

The Kinetics Human Action Video Dataset

Will Kay

wkay@google.com

João Carreira

joaoluis@google.com

Karen Simonyan

simonyan@google.com

Brian Zhang

brianzhang@google.com

Chloe Hillier

chillier@google.com

Sudheendra Vijayanarasimhan

svnaras@google.com

Fabio Viola

fviola@google.com

Tim Green

tfgg@google.com

Trevor Back

back@google.com

Paul Natsev

natsev@google.com

Mustafa Suleyman

mustafasul@google.com

Andrew Zisserman

zisserman@google.com

Abstract

We describe the DeepMind Kinetics human action video dataset. The dataset contains 400 human action classes, with at least 400 video clips for each action. Each clip lasts around 10s and is taken from a different YouTube video. The actions are human focussed and cover a broad range of classes including human-object interactions such as playing instruments, as well as human-human interactions such as shaking hands. We describe the statistics of the dataset, how it was collected, and give some baseline performance figures for neural network architectures trained and tested for human action classification on this dataset. We also carry out a preliminary analysis of whether imbalance in the dataset leads to bias in the classifiers.

1. Introduction

In this paper we introduce a new, large, video dataset for human action classification. We developed this dataset principally because there is a lack of such datasets for human action classification, and we believe that having one will facilitate research in this area – both because the dataset is large enough to train deep networks from scratch, and also because the dataset is challenging enough to act as a performance benchmark where the advantages of different architectures can be teased apart.

Our aim is to provide a large scale high quality dataset, covering a diverse range of human actions, that can be used for human action *classification*, rather than temporal localization. Since the use case is classification, only short clips of around 10s containing the action are included, and there are no untrimmed videos. However, the clips also contain sound so the dataset can potentially be used for many

purposes, including multi-modal analysis. Our inspiration in providing a dataset for classification is ImageNet [18], where the significant benefits of first training deep networks on this dataset for classification, and then using the trained network for other purposes (detection, image segmentation, non-visual modalities (e.g. sound, depth), etc) are well known.

The Kinetics dataset can be seen as the successor to the two human action video datasets that have emerged as the standard benchmarks for this area: HMDB-51 [15] and UCF-101 [20]. These datasets have served the community very well, but their usefulness is now expiring. This is because they are simply not large enough or have sufficient variation to train and test the current generation of human action classification models based on deep learning. Coincidentally, one of the motivations for introducing the HMDB dataset was that the then current generation of action datasets was too small. The increase then was from 10 to 51 classes, and we in turn increase this to 400 classes.

Table 1 compares the size of Kinetics to a number of recent human action datasets. In terms of variation, although the UCF-101 dataset contains 101 actions with 100+ clips for each action, all the clips are taken from only 2.5k distinct videos. For example there are 7 clips from one video of the same person brushing their hair. This means that there is far less variation than if the action in each clip was performed by a different person (and different viewpoint, lighting, etc). This problem is avoided in Kinetics as each clip is taken from a different video.

The clips are sourced from YouTube videos. Consequently, for the most part, they are not professionally videoed and edited material (as in TV and film videos). There can be considerable camera motion/shake, illumination variations, shadows, background clutter, etc. More im-

Dataset	Year	Actions	Clips	Total	Videos
HMDB-51 [15]	2011	51	min 102	6,766	3,312
UCF-101 [20]	2012	101	min 101	13,320	2,500
ActivityNet-200 [3]	2015	200	avg 141	28,108	19,994
Kinetics	2017	400	min 400	306,245	306,245

Table 1: Statistics for recent human action recognition datasets. ‘Actions’, specifies the number of action classes; ‘Clips’, the number of clips per class; ‘Total’, is the total number of clips; and ‘Videos’, the total number of videos from which these clips are extracted.

portantly, there are a great variety of performers (since each clip is from a different video) with differences in *how* the action is performed (e.g. its speed), clothing, body pose and shape, age, and camera framing and viewpoint.

Our hope is that the dataset will enable a new generation of neural network architectures to be developed for video. For example, architectures including multiple streams of information (RGB/appearance, optical flow, human pose, object category recognition), architectures using attention, etc. That will enable the virtues (or otherwise) of the new architectures to be demonstrated. Issues such as the tension between static and motion prediction, and the open question of the best method of temporal aggregation in video (recurrent vs convolutional) may finally be resolved.

The rest of the paper is organized as: Section 2 gives an overview of the new dataset; Section 3 describes how it was collected and discusses possible imbalances in the data and their consequences for classifier bias. Section 4 gives the performance of a number of ConvNet architectures that are trained and tested on the dataset. Our companion paper [5] explores the benefit of pre-training an action classification network on Kinetics, and then using the features from the network for action classification on other (smaller) datasets.

The URLs of the YouTube videos and temporal intervals of the dataset can be obtained from <http://deepmind.com/kinetics>.

2. An Overview of the Kinetics Dataset

Content: The dataset is focused on human actions (rather than activities or events). The list of action classes covers: *Person Actions (singular)*, e.g. drawing, drinking, laughing, pumping fist; *Person-Person Actions*, e.g. hugging, kissing, shaking hands; and, *Person-Object Actions*, e.g. opening present, mowing lawn, washing dishes. Some actions are fine grained and require temporal reasoning to distinguish, for example different types of swimming. Other actions require more emphasis on the object to distinguish, for example playing different types of wind instruments.

There is not a deep hierarchy, but instead there are several (non-exclusive) parent-child groupings, e.g. Music (playing drums, trombone, violin, ...); Personal Hygiene (brushing teeth, cutting nails, washing hands, ...); Dancing

(ballet, macarena, tap, ...); Cooking (cutting, frying, peeling, ...). The full list of classes is given in the appendix, together with parent-child groupings. Figure 1 shows clips from a sample of classes.

Statistics: The dataset has 400 human action classes, with 400–1150 clips for each action, each from a unique video. Each clip lasts around 10s. The current version has 306,245 videos, and is divided into three splits, one for training having 250–1000 videos per class, one for validation with 50 videos per class and one for testing with 100 videos per class. The statistics are given in table 2. The clips are from YouTube videos and have a variable resolution and frame rate.

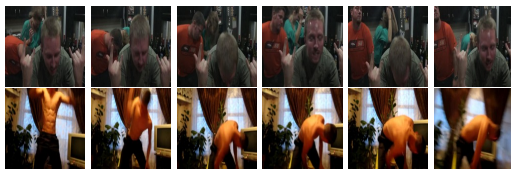
Train	Validation	Test
250–1000	50	100

Table 2: Kinetics Dataset Statistics. The number of clips for each class in the train/val/test partitions.

Non-exhaustive annotation. Each class contains clips illustrating that action. However, a particular clip can contain several actions. Interesting examples in the dataset include: “texting” while “driving a car”; “Hula hooping” while “playing ukulele”; “brushing teeth” while “dancing” (of some type). In each case both of the actions are Kinetics classes, and the clip will probably only appear under only one of these classes not both, i.e. clips do not have complete (exhaustive) annotation. For this reason when evaluating classification performance, a top-5 measure is more suitable than top-1. This is similar to the situation in ImageNet [18], where one of the reasons for using a top-5 measure is that images are only labelled for a single class, although it may contain multiple classes.

3. How the Dataset was Built

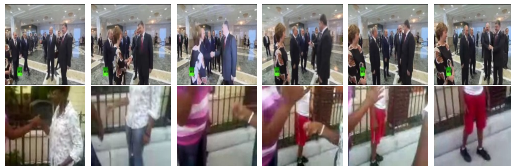
In this section we describe the collection process: how candidate videos were obtained from YouTube, and then the processing pipeline that was used to select the candidates



(a) headbanging



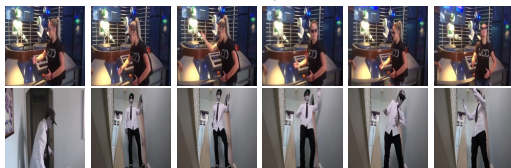
(b) stretching leg



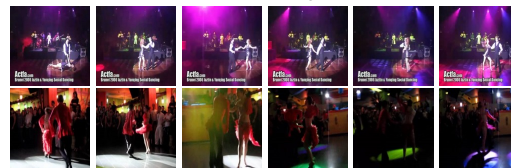
(c) shaking hands



(d) tickling



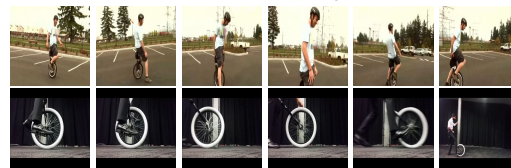
(e) robot dancing



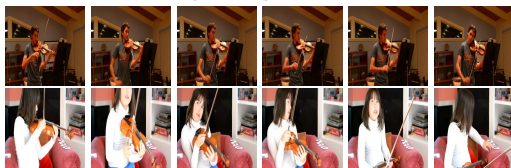
(f) salsa dancing



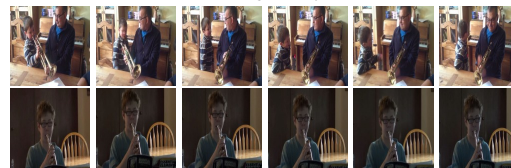
(g) riding a bike



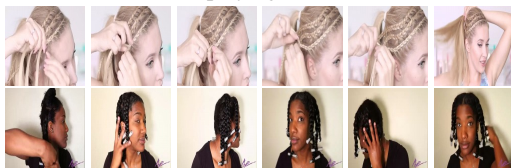
(h) riding unicycle



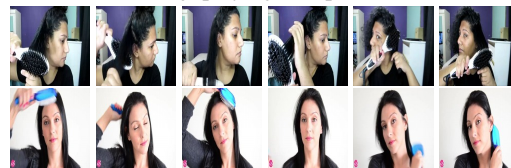
(i) playing violin



(j) playing trumpet



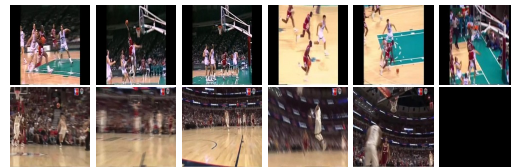
(k) braiding hair



(l) brushing hair



(m) dribbling basketball



(n) dunking basketball

Figure 1: Example classes from the Kinetics dataset. Best seen in colour and with zoom. Note that in some cases a single image is not enough for recognizing the action (e.g. “headbanging”) or distinguishing classes (“dribbling basketball” vs “dunking basketball”). The dataset contains: Singular Person Actions (e.g. “robot dancing”, ”stretching leg”); Person-Person Actions (e.g. “shaking hands”, ”tickling”); Person-Object Actions (e.g. “riding a bike”); same verb different objects (e.g. “playing violin”, “playing trumpet”); and same object different verbs (e.g. “dribbling basketball”, “dunking basketball”). These are realistic (amateur) videos – there is often significant camera shake, for instance.

and clean up the dataset. We then discuss possible biases in the dataset due to the collection process.

Overview: clips for each class were obtained by first searching on YouTube for candidates, and then using Amazon Mechanical Turkers (AMT) to decide if the clip contains the action or not. Three or more confirmations (out of five) were required before a clip was accepted. The dataset was de-duped, by checking that only one clip is taken from each video, and that clips do not contain common video material. Finally, classes were checked for overlap and de-noised.

We now describe these stages in more detail.

3.1. Stage 1: Obtaining an action list

Curating a large list of human actions is challenging, as there is no single listing available at this scale with suitable visual action classes. Consequently, we had to combine numerous sources together with our own observations of actions that surround us. These sources include: (i) **Action datasets** – existing datasets like ActivityNet [3], HMDB [15], UCF101 [20], MPII Human Pose [2], ACT [25] have useful classes and a suitable sub set of these were used; (ii) **Motion capture** – there are a number of motion capture datasets which we looked through and extracted file titles. These titles described the motion within the file and were often quite creative; and, (iii) **Crowd-sourced** – we asked Mechanical Turk workers to come up with a more appropriate action if the label we had presented to them for a clip was incorrect.

3.2. Stage 2: Obtaining candidate clips

The chosen method and steps are detailed below which combine a number of different internal efforts:

Step 1: obtaining videos. Videos are drawn from the YouTube corpus by matching video titles with the Kinetics actions list.

Step 2: temporal positioning within a video. Image classifiers are available for a large number of human actions. These classifiers are obtained by tracking user actions on Google Image Search. For example, for a search query “climbing tree”, user relevance feedback on images is collected by aggregating across the multiple times that that search query is issued. This relevance feedback is used to select a high-confidence set of images that can be used to train a “climbing tree” image classifier. These classifiers are run at the frame level over the videos found in step 1, and clips extracted around the top k responses (where $k = 2$).

It was found that the action list had a better match to relevant classifiers if action verbs are formatted to end with

‘ing’. Thinking back to image search, this makes sense as typically if you are searching for an example of someone performing an action you would issue queries like ‘running man’ or ‘brushing hair’ over other tenses like ‘man ran’ or ‘brush hair’.

The output of this stage is a large number of videos and a position in all of them where one of the actions is potentially occurring. 10 second clips are created by taking 5 seconds either side of that position (there are length exceptions when the position is within 5 seconds of the start or end of the video leading to a shorter clip length). The clips are then passed onto the next stage of cleanup through human labelling.

3.3. Stage 3: Manual labelling process

The key aim of this stage was to identify whether the supposed action was actually occurring during a clip or not. A human was required in the loop for this phase and we chose to use Amazon’s Mechanical Turk (AMT) for the task due to the large numbers of high quality workers using the platform.

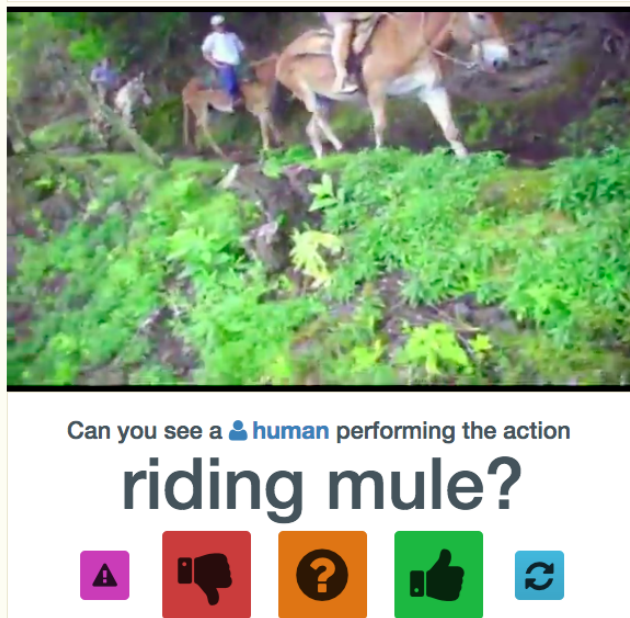
A single-page webapp was built for the labelling task and optimised to maximise the number of clips presented to the workers whilst maintaining a high quality of annotation. The labelling interface is shown in figure 2. The user interface design and theme were chosen to differentiate the task from many others on the platform as well as make the task as stimulating and engaging as possible. This certainly paid off as the task was one of the highest rated on the platform and would frequently get more than 400 distinct workers as soon as a new run was launched.

The workers were given clear instructions at the beginning. There were two screens of instruction, the second reinforcing the first. After acknowledging they understood the task they were presented with a media player and several response icons. The interface would fetch a set of videos from the available pool for the worker at that moment and embed the first clip. The task consisted of 20 videos each with a different class where possible; we randomised all the videos and classes to make it more interesting for the workers and prevent them from becoming stuck on classes with low yields. Two of the video slots were used by us to inject groundtruth clips. This allowed us to get an estimate of the accuracy for each worker. If a worker fell below a 50% success rating on these, we showed them a ‘low accuracy’ warning screen. This helped address many low accuracies.

In the labelling interface, workers were asked the question “Can you see a human performing the action `class-name`?”. The following response options were available on the interface as icons:

- Yes, this contains a true example of the action
- No, this does not contain an example of the action

Evaluating Actions in Videos



Instructions

We would like to find videos that contain real humans performing actions e.g. scrubbing their face, jumping, kissing someone etc.

Please click on the most appropriate button after watching each video:



Yes, this contains a true example of the action



No, this does not contain an example of the action



You are unsure if there is an example of the action



Replay the video



Video does not play, does not contain a human, is an image, cartoon or a computer game.



We have turned off the audio, you need to judge the clip using the visuals only.

Figure 2: Labeling interface used in Mechanical Turk.

- You are unsure if there is an example of the action
- Replay the video
- Video does not play, does not contain a human, is an image, cartoon or a computer game.

When a worker responded with ‘Yes’ we also asked the question “Does the action last for the whole clip?” in order to use this signal later during model training.

Note, the AMT workers didn’t have access to the audio to ensure that the video can be classified purely based on its visual content.

In order for a clip to be added to the dataset, it needed to receive at least 3 positive responses from workers. We allowed each clip to be annotated 5 times except if it had been annotated by more than 2 of a specific response. For example, if 3 out of 3 workers had said it did not contain an example of the action we would immediately remove it from the pool and not continue until 5 workers had annotated it.

Due to the large scale of the task it was necessary to quickly remove classes that were made up of low quality or completely irrelevant candidates. Failing to do this would have meant that we spent a lot of money paying workers to mark videos as negative or bad. Accuracies for each class were calculated after 20 clips from that class had been annotated. We adjusted the accuracy threshold between runs but would typically start at a high accuracy of 50% (1 in 2 videos were expected to contain the action).

Following annotating, the video ids, clip times and labels were exported from the database and handed on to be used for model training.

What we learnt: We found that more specific classes like ‘riding mule’ were producing much less noise than more general classes like ‘riding’. However, occasionally using more general classes was a benefit as they could subsequently be split into a few distinct classes that were not previously present and the candidates resent out to workers e.g. ‘gardening’ was split into ‘watering plants’, ‘trimming trees’ and ‘planting trees’.

The amount of worker traffic that the task generated meant that we could not rely on direct fetching and writes to the database even with appropriate indexes and optimised queries. We therefore created many caches which were made up of groups of clips for each worker. When a worker started a new task, the interface would fetch a set of clips for that specific worker. The cache was replenished often by background processes as clips received a sufficient number of annotations. This also negated labelling collisions where previously > 1 worker might pick up the same video to annotate and we would quickly exceed 5 responses for any 1 clip.

3.4. Stage 4: Cleaning up and de-noising

One of the dataset design goals was having a single clip from each given video sequence, different from existing datasets which slice videos containing repetitive actions into many (correlated) training examples. We also employed mechanisms for identifying structural problems as we grew the dataset, such as repeated classes due to synonymy or different word order (e.g. riding motorbike, riding motorcycle), classes that are too general and co-occur with many others (e.g. talking) and which are problematic for typical 1-of-K classification learning approaches (instead of multi-label classification). We will now describe these procedures.

De-duplicating videos. We de-duplicated videos using two complementary approaches. First, in order to have only one clip from each YouTube link, we randomly selected a single clip from amongst those validated by Turkers for that video. This stage filtered out around 20% of Turker-approved examples, but we visually found that it still left many duplicates. The reason is that YouTube users often create videos reusing portions of other videos, for example as part of video compilations or promotional adverts. Sometimes they are cropped, resized and generally pre-processed in different ways (but, nevertheless, the image classifier could localize the same clip). So even though each clip is from a distinct video there were still duplications.

We devised a process for de-duplicating across YouTube links which operated independently for each class. First we computed Inception-V1 [12] feature vectors (taken after last average pooling layer) on 224×224 center crops of 25 uniformly sampled frames from each video, which we then averaged. Afterwards we built a class-wise matrix having all cosine similarities between these feature vectors and thresholded it. Finally, we computed connected components and kept a random example from each. We found this to work well for most classes using the same threshold of 0.97, but adjusted it in a few cases where classes were visually similar, such as some taking place in the snow or in the water. This process reduced the number of Turker-approved examples by a further 15%.

Detecting noisy classes. Classes can be ‘noisy’ in that they may overlap with other classes or they may contain several quite distinct (in terms of the action) groupings due to an ambiguity in the class name. For example, ‘skipping’ can be ‘skipping with a rope’ and also ‘skipping stones across water’. We trained two-stream action classifiers [19] repeatedly throughout the dataset development to identify these noise classes. This allowed us to find the top confusions for each class, which sometimes were clear even by just verifying the class names (but went unnoticed due

to the scale of the dataset), and other times required eyeballing the data to understand if the confusions were alright and the classes were just difficult to distinguish because of shortcomings of the model. We merged, split or outright removed classes based on these detected confusions.

Final filtering. After all the data was collected, de-duplicated and the classes were selected, we ran a final manual clip filtering stage. Here the class scores from the two-stream model were again useful as they allowed sorting the examples from most confident to least confident – a measure of how prototypical they were. We found that noisy examples were often among the lowest ranked examples and focused on those. The ranking also made adjacent any remaining duplicate videos, which made it easier to filter out those too.

3.5. Discussion: dataset bias I

We are familiar with the notion of dataset bias leading to lack of generalization: where a classifier trained on one dataset, e.g. Caltech 256 [10], does not perform well when tested on another, e.g. PASCAL VOC [8]. Indeed it is even possible to train a classifier to identify which dataset an image belongs to [22].

There is another sense of bias which could arise from unbalanced categories *within* a dataset. For example, gender imbalance in a training set could lead to a corresponding performance bias for classifiers trained on this set. There are precedents for this, e.g. in publicly available face detectors not being race agnostic¹, and more recently in learning a semantic bias in written texts [4]. It is thus an important question as to whether Kinetics leads to such bias.

To this end we carried out a preliminary study on (i) whether the data for each action class of Kinetics is gender balanced, and (ii) if, there is an imbalance, whether it leads to a biased performance of the action classifiers.

The outcome of (i) is that in 340 action classes out of the 400, the data is either not dominated by a single gender, or it is mostly not possible to determine the gender – the latter arises in classes where, for example, only hands appear, or the ‘actors’ are too small or heavily clothed. The classes that do show gender imbalance include ‘shaving beard’ and ‘dunking basketball’, that are mostly male, and ‘filling eyebrows’ and ‘cheerleading’, that are mostly female.

The outcome of (ii) for these classes we found little evidence of classifier bias for action classes with gender imbalance. For example in ‘playing poker’, which tends to have more male players, all videos with female players are correctly classified. The same happens for ‘Hammer throw’. We can conjecture that this lack of bias is because the classifier is able to make use of both the objects involved in

¹<https://www.media.mit.edu/posts/media-lab-student-recognized-for-fighting-bias-in-machine-learning/>

an action as well as the motion patterns, rather than simply physical appearance.

Imbalance can also be examined on other ‘axes’, for example age and race. Again, in a preliminary investigation we found very little clear bias. There is one exception where there is clear bias to babies – in ‘crying’, where many of the videos of non-babies crying are misclassified; another example is ‘wrestling’, where the opposite happens: adults wrestling in a ring seem to be better classified than children wrestling in their homes, but it is hard to tell whether the deciding factor is age or the scenes where the actions happen. Nevertheless, these issues of dataset imbalance and any resulting classifier bias warrant a more thorough investigation, and we return to this in section 5.

3.6. Discussion: dataset bias II

Another type of bias could arise because classifiers are involved in the dataset collection pipeline: it could be that these classifiers lead to a reduction in the visual variety of the clips obtained, which in turn leads to a bias in the action classifier trained on these clips. In more detail, although the videos are selected based on their title (which is provided by the person uploading the video to YouTube), the *position* of the candidate clip within the video is provided by an image (RGB) classifier, as described above. In practice, using a classifier at this point does not seem to constrain the variety of the clips – since the video is about the action, the particular frame chosen as part of the clip may not be crucial; and, in any case, the clip contains hundreds of more frames where the appearance (RGB) and motion can vary considerably. For these reasons we are not so concerned about the intermediate use of image classifiers.

4. Benchmark Performance

In this section we first briefly describe three standard ConvNet architectures for human action recognition in video. We then use these architectures as baselines and compare their performance by training and testing on the Kinetics dataset. We also include their performance on UCF-101 and HMDB-51.

We consider three typical approaches for video classification: ConvNets with an LSTM on top [7, 26]; two-stream networks [9, 19]; and a 3D ConvNet [13, 21, 23]. There have been many improvements over these basic architectures, e.g. [9], but our intention here is not to perform a thorough study on what is the very best architecture on Kinetics, but instead to provide an indication of the level of difficulty of the dataset. A rough graphical overview of the three types of architectures we compare is shown in figure 3, and the specification of their temporal interfaces is given in table 3.

For the experiments on the Kinetics dataset all three architectures are trained from scratch using Kinetics. How-

ever, for the experiments on UCF-101 and HMDB-51 the architectures (apart from the 3D ConvNet) are pre-trained on ImageNet (since these datasets are too small to train the architectures from scratch).

4.1. ConvNet+LSTM

The high performance of image classification networks makes it appealing to try to reuse them with as minimal change as possible for video. This can be achieved by using them to extract features independently from each frame then pooling their predictions across the whole video [14]. This is in the spirit of bag of words image modeling approaches [16, 17, 24], but while convenient in practice, it has the issue of entirely ignoring temporal structure (e.g. models can’t potentially distinguish opening from closing a door).

In theory, a more satisfying approach is to add a recurrent layer to the model [7, 26], such as an LSTM, which can encode state, and capture temporal ordering and long range dependencies. We position an LSTM layer with batch normalization (as proposed by Cooijmans *et al.* [6]) after the last average pooling layer of a ResNet-50 model [11], with 512 hidden units. We then add a fully connected layer on top of the output of the LSTM for the multi-way classification. At test time the classification is taken from the model output for the last frame.

4.2. Two-Stream networks

LSTMs on features from the last layers of ConvNets can model high-level variation, but may not be able to capture fine low-level motion which is critical in many cases. It is also expensive to train as it requires unrolling the network through multiple frames for backpropagation-through-time.

A different, very practical approach, introduced by Simonyan and Zisserman [19], models short temporal snapshots of videos by averaging the predictions from a single RGB frame and a stack of 10 externally computed optical flow frames, after passing them through two replicas of an ImageNet-pretrained ConvNet. The flow stream has an adapted input convolutional layer with twice as many input channels as flow frames (because flow has two channels, horizontal and vertical), and at test time multiple snapshots are sampled from the video and the action prediction is averaged. This was shown to get very high performance on existing benchmarks, while being very efficient to train and test.

4.3. 3D ConvNets

3D ConvNets [13, 21, 23] seem like a natural approach to video modeling. They are just like standard 2D convolutional networks, but with spatio-temporal filters, and have a very interesting characteristic: they directly create hierarchical representations of spatio-temporal data. One issue with these models is that they have many more parameters

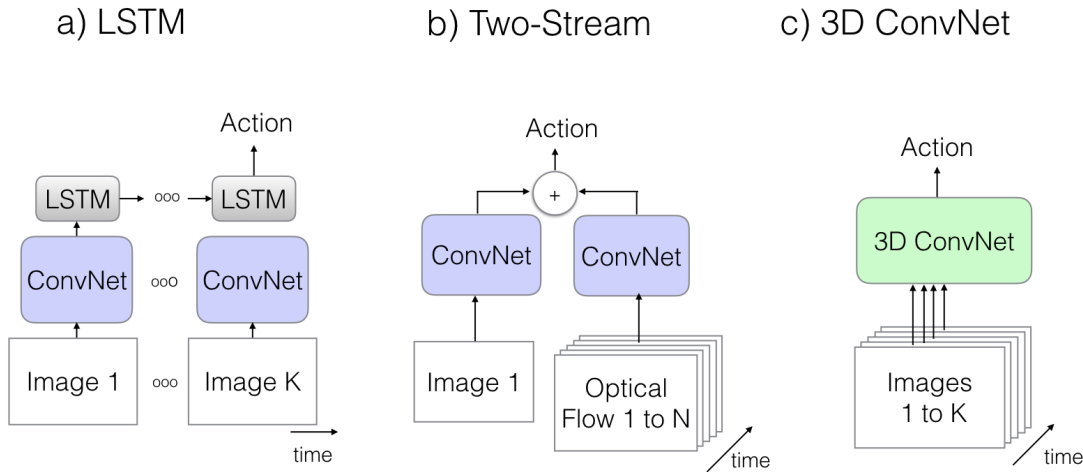


Figure 3: Video architectures used as baseline human action classifiers.

than 2D ConvNets because of the additional kernel dimension, and this makes them harder to train. Also, they seem to preclude the benefits of ImageNet pre-training and previous work has defined relatively shallow custom architectures and trained them from scratch [13, 14, 21, 23]. Results on benchmarks have shown promise but have not yet matched the state-of-the-art, possibly because they require more training data than their 2D counterparts. Thus 3D ConvNets are a good candidate for evaluation on our larger dataset.

For this paper we implemented a small variation of C3D [23], which has 8 convolutional layers, 5 pooling layers and 2 fully connected layers at the top. The inputs to the model are short 16-frame clips with 112×112 -pixel crops. Differently from the original paper we use batch normalization after all convolutional and fully connected layers. Another difference to the original model is in the first pooling layer, where we use a temporal stride of 2 instead of 1, which reduces the memory footprint and allows for bigger batches – this was important for batch normalization (especially after the fully connected layers, where there is no weight tying). Using this stride we were able to train with 15 videos per batch per GPU using standard K40 GPUs.

At test time, we split the video uniformly into crops of 16 frames and apply the classifier separately on each. We then average the class scores, as in the original paper.

4.4. Implementation details

The ConvNet+LSTM and Two-Stream architectures use ResNet-50 as the base architecture. In the case of the Two-Stream architecture, a separate ResNet-50 is trained independently for each stream. As noted earlier, for these architectures the ResNet-50 model is pre-trained on ImageNet for the experiments on UCF-101 and HMDB-51, and trained from scratch for experiments on Kinetics. The 3D-ConvNet is not pre-trained.

We trained the models on videos using standard SGD with momentum in all cases, with synchronous parallelization across 64 GPUs for all models. We trained models on Kinetics for up to 100k steps, with a 10x reduction of learning rate when validation loss saturated, and tuned weight decay and learning rate hyperparameters on the validation set of Kinetics. All the models were implemented in TensorFlow [1].

The original clips have variable resolution and frame rate. In our experiments they are all normalized so that the larger image side is 340 pixels wide for models using ResNet-50 and 128 pixels wide for the 3D ConvNet. We also resample the videos so they have 25 frames per second.

Data augmentation is known to be of crucial importance for the performance of deep architectures. We used random cropping both spatially – randomly cropping a 299×299

Method	#Params	Training		Testing	
		# Input Frames	Temporal Footprint	# Input Frames	Temporal Footprint
(a) ConvNet+LSTM	29M	25 rgb	5s	50 rgb	10s
(b) Two-Stream	48M	1 rgb, 10 flow	0.4s	25 rgb, 250 flow	10s
(c) 3D-ConvNet	79M	16 rgb	0.64s	240 rgb	9.6s

Table 3: Number of parameters and temporal input sizes of the models. ConvNet+LSTM and Two-Stream use ResNet-50 ConvNet modules.

Architecture	UCF-101			HMDB-51			Kinetics		
	RGB	Flow	RGB+Flow	RGB	Flow	RGB+Flow	RGB	Flow	RGB+Flow
(a) ConvNet+LSTM	84.3	–	–	43.9	–	–	57.0 / 79.0	–	–
(b) Two-Stream	84.2	85.9	92.5	51.0	56.9	63.7	56.0 / 77.3	49.5 / 71.9	61.0 / 81.3
(c) 3D-ConvNet	51.6	–	–	24.3	–	–	56.1 / 79.5	–	–

Table 4: Baseline comparisons across datasets: (left) training and testing on split 1 of UCF-101; (middle) training and testing on split 1 of HMDB-51; (right) training and testing on Kinetics (showing top-1/top-5 performance). ConvNet+LSTM and Two-Stream use ResNet-50 ConvNet modules, pretrained on ImageNet for UCF-101 and HMDB-51 examples but not for the Kinetics experiments. Note that the Two-Stream architecture numbers on individual RGB and Flow streams can be interpreted as a simple baseline which applies a ConvNet independently on 25 uniformly sampled frames then averages the predictions.

patch (respectively 112×112 for the 3D ConvNet) – and temporally, when picking the starting frame among those early enough to guarantee a desired number of frames. For shorter videos, we looped the video as many times as necessary to satisfy each model’s input interface. We also applied random left-right flipping consistently for each video during training.

At test time, we sample from up to 10 seconds of video, again looping if necessary. Better performance could be obtained by also considering left-right flipped videos at test time and by adding additional augmentation, such as photometric, during training. We leave this to future work.

4.5. Baseline evaluations

In this section we compare the performance of the three baseline architectures whilst varying the dataset used for training and testing.

Table 4 shows the classification accuracy when training and testing on either UCF-101, HMDB-51 or Kinetics. We train and test on split 1 of UCF-101 and HMDB-51, and on the train/val set and held-out test set of Kinetics.

There are several noteworthy observations. First, the performance is far lower on Kinetics than on UCF-101, an indication of the different levels of difficulty of the two datasets. On the other hand, the performance on HMDB-51 is worse than on Kinetics – it seems to have a truly difficult test set, and it was designed to be difficult to appearance-centered methods, while having little training data. The parameter-rich 3D-ConvNet model is not pre-trained on ImageNet,

unlike the other baselines. This translates into poor performance on all datasets but especially on UCF-101 and HMDB-51 – on Kinetics it is much closer to the performance of the other models, thanks to the much larger training set of Kinetics.

- **Class difficulty.** We include a full list of Kinetics classes sorted by classification accuracy under the two-stream model in figure 4. Eating classes are among the hardest, as they sometimes require distinguishing what is being eaten, such as hotdogs, chips and doughnuts – and these may appear small and already partially consumed, in the video. Dancing classes are also hard, as well as classes centered on a specific body part, such as “massaging feet”, or “shaking head”.
- **Class confusion.** The top 10 class confusions are provided in table 5. They mostly correspond to fine-grained distinctions that one would expect to be hard, for example ‘long jump’ and ‘triple jump’, confusing burger with doughnuts. The confusion between ‘swing dancing’ and ‘salsa dancing’ raises the question of how accurate motion modeling is in the two-stream model, since ‘swing dancing’ is typically much faster-paced and has a peculiar style that makes it easy for humans to distinguish from salsa.
- **Classes where motion matters most.** We tried to analyze for which classes motion is more important and

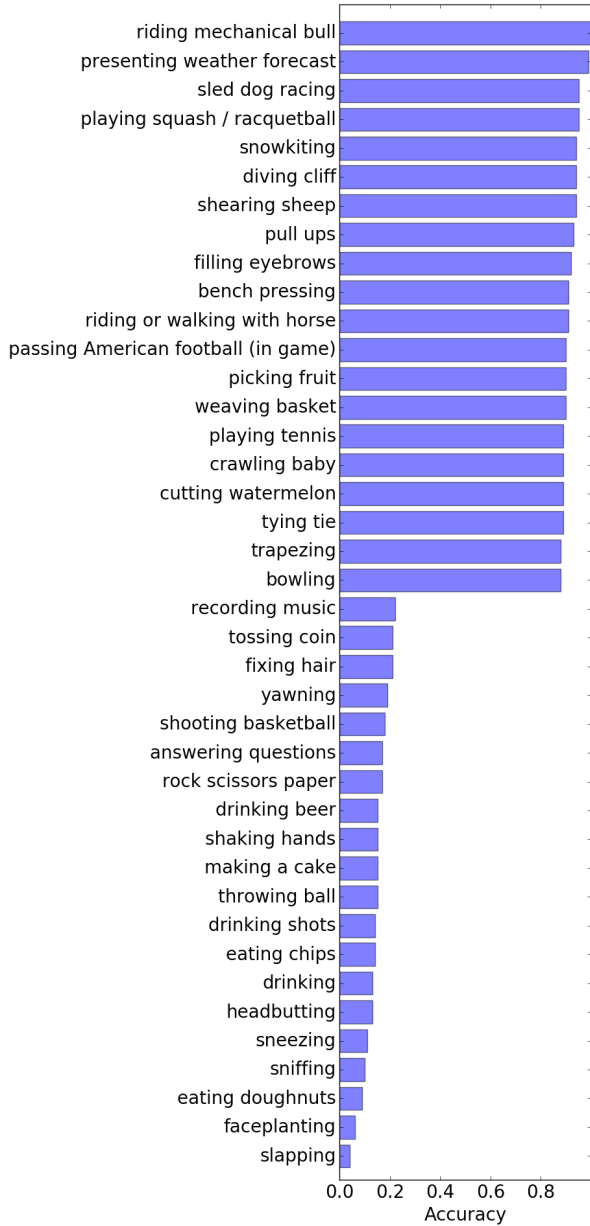


Figure 4: List of 20 easiest and 20 hardest Kinetics classes sorted by class accuracies obtained using the two-stream model.

which ones were recognized correctly using just appearance information, by comparing the recognition accuracy ratios when using the flow and RGB streams of the two-stream model in isolation. We show the five classes where this ratio is largest and smallest in table 6.

5. Conclusion

We have described the Kinetics Human Action Video dataset, which has an order of magnitude more videos than previous datasets of its type. We have also discussed the procedures we employed collecting the data and for ensuring its quality. We have shown that the performance of standard existing models on this dataset is much lower than on UCF-101 and on par with HMDB-51, whilst allowing large models such as 3D ConvNets to be trained from scratch, unlike the existing human action datasets.

We have also carried out a preliminary analysis of dataset imbalance and whether this leads to bias in the classifiers trained on the dataset. We found little evidence that the resulting classifiers demonstrate bias along sensitive axes, such as across gender. This is however a complex area that deserves further attention. We leave a thorough analysis for future work, in collaboration with specialists from complementary areas, namely social scientists and critical humanists.

We will release trained baseline models (in TensorFlow), so that they can be used, for example, to generate features for new action classes.

Acknowledgements:

The collection of this dataset was funded by DeepMind. We are very grateful for help from Andreas Kirsch, John-Paul Holt, Danielle Breen, Jonathan Fildes, James Besley and Brian Carver. We are grateful for advice and comments from Tom Duerig, Juan Carlos Niebles, Simon Osindero, Chuck Rosenberg and Sean Legassick; we would also like to thank Sandra and Aditya for data clean up.

References

- [1] M. Abadi, A. Agarwal, P. Barham, E. Brevdo, Z. Chen, C. Citro, G. S. Corrado, A. Davis, J. Dean, M. Devin, et al. Tensorflow: Large-scale machine learning on heterogeneous distributed systems. *arXiv preprint arXiv:1603.04467*, 2016.
- [2] M. Andriluka, L. Pishchulin, P. Gehler, and B. Schiele. 2d human pose estimation: New benchmark and state of the art analysis. In *Computer Vision and Pattern Recognition (CVPR), 2014 IEEE Conference on*. IEEE, 2014.
- [3] F. Caba Heilbron, V. Escorcia, B. Ghanem, and J. C. Niebles. Activitynet: A large-scale video benchmark for human activity understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015.
- [4] A. Caliskan, J. J. Bryson, and A. Narayanan. Semantics derived automatically from language corpora contain human-like biases. *Science*, 356(6334):183–186, 2017.
- [5] J. Carreira and A. Zisserman. Quo vadis, action recognition? new models and the kinetics dataset. In *IEEE International Conference on Computer Vision and Pattern Recognition CVPR*, 2017.
- [6] T. Cooijmans, N. Ballas, C. Laurent, and A. Courville.

Class 1	Class 2	confusion
‘riding mule’	‘riding or walking with horse’	40%
‘hockey stop’	‘ice skating’	36%
‘swing dancing’	‘salsa dancing’	36%
‘strumming guitar’	‘playing guitar’	35%
‘shooting basketball’	‘playing basketball’	32%
‘cooking sausages’	‘cooking chicken’	29%
‘sweeping floor’	‘mopping floor’	27%
‘triple jump’	‘long jump’	26%
‘doing aerobics’	‘zumba’	26%
‘petting animal (not cat)’	‘feeding goats’	25%
‘shaving legs’	‘waxing legs’	25%
‘snowboarding’	‘skiing (not slalom or crosscountry)’	22%

Table 5: Top-12 class confusions in Kinetics, using the two-stream model.

Class	Flow/RGB accuracy ratio
‘rock scissors paper’	5.3
‘sword fighting’	3.1
‘robot dancing’	3.1
‘air drumming’	2.8
‘exercising arm’	2.5
‘making a cake’	0.1
‘cooking sausages’	0.1
‘sniffing’	0.1
‘eating cake’	0.0
‘making a sandwich’	0.0

Table 6: Classes with largest and smallest ratios of recognition accuracy when using flow and RGB. The highest ratios correspond to when flow does better, the smallest to when RGB does better. We also evaluated the ratios of rgb+flow to rgb accuracies and the ordering was quite similar.

- Recurrent batch normalization. *arXiv preprint arXiv:1603.09025*, 2016.
- [7] J. Donahue, L. Anne Hendricks, S. Guadarrama, M. Rohrbach, S. Venugopalan, K. Saenko, and T. Darrell. Long-term recurrent convolutional networks for visual recognition and description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2625–2634, 2015.
- [8] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–136, 2015.
- [9] C. Feichtenhofer, A. Pinz, and A. Zisserman. Convolutional two-stream network fusion for video action recognition. In *IEEE International Conference on Computer Vision and Pattern Recognition CVPR*, 2016.
- [10] G. Griffin, A. Holub, and P. Perona. Caltech-256 object category dataset. 2007.
- [11] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Computer Vision and Pattern Recognition (CVPR), 2016 IEEE Conference on*, 2016.
- [12] S. Ioffe and C. Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- [13] S. Ji, W. Xu, M. Yang, and K. Yu. 3d convolutional neural networks for human action recognition. *IEEE transactions on pattern analysis and machine intelligence*, 35(1):221–231, 2013.
- [14] A. Karpathy, G. Toderici, S. Shetty, T. Leung, R. Sukthankar, and L. Fei-Fei. Large-scale video classification with convolutional neural networks. In *Proceedings of the IEEE conference on Computer Vision and Pattern Recognition*, pages 1725–1732, 2014.
- [15] H. Kuehne, H. Jhuang, E. Garrote, T. Poggio, and T. Serre. HMDB: a large video database for human motion recognition. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2011.
- [16] I. Laptev, M. Marszalek, C. Schmid, and B. Rozenfeld. Learning realistic human actions from movies. In *Computer Vision and Pattern Recognition, 2008. CVPR 2008. IEEE Conference on*, pages 1–8. IEEE, 2008.
- [17] J. C. Niebles, H. Wang, and L. Fei-Fei. Unsupervised learning of human action categories using spatial-temporal words. *International journal of computer vision*, 79(3):299–318, 2008.
- [18] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, S. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. Berg, and F. Li. Imagenet large scale visual recognition challenge. *IJCV*, 2015.
- [19] K. Simonyan and A. Zisserman. Two-stream convolutional networks for action recognition in videos. In *Advances in Neural Information Processing Systems*, pages 568–576, 2014.
- [20] K. Soomro, A. R. Zamir, and M. Shah. Ucf101: A dataset of 101 human actions classes from videos in the wild. *arXiv preprint arXiv:1212.0402*, 2012.
- [21] G. W. Taylor, R. Fergus, Y. LeCun, and C. Bregler. Convolutional learning of spatio-temporal features. In *European con-*

- ference on computer vision*, pages 140–153. Springer, 2010.
- [22] A. Torralba and A. A. Efros. Unbiased look at dataset bias. In *Computer Vision and Pattern Recognition (CVPR), 2011 IEEE Conference on*, pages 1521–1528. IEEE, 2011.
- [23] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features with 3d convolutional networks. In *2015 IEEE International Conference on Computer Vision (ICCV)*, pages 4489–4497. IEEE, 2015.
- [24] H. Wang and C. Schmid. Action recognition with improved trajectories. In *International Conference on Computer Vision*, 2013.
- [25] X. Wang, A. Farhadi, and A. Gupta. Actions ~ transformations. In *CVPR*, 2016.
- [26] J. Yue-Hei Ng, M. Hausknecht, S. Vijayanarasimhan, O. Vinyals, R. Monga, and G. Toderici. Beyond short snippets: Deep networks for video classification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4694–4702, 2015.

A. List of Kinetics Human Action Classes

This is the list of classes included in the human action video dataset. The number of clips for each action class is given by the number in brackets following each class name.

1. abseiling (1146)
2. air drumming (1132)
3. answering questions (478)
4. applauding (411)
5. applying cream (478)
6. archery (1147)
7. arm wrestling (1123)
8. arranging flowers (583)
9. assembling computer (542)
10. auctioning (478)
11. baby waking up (611)
12. baking cookies (927)
13. balloon blowing (826)
14. bandaging (569)
15. barbequing (1070)
16. bartending (601)
17. beatboxing (943)
18. bee keeping (430)
19. belly dancing (1115)
20. bench pressing (1106)
21. bending back (635)
22. bending metal (410)
23. biking through snow (1052)
24. blasting sand (713)
25. blowing glass (1145)
26. blowing leaves (405)
27. blowing nose (597)
28. blowing out candles (1150)
29. bobsledding (605)
30. bookbinding (914)
31. bouncing on trampoline (690)
32. bowling (1079)
33. braiding hair (780)
34. breeding or breadcrumbing (454)
35. breakdancing (948)
36. brush painting (532)

37. brushing hair (934)
38. brushing teeth (1149)
39. building cabinet (431)
40. building shed (427)
41. bungee jumping (1056)
42. busking (851)
43. canoeing or kayaking (1146)
44. capoeira (1092)
45. carrying baby (558)
46. cartwheeling (616)
47. carving pumpkin (711)
48. catching fish (671)
49. catching or throwing baseball (756)
50. catching or throwing frisbee (1060)
51. catching or throwing softball (842)
52. celebrating (751)
53. changing oil (714)
54. changing wheel (459)
55. checking tires (555)
56. cheerleading (1145)
57. chopping wood (916)
58. clapping (491)
59. clay pottery making (513)
60. clean and jerk (902)
61. cleaning floor (874)
62. cleaning gutters (598)
63. cleaning pool (447)
64. cleaning shoes (706)
65. cleaning toilet (576)
66. cleaning windows (695)
67. climbing a rope (413)
68. climbing ladder (662)
69. climbing tree (1120)
70. contact juggling (1135)
71. cooking chicken (1000)
72. cooking egg (618)
73. cooking on campfire (403)
74. cooking sausages (467)
75. counting money (674)
76. country line dancing (1015)
77. cracking neck (449)
78. crawling baby (1150)
79. crossing river (951)
80. crying (1037)
81. curling hair (855)
82. cutting nails (560)
83. cutting pineapple (712)
84. cutting watermelon (767)
85. dancing ballet (1144)
86. dancing charleston (721)
87. dancing gangnam style (836)
88. dancing macarena (958)
89. deadlifting (805)
90. decorating the christmas tree (612)
91. digging (404)
92. dining (671)
93. disc golfing (565)
94. diving cliff (1075)
95. dodgeball (595)
96. doing aerobics (461)
97. doing laundry (461)
98. doing nails (949)
99. drawing (445)
100. dribbling basketball (923)
101. drinking (599)
102. drinking beer (575)
103. drinking shots (403)
104. driving car (1118)
105. driving tractor (922)
106. drop kicking (716)
107. drumming fingers (409)
108. dunking basketball (1105)
109. dying hair (1072)
110. eating burger (864)
111. eating cake (494)
112. eating carrots (516)
113. eating chips (749)
114. eating doughnuts (528)
115. eating hotdog (570)
116. eating ice cream (927)

117. eating spaghetti (1145)
118. eating watermelon (550)
119. egg hunting (500)
120. exercising arm (416)
121. exercising with an exercise ball (438)
122. extinguishing fire (602)
123. faceplanting (441)
124. feeding birds (1150)
125. feeding fish (973)
126. feeding goats (1027)
127. filling eyebrows (1085)
128. finger snapping (825)
129. fixing hair (676)
130. flipping pancake (720)
131. flying kite (1063)
132. folding clothes (695)
133. folding napkins (874)
134. folding paper (940)
135. front raises (962)
136. frying vegetables (608)
137. garbage collecting (441)
138. gargling (430)
139. getting a haircut (658)
140. getting a tattoo (737)
141. giving or receiving award (953)
142. golf chipping (699)
143. golf driving (836)
144. golf putting (1081)
145. grinding meat (415)
146. grooming dog (613)
147. grooming horse (645)
148. gymnastics tumbling (1143)
149. hammer throw (1148)
150. headbanging (1090)
151. headbutting (640)
152. high jump (954)
153. high kick (825)
154. hitting baseball (1071)
155. hockey stop (468)
156. holding snake (430)
157. hopscotch (726)
158. hoverboarding (564)
159. hugging (517)
160. hula hooping (1129)
161. hurdling (622)
162. hurling (sport) (836)
163. ice climbing (845)
164. ice fishing (555)
165. ice skating (1140)
166. ironing (535)
167. javelin throw (912)
168. jetskiing (1140)
169. jogging (417)
170. juggling balls (923)
171. juggling fire (668)
172. juggling soccer ball (484)
173. jumping into pool (1133)
174. jumpstyle dancing (662)
175. kicking field goal (833)
176. kicking soccer ball (544)
177. kissing (733)
178. kitesurfing (794)
179. knitting (691)
180. krumping (657)
181. laughing (926)
182. laying bricks (432)
183. long jump (831)
184. lunge (759)
185. making a cake (463)
186. making a sandwich (440)
187. making bed (679)
188. making jewelry (658)
189. making pizza (1147)
190. making snowman (756)
191. making sushi (434)
192. making tea (426)
193. marching (1146)
194. massaging back (1113)
195. massaging feet (478)
196. massaging legs (592)

197. massaging person (672)
198. milking cow (980)
199. mopping floor (606)
200. motorcycling (1142)
201. moving furniture (426)
202. mowing lawn (1147)
203. news anchoring (420)
204. opening bottle (732)
205. opening present (866)
206. paragliding (800)
207. parasailing (762)
208. parkour (504)
209. passing American football (in game) (863)
210. passing American football (not in game) (1045)
211. peeling apples (592)
212. peeling potatoes (457)
213. petting animal (not cat) (757)
214. petting cat (756)
215. picking fruit (793)
216. planting trees (557)
217. plastering (428)
218. playing accordion (925)
219. playing badminton (944)
220. playing bagpipes (838)
221. playing basketball (1144)
222. playing bass guitar (1135)
223. playing cards (737)
224. playing cello (1081)
225. playing chess (850)
226. playing clarinet (1022)
227. playing controller (524)
228. playing cricket (949)
229. playing cymbals (636)
230. playing didgeridoo (787)
231. playing drums (908)
232. playing flute (475)
233. playing guitar (1135)
234. playing harmonica (1006)
235. playing harp (1149)
236. playing ice hockey (917)
237. playing keyboard (715)
238. playing kickball (468)
239. playing monopoly (731)
240. playing organ (672)
241. playing paintball (1140)
242. playing piano (691)
243. playing poker (1134)
244. playing recorder (1148)
245. playing saxophone (916)
246. playing squash or racquetball (980)
247. playing tennis (1144)
248. playing trombone (1149)
249. playing trumpet (989)
250. playing ukulele (1146)
251. playing violin (1142)
252. playing volleyball (804)
253. playing xylophone (746)
254. pole vault (984)
255. presenting weather forecast (1050)
256. pull ups (1121)
257. pumping fist (1009)
258. pumping gas (544)
259. punching bag (1150)
260. punching person (boxing) (483)
261. push up (614)
262. pushing car (1069)
263. pushing cart (1150)
264. pushing wheelchair (465)
265. reading book (1148)
266. reading newspaper (424)
267. recording music (415)
268. riding a bike (476)
269. riding camel (716)
270. riding elephant (1104)
271. riding mechanical bull (698)
272. riding mountain bike (495)
273. riding mule (476)
274. riding or walking with horse (1131)
275. riding scooter (674)
276. riding unicycle (864)

277. ripping paper (605)
278. robot dancing (893)
279. rock climbing (1144)
280. rock scissors paper (424)
281. roller skating (960)
282. running on treadmill (428)
283. sailing (867)
284. salsa dancing (1148)
285. sanding floor (574)
286. scrambling eggs (816)
287. scuba diving (968)
288. setting table (478)
289. shaking hands (640)
290. shaking head (885)
291. sharpening knives (424)
292. sharpening pencil (752)
293. shaving head (971)
294. shaving legs (509)
295. shearing sheep (988)
296. shining shoes (615)
297. shooting basketball (595)
298. shooting goal (soccer) (444)
299. shot put (987)
300. shoveling snow (879)
301. shredding paper (403)
302. shuffling cards (828)
303. side kick (991)
304. sign language interpreting (446)
305. singing (1147)
306. situp (817)
307. skateboarding (1139)
308. ski jumping (1051)
309. skiing (not slalom or crosscountry) (1140)
310. skiing crosscountry (477)
311. skiing slalom (539)
312. skipping rope (488)
313. skydiving (505)
314. slacklining (790)
315. slapping (465)
316. sled dog racing (775)
317. smoking (1105)
318. smoking hookah (857)
319. snatch weight lifting (943)
320. sneezing (505)
321. sniffing (399)
322. snorkeling (1012)
323. snowboarding (937)
324. snowkiting (1145)
325. snowmobiling (601)
326. somersaulting (993)
327. spinning poi (1134)
328. spray painting (908)
329. spraying (470)
330. springboard diving (406)
331. squat (1148)
332. sticking tongue out (770)
333. stomping grapes (444)
334. stretching arm (718)
335. stretching leg (829)
336. strumming guitar (472)
337. surfing crowd (876)
338. surfing water (751)
339. sweeping floor (604)
340. swimming backstroke (1077)
341. swimming breast stroke (833)
342. swimming butterfly stroke (678)
343. swing dancing (512)
344. swinging legs (409)
345. swinging on something (482)
346. sword fighting (473)
347. tai chi (1070)
348. taking a shower (378)
349. tango dancing (1114)
350. tap dancing (947)
351. tapping guitar (815)
352. tapping pen (703)
353. tasting beer (588)
354. tasting food (613)
355. testifying (497)
356. texting (704)

- 357. throwing axe (816)
- 358. throwing ball (634)
- 359. throwing discus (1104)
- 360. tickling (610)
- 361. tobogganing (1147)
- 362. tossing coin (461)
- 363. tossing salad (463)
- 364. training dog (481)
- 365. trapezing (786)
- 366. trimming or shaving beard (981)
- 367. trimming trees (665)
- 368. triple jump (784)
- 369. tying bow tie (387)
- 370. tying knot (not on a tie) (844)
- 371. tying tie (673)
- 372. unboxing (858)
- 373. unloading truck (406)
- 374. using computer (937)
- 375. using remote controller (not gaming) (549)
- 376. using segway (387)
- 377. vault (562)
- 378. waiting in line (430)
- 379. walking the dog (1145)
- 380. washing dishes (1048)
- 381. washing feet (862)
- 382. washing hair (423)
- 383. washing hands (916)
- 384. water skiing (763)
- 385. water sliding (420)
- 386. watering plants (680)
- 387. waxing back (537)
- 388. waxing chest (760)
- 389. waxing eyebrows (720)
- 390. waxing legs (948)
- 391. weaving basket (743)
- 392. welding (759)
- 393. whistling (416)
- 394. windsurfing (1114)
- 395. wrapping present (861)
- 396. wrestling (488)

- 397. writing (735)
- 398. yawning (398)
- 399. yoga (1140)
- 400. zumba (1093)

B. List of Parent-Child Groupings

These lists are not exclusive and are not intended to be comprehensive. Rather, they are a guide for related human action classes.

arts and crafts (12)

- arranging flowers
- blowing glass
- brush painting
- carving pumpkin
- clay pottery making
- decorating the christmas tree
- drawing
- getting a tattoo
- knitting
- making jewelry
- spray painting
- weaving basket

athletics – jumping (6)

- high jump
- hurdling
- long jump
- parkour
- pole vault
- triple jump

athletics – throwing + launching (9)

- archery
- catching or throwing frisbee
- disc golfing
- hammer throw
- javelin throw
- shot put
- throwing axe
- throwing ball
- throwing discus

auto maintenance (4)

- changing oil
- changing wheel
- checking tires
- pumping gas

ball sports (25)

- bowling
- catching or throwing baseball

catching or throwing softball
dodgeball
dribbling basketball
dunking basketball
golf chipping
golf driving
golf putting
hitting baseball
hurling (sport)
juggling soccer ball
kicking field goal
kicking soccer ball
passing American football (in game)
passing American football (not in game)
playing basketball
playing cricket
playing kickball
playing squash or racquetball
playing tennis
playing volleyball
shooting basketball
shooting goal (soccer)
shot put

body motions (16)

air drumming
applauding
baby waking up
bending back
clapping
cracking neck
drumming fingers
finger snapping
headbanging
headbutting
pumping fist
shaking head
stretching arm
stretching leg
swinging legs

cleaning (13)

cleaning floor
cleaning gutters
cleaning pool
cleaning shoes
cleaning toilet
cleaning windows
doing laundry
making bed
mopping floor
setting table
shining shoes

sweeping floor
washing dishes

cloths (8)

bandaging
doing laundry
folding clothes
folding napkins
ironing
making bed
tying bow tie
tying knot (not on a tie)
tying tie

communication (11)

answering questions
auctioning
bartending
celebrating
crying
giving or receiving award
laughing
news anchoring
presenting weather forecast
sign language interpreting
testifying

cooking (22)

baking cookies
barbequing
breading or breadcrumbing
cooking chicken
cooking egg
cooking on campfire
cooking sausages
cutting pineapple
cutting watermelon
flipping pancake
frying vegetables
grinding meat
making a cake
making a sandwich
making pizza
making sushi
making tea
peeling apples
peeling potatoes
picking fruit
scrambling eggs
tossing salad

dancing (18)

belly dancing

breakdancing
capoeira
cheerleading
country line dancing
dancing ballet
dancing charleston
dancing gangnam style
dancing macarena
jumpstyle dancing
krumping
marching
robot dancing
salsa dancing
swing dancing
tango dancing
tap dancing
zumba

eating + drinking (17)

bartending
dining
drinking
drinking beer
drinking shots
eating burger
eating cake
eating carrots
eating chips
eating doughnuts
eating hotdog
eating ice cream
eating spaghetti
eating watermelon
opening bottle
tasting beer
tasting food

electronics (5)

assembling computer
playing controller
texting
using computer
using remote controller (not gaming)

garden + plants (10)

blowing leaves
carving pumpkin
chopping wood
climbing tree
decorating the christmas tree
egg hunting
mowing lawn
planting trees

trimming trees
watering plants

golf (3)

golf chipping
golf driving
golf putting

gymnastics (5)

bouncing on trampoline
cartwheeling
gymnastics tumbling
somersaulting
vault

hair (14)

braiding hair
brushing hair
curling hair
dying hair
fixing hair
getting a haircut
shaving head
shaving legs
trimming or shaving beard
washing hair
waxing back
waxing chest
waxing eyebrows
waxing legs

hands (9)

air drumming
applauding
clapping
cutting nails
doing nails
drumming fingers
finger snapping
pumping fist
washing hands

head + mouth (17)

balloon blowing
beatboxing
blowing nose
blowing out candles
brushing teeth
gargling
headbanging
headbutting
shaking head
singing

smoking
smoking hookah
sneezing
sniffing
sticking tongue out
whistling
yawning

heights (15)

abseiling
bungee jumping
climbing a rope
climbing ladder
climbing tree
diving cliff
ice climbing
jumping into pool
paragliding
rock climbing
skydiving
slacklining
springboard diving
swinging on something
trapezing

interacting with animals (19)

bee keeping
catching fish
feeding birds
feeding fish
feeding goats
grooming dog
grooming horse
holding snake
ice fishing
milking cow
petting animal (not cat)
petting cat
riding camel
riding elephant
riding mule
riding or walking with horse
shearing sheep
training dog
walking the dog

juggling (6)

contact juggling
hula hooping
juggling balls
juggling fire
juggling soccer ball
spinning poi

makeup (5)

applying cream
doing nails
dying hair
filling eyebrows
getting a tattoo

martial arts (10)

arm wrestling
capoeira
drop kicking
high kick
punching bag
punching person
side kick
sword fighting
tai chi
wrestling

miscellaneous (9)

digging
extinguishing fire
garbage collecting
laying bricks
moving furniture
spraying
stomping grapes
tapping pen
unloading truck

mobility – land (20)

crawling baby
driving car
driving tractor
faceplanting
hoverboarding
jogging
motorcycling
parkour
pushing car
pushing cart
pushing wheelchair
riding a bike
riding mountain bike
riding scooter
riding unicycle
roller skating
running on treadmill
skateboarding
surfing crowd
using segway
waiting in line

mobility – water (10)

crossing river
diving cliff
jumping into pool
scuba diving
snorkeling
springboard diving
swimming backstroke
swimming breast stroke
swimming butterfly stroke
water sliding

music (29)

beatboxing
busking
playing accordion
playing bagpipes
playing bass guitar
playing cello
playing clarinet
playing cymbals
playing didgeridoo
playing drums
playing flute
playing guitar
playing harmonica
playing harp
playing keyboard
playing organ
playing piano
playing recorder
playing saxophone
playing trombone
playing trumpet
playing ukulele
playing violin
playing xylophone
recording music
singing
strumming guitar
tapping guitar
whistling

paper (12)

bookbinding
counting money
folding napkins
folding paper
opening present
reading book
reading newspaper
ripping paper

shredding paper
unboxing
wrapping present
writing

personal hygiene (6)

brushing teeth
taking a shower
trimming or shaving beard
washing feet
washing hair
washing hands

playing games (13)

egg hunting
flying kite
hopscotch
playing cards
playing chess
playing monopoly
playing paintball
playing poker
riding mechanical bull
rock scissors paper
shuffling cards
skipping rope
tossing coin

racquet + bat sports (8)

catching or throwing baseball
catching or throwing softball
hitting baseball
hurling (sport)
playing badminton
playing cricket
playing squash or racquetball
playing tennis

snow + ice (18)

biking through snow
bobsledding
hockey stop
ice climbing
ice fishing
ice skating
making snowman
playing ice hockey
shoveling snow
ski jumping
skiing (not slalom or crosscountry)
skiing crosscountry
skiing slalom
sled dog racing

snowboarding
snowkiting
snowmobiling
tobogganing

swimming (3)

swimming backstroke
swimming breast stroke
swimming butterfly stroke

touching person (11)

carrying baby
hugging
kissing
massaging back
massaging feet
massaging legs
massaging person's head
shaking hands
slapping
tickling

using tools (13)

bending metal
blasting sand
building cabinet
building shed
changing oil
changing wheel
checking tires
plastering
pumping gas
sanding floor
sharpening knives
sharpening pencil
welding

water sports (8)

canoeing or kayaking
jetskiing
kitesurfing
parasailing
sailing
surfing water
water skiing
windsurfing

waxing (4)

waxing back
waxing chest
waxing eyebrows
waxing legs