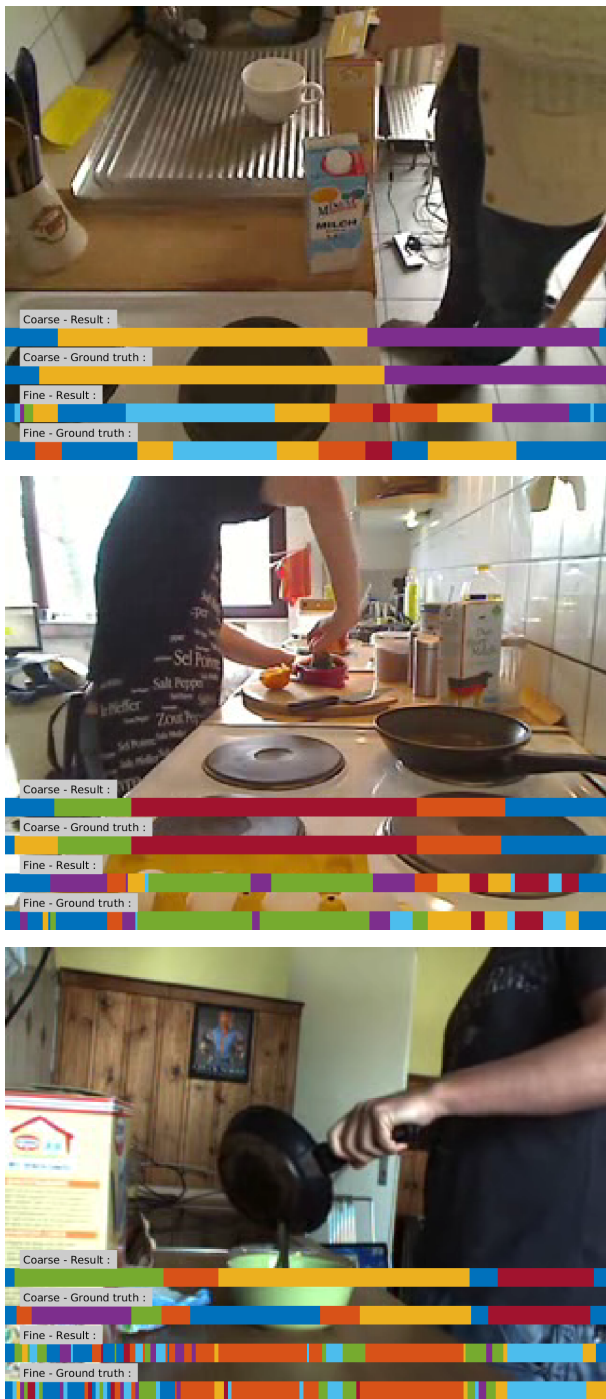
Fig. 9: Examples of coarse and fine segmentation results for Breakfast. The upper two bar shows the recognized sequence and the respective ground truth for the coarse annotations and the lower colorbar for the fine annotation. Although the fine grained units are usually correctly classified, the overall alignment error of fine-grained units leads to a lower frame classification accuracy as in case of the coarse units.

## 5.5 Runtime analysis

We report the runtime for activity classification and segmentation on Breakfast Fine over all four splits without the computation of the FV representation for each frame. The grammar with 48 coarse action units requires 9.1 hours for training and 1.3 hours for testing. For 178 fine-grained action units, the training reduces to 0.86 hours since each HMM consists of less states and is trained on shorter sequences. The inference time, however, increases to 3.85 hours since the grammar comprises more and longer valid paths.

## 5.6 Impact of training data

We have observed that activity recognition with coarse action units is lower on Breakfast Fine than on Breakfast due to the smaller amount of the training data. As the amount of training data available plays an important role in the context of generative models, we asses how the amount of training data influences the overall recognition accuracy. To provide a more in-depth analysis of the impact of the amount of training data, we reduce the input training data for each HMM to 50, 25, 10, 5 and 3 samples per coarse action unit on Breakfast Fine.

Figure 10 plots the activity classification accuracy and the segmentation accuracy measured as mean over frames or mean over classes. The activity classification accuracy drops to 40% when reducing the number of samples from 50 to 25 samples. If the number of samples per unit is less equal to 10, the activity recognition is not anymore better than chance. For the segmentation task, the same observation can be made. This shows the strong influence of the amount of training data for the approach and thus the need for large datasets like Breakfast and Breakfast Fine to study such approaches.

## 5.7 Bootstrapping

We have observed that a large amount of training data is needed to successfully train a generative model for activity classification and segmentation. The acquisition and annotation of training data, however, is very time consuming. An alternative is semi-supervised learning where only a subset is fully annotated at the frame level and most of the data is only weakly annotated with the order of action units as they appear in a video but without any alignment with the frames. The weak annotation can be considered as a transcript of the video.

For semi-supervised learning, we use bootstrapping, i.e. we first train the approach on the fully annotated
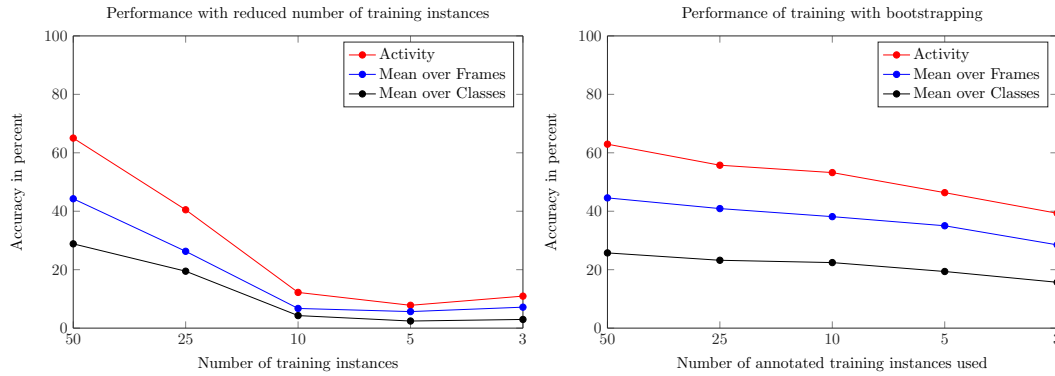
Fig. 10: Results for activity recognition and segmentation with reduced training data without and with additional bootstrapping

training data as initialization and use the transcripts of the rest of the training data to reestimate the model parameters. After the model is trained on the annotated data, we infer the segmentation on the rest of the training data where the path in our grammar is given by the transcript of each video. The model parameters are then reestimated on the entire training set.

For evaluation, we reduce the amount of fully annotated training data to 50, 25, 10, 5 and 3 samples per coarse action unit class on Breakfast Fine. Figure 10 plots the activity classification accuracy and the segmentation accuracy measured as mean over frames or mean over classes. The plot shows that even with a much smaller amount of fully annotated data the activity classification accuracy decreases gently to about 40%. The same holds for the segmentation accuracy. Even though only 3 samples were used for initialization, about 30% of all frames are correctly classified after model reestimation based on the transcribed data.

## 6 Evaluation on other datasets

We assess the performance of the proposed generative model also on four other available datasets. We first discuss its application for the action detection task on the Cha LAP dataset as an example how data augmentation in combination with majority voting can help to apply the approach even to datasets with few samples. Second, we evaluate the segmentation performance on a broad list of available datasets.

### 6.1 Cha LAP dataset

The action detection task for the Cha learning challenge differs from the so far discussed scenario since the videos contain multiple activities at the same time, i.e., the

| ChaLAP - Jaccard | | | |
|---|---|---|---|
| GMMs | 16 | 32 | 64 | 128 |
| $D' = 16$ | 0.456 | 0.492 | 0.436 | 0.339 |
| $D' = 32$ | 0.475 | 0.499 | 0.450 | 0.361 |
| $D' = 64$ | 0.416 | 0.474 | 0.418 | 0.331 |

Table 10: Jaccard index for Cha LAP

| ChaLAP - Jaccard - Benchmarks | |
|---|---|
| Suh [29] | 0.4226 |
| Pei [20] | 0.5011 |
| Peng [21] | 0.5071 |
| Wang[40] | 0.5385 |
| proposed | 0.5239 |

Table 11: Jaccard index for Cha LAP

activity of each actor needs to be inferred. We therefore detect and track each actor to acquire the actor specific bounding boxes. The related dense trajectories are then sampled for each tracked bounding box. The dataset is very small and contains only 7 sequences for training. The segmentation quality is measured by the Jaccard index.

In Table 10, we report the segmentation accuracy for 16, 32, 64 and 128 GMM components and reduced feature dimensionality of 16, 32 and 64. The Jaccard index ranges form 0.33 to 0.5 for the different parameter combinations. Due to the small size of the dataset, it is difficult to determine the optimal parameters. We therefore propose two approaches to compensate for the lack of training data. First, we augment the training data by mirroring the videos. Averaged over all parameter settings, the accuracy increases from 0.429 to 0.452. Second, we segment the tracked bounding boxes with all parameter settings and assign a label for each frame by majority vote.

As shown in Table 10, we have 12 different parameter combinations resulting in 12 segmentation hypotheses of the original video. As for the training data, we

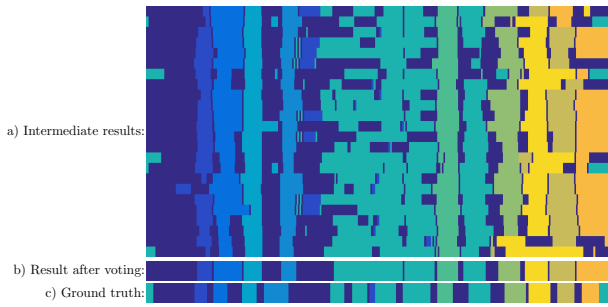a) Intermediate results:

b) Result after voting:

c) Ground truth:

Fig. 11: Example of a segmentation result for one actor on Cha LAP. a) shows the 24 segmentation hypotheses, b) shows the segmentation based on majority voting and c) shows the ground truth annotation.

also mirror the test video, which results in additional 12 segmentation hypotheses. An example for the 24 segmentation hypotheses for one actor is shown in Figure 11. The final segmentation is obtained by majority voting, i.e. selecting for each frame the class label that occurred the most often in all hypotheses.

Due to the voting, we do not have to choose a particular parameter setting and achieve a Jaccard index of 0.5239. As shown in Table 11, our voting procedure reaches state-of-the-art accuracy for this task despite of the small amount of training data.

### 6.2 Other datasets

We also evaluate the parsing and segmentation performance of the proposed generative recognition approach on other available complex activity datasets that are labeled at one or more levels of granularity as listed in Table 1. The datasets used for this evaluation are Toy assembly [36], CMU MMAC [33], MPII Cooking [25] and 50 Salads [34]. Sample frames for each of these datasets are shown in Figure 12. Since we observed that the accuracy of the proposed method strongly depends on the amount of training data, we report the number of training samples per class in Table 12. Depending on the benchmark, different measures have been proposed. For all datasets, we report the segmentation accuracy as mean over class (MoC). In addition, we report the accuracy according to the measure that was originally proposed by the corresponding dataset.

We compare our approach to the best reported results on each benchmark in Table 13. It shows that the proposed approach underperforms the best segmentation results in case of the small Toy assembly dataset. For the other datasets, which contain at least 4 hours of video, our approach significantly outperforms the state-of-the-art in terms of segmentation accuracy.

|  | Train samples used per class |
|---|---|
| Toy assembly | 15-20 samples |
| CMU MMAC | 30-40 samples |
| MPII Cooking | 12-30 samples |
| 50 Salads | 30-35 samples |

Table 12: Number of used training samples per action unit

## 7 Conclusion

We proposed a challenging large-scale dataset for human activity recognition in the wild as well as a generative approach for activity recognition and segmentation. The dataset offers the opportunity to evaluate the performance of approaches for activity classification as well as segmentation on realistic, large-scale data. The labeling at two levels of temporal granularity further allows to investigate the impact of granularity. On this dataset, we thoroughly evaluated the proposed generative approach that models activities by a grammar and action units as Hidden Markov Models in combination with a compact video representation based on Fisher Vectors. The experimental evaluation showed that the approach outperforms the state-of-the-art and that generative approaches can provide high quality activity classification and segmentation results, but they need sufficient training data. While we also discussed approaches to overcome the lack of annotated training data, the two Breakfast datasets will allow to study temporally structured models more in detail.

## References

1. Baird, J., Baldwin, D.: Making Sense of Human Behavior: Action Parsing and Intentional Inference. In: Intentions and Intentionality. MIT Press (2001)
2. Bhattacharya, S., Kalayeh, M., Sukthankar, R., Sha, M.: Recognition of Complex Events: Exploiting Temporal Dynamics between Underlying Concepts. In: CVPR (2014)
3. Blank, M., Gorelick, L., Shechtman, E., Irani, M., Basri, R.: Actions as space-time shapes. In: ICCV (2005)
4. Chaquet, J.M., Carmona, E.J., Fernndez-Caballero, A.: A survey of video datasets for human action and activity recognition. CVIU **117**(6), 633 – 659 (2013)
5. Chen, C., Aggarwal, J.: Modeling human activities as speech. In: CVPR (2011)
6. Cheng, Y., Fan, Q., Pankanti, S., Choudhary, A.: Temporal Sequence Modeling For Video Event Detection . In: CVPR (2014)
7. Gaidon, A., Harchaoui, Z., Schmid, C.: Actom Sequence Models for Efficient Action Detection. In: CVPR, pp. 3201–3208 (2011)
8. Gaidon, A., Harchaoui, Z., Schmid, C.: Activity representation with motion hierarchies. IJCV **107**(3), 219–238 (2014)
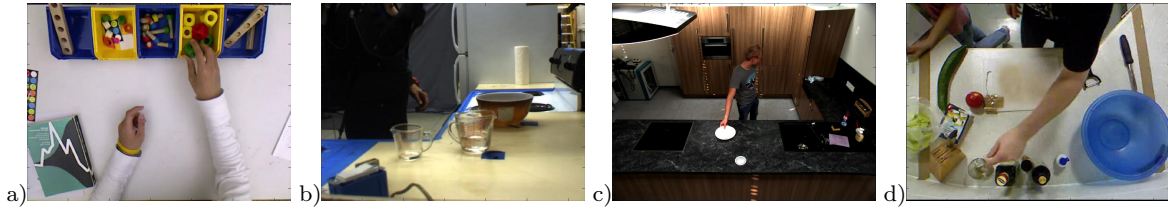
Fig. 12: Sample frames from the datasets used for performance evaluation: a) Toy assembly [36], b) CMU MMAC [33], c) MPII Cooking [24], d) 50 Salads [34].

| | | Segmentation | | | |
|---|---|---|---|---|---|
| GMM= | Toy assembly | CMU MMAC | MPII Cooking | 50 Salads | Breakfast |
| 16 | 50.3 / *64.3* | 53.8 / *60.8* | 46.5 / *58.5* | 81.6 | 36.2 / *54.2* |
| 32 | 48.6 / *63.1* | 53.7 / *60.7* | 53.9 / *68.5* | 80.4 | 36.9 / *54.4* |
| 64 | 56.7 / *67.5* | 53.0 / *60.3* | 51.6 / *63.9* | **83.8** | **38.1** / ***56.3*** |
| 128 | 60.5 / *70.8* | 52.5 / *60.4* | 53.9 / *66.8* | 82.0 | 34.0 / *51.2* |
| 256 | 63.5 / *72.2* | **58.8 / 67.1** | **57.3 / 71.7** | **83.8** | 32.7 / *50.7* |
| Best | – / ***91.0*** [36] | – / *59.0* [36] | – / *54.3* [17] | 67.6 [34] | – / *28.8* [12] |

Table 13: Overview of the segmentation results for all datasets. Accuracy is computed as the mean over all classes. For comparison, we also report the frame-based accuracy (*italic*) for the Toy assembly and CMU MMAC and midpoint hit accuracy (*also italic*) for the MPII Cooking dataset as used by the authors in the original studies.

9. Guerra-Filho, G., Fermüller, C., Aloimonos, Y.: Discovering a language for human activity. In: Proc. of the AAAI Symposium on Anticipatory Cognitive Embodied Systems (2005)
10. Jain, M., van Gemert, J.C., Snoek, C.G.: What do 15,000 object categories tell us about classifying and localizing actions? In: CVPR, pp. 46–55 (2015)
11. Karpathy, A., Toderici, G., Shetty, S., Leung, T., Sukthankar, R., Fei-Fei, L.: Large-scale video classification with convolutional neural networks. In: CVPR, pp. 1725–1732 (2014)
12. Kuehne, H., Arslan, A., Serre, T.: The Language of Actions: Recovering the Syntax and Semantics of Goal-Directed Human Activities. In: CVPR (2014)
13. Kuehne, H., Gall, J., Serre, T.: An end-to-end generative framework for video segmentation and recognition. In: WACV (2016)
14. Kuehne, H., Jhuang, H., Garrote, E., Poggio, T., Serre, T.: HMDB: a large video database for human motion recognition. In: ICCV (2011)
15. Lillo, I., Soto, A., Niebles, J.: Discriminative Hierarchical Modeling of Spatio-Temporally Composable Human Activities. In: CVPR (2014)
16. Messing, R., Pal, C., Kautz, H.: Activity recognition using the velocity histories of tracked keypoints. In: ICCV, pp. 104–111 (2009)
17. Ni, B., Paramathayalan, V., Moulin, P.: Multiple granularity analysis for fine-grained action detection. In: CVPR (2014)
18. Niebles, J., Chen, C.W., Fei-Fei, L.: Modeling temporal structure of decomposable motion segments for activity classification. In: ECCV, pp. 392–405. Springer Berlin Heidelberg (2010)
19. Oneata, D., Verbeek, J., Schmid, C.: Action and Event Recognition with Fisher Vectors on a Compact Feature Set. In: ICCV, pp. 1817–1824. IEEE (2013)
20. Pei, Y., Ni, B., Atmosukarto, I.: ECCV, chap. Mixture of Heterogeneous Attribute Analyzers for Human Action Detection. Springer International Publishing (2015)
21. Peng, X., Wang, L., Cai, Z., Qiao, Y.: ECCV, chap. Action and Gesture Temporal Spotting with Super Vector Representation. Springer International Publishing (2015)
22. Pirsiavash, H., Ramanan, D.: Parsing videos of actions with segmental grammars. In: CVPR (2014)
23. Rao, C., Yilmaz, A., Shah, M.: View-Invariant Representation and Recognition of Actions. IJCV **50**(2), 203–226 (2002)
24. Rohrbach, M.: MPII Cooking Activities Dataset (2013). URL http://www.d2.mpi-inf.mpg.de/mpii-cooking
25. Rohrbach, M., Amin, S., Andriluka, M., Schiele, B.: A Database for Fine Grained Activity Detection of Cooking Activities. In: CVPR (2012)
26. Rohrbach, M., Regneri, M., Andriluka, M., Amin, S., Pinkal, M., Schiele, B.: ECCV, chap. Script Data for Attribute-Based Recognition of Composite Activities. Springer Berlin Heidelberg (2012)
27. Ryoo, M.S., Aggarwal, J.K.: Semantic Representation and Recognition of Continued and Recursive Human Activities. IJCV **82**(1), 1–24 (2009)
28. Schuldt, C., Laptev, I., Caputo, B.: Recognizing human actions: A local svm approach. In: ICPR (2004)
29. Shu, Z., Yun, K., Samaras, D.: ECCV, chap. Action Detection with Improved Dense Trajectories and Sliding Window. Springer International Publishing (2015)
30. Sminchisescu, C., Kanaujia, A., Li, Z., Metaxas, D.: Conditional models for contextual human motion recognition. In: ICCV, vol. 2, pp. 1808–1815 (2005)
31. Soomro, K., Zamir, A.R., Shah, M.: UCF101: A dataset of 101 human actions classes from videos in the wild. CoRR **abs/1212.0402** (2012). URL http://arxiv.org/abs/1212.0402
32. Soran, B., Farhadi, A., Shapiro, L.: Generating notifications for missing actions: Don't forget to turn the lights off! In: ICCV (2015)
33. Spriggs, E.H., De la Torre, F., Hebert, M.: Temporal Segmentation and Activity Classification from First-person Sensing. In: IEEE Workshop on Egocentric Vision at CVPR (2009)

34. Stein, S., McKenna, S.J.: Combining embedded accelerometers with computer vision for recognizing food preparation activities. In: UbiComp. ACM (2013)
35. Vedaldi, A., Fulkerson, B.: VLFeat: An open and portable library of computer vision algorithms. `http://www.vlfeat.org/` (2008)
36. Vo, N., Bobick, A.: From Stochastic Grammar to Bayes Network: Probabilistic Parsing of Complex Activity. In: CVPR (2014)
37. Wang, H., Klser, A., Schmid, C., Liu, C.L.: Dense trajectories and motion boundary descriptors for action recognition. IJCV **103**(1), 60–79 (2013)
38. Wang, H., Schmid, C.: Action Recognition with Improved Trajectories. In: ICCV, pp. 3551–3558 (2013)
39. Wang, L., Qiao, Y., Tang, X.: Action recognition with trajectory-pooled deep-convolutional descriptors. In: CVPR (2015)
40. Wang, Z., Wang, L., Du, W., Qiao, Y.: Exploring fisher vector and deep networks for action spotting. In: IEEE Conf. on Computer Vision and Pattern Recognition Workshops, (2015)
41. Wood, F., Archambeau, C., Gasthaus, J., James, L., Teh, Y.W.: A stochastic memoizer for sequence data. In: Proceedings of the 26th Annual International Conference on Machine Learning, ICML. ACM (2009)
42. Wu, C., Zhang, J., Savarese, S., Saxena, A.: Watch-n-patch: Unsupervised understanding of actions and relations. In: CVPR (2015)
43. Young, S., Russell, N., Thornton, J.: Token passing: a simple conceptual model for connected speech recognition systems. Tech. rep., Cambridge University (1989)
44. Young, S.J., Evermann, G., Gales, M.J.F., Hain, T., Kershaw, D., Moore, G., Odell, J., Ollason, D., Povey, D., Valtchev, V., Woodland, P.C.: The HTK Book, version 3.4. Cambridge University Engineering Department (2006)
45. Yuan, F., Prinet, V., Yuan, J.: Middle-level representation for human activities recognition: The role of spatiotemporal relationships. In: Trends and Topics in Computer Vision, *Lecture Notes in Computer Science*, vol. 6553, pp. 168–180. Springer Berlin Heidelberg (2012)
46. Zacks, J.M., Kumar, S., Abrams, R.A., Mehta, R.: Using movement and intentions to understand human activity. Cognition, International Journal of Cognitive Science **112**(2), 201 – 216 (2009)
47. Zacks, J.M., Speer, N.K., Swallow, K.M., Braver, T.S., Reynolds, J.R.: Event perception: a mind-brain perspective. Psychological bulletin **133**(2), 273 – 293 (2007)