

A Short Note about Kinetics-600

João Carreira

joaoluis@google.com

Eric Noland

enoland@google.com

Andras Banki-Horvath

bhandras@google.com

Chloe Hillier

chillier@google.com

Andrew Zisserman

zisserman@google.com

Abstract

We describe an extension of the DeepMind Kinetics human action dataset from 400 classes, each with at least 400 video clips, to 600 classes, each with at least 600 video clips. In order to scale up the dataset we changed the data collection process so it uses multiple queries per class, with some of them in a language other than english – portuguese. This paper details the changes between the two versions of the dataset and includes a comprehensive set of statistics of the new version as well as baseline results using the I3D neural network architecture. The paper is a companion to the release of the ground truth labels for the public test set.

1. Introduction

The release of the Kinetics dataset [6] in 2017 led to marked improvements in state-of-the-art performance on a variety of action recognition datasets: UCF-101 [9], HMDB-51 [7], Charades [8], AVA [3], Thumos [5], among others. Video models pre-trained on Kinetics generalized well when transferred to different video tasks on smaller video datasets, similar to what happened to image classifiers trained on ImageNet.

The goal of the Kinetics project from the start was to replicate the size of ImageNet, which has 1000 classes, each with 1000 image examples. This proved difficult initially and the first version of the dataset had 400 classes, each with 400 video clip examples. There were two main bottlenecks and they were related: (a) identifying relevant candidate YouTube videos for each action class, and (b) finding classes having many candidates. Problem (b) was particularly acute and exposed inefficiencies with the way videos were selected – querying YouTube for simple variations of the class names, by varying singular/plural of nouns, adding articles (e.g. “catching a ball” / “catching ball”), etc. These problems have now been overcome, as described in the sequel.

The new version of the dataset, called Kinetics-600, follows the same principles as Kinetics-400: (i) The clips are from YouTube video, last 10s, and have a variable resolution and frame rate; (ii) for an action class, all clips are from different YouTube videos. Kinetics-600 represents a 50% increase in number of classes, from 400 to 600, and a 60% increase in the number of video clips, from around 300k to around 500k. The statistics of the two dataset versions are detailed in table 1.

In the new Kinetics-600 dataset there is a standard test set, for which labels have been publicly released, and also a held-out test set (where the labels are not released). We encourage researchers to report results on the standard test set, unless they want to compare with participants of the Activity-Net kinetics challenge. Performance on the combination of standard test set plus held-out test should be used in that case, and can be measured only through the challenge evaluation website¹.

The URLs of the YouTube videos and temporal intervals of both Kinetics-600 and Kinetics-400 can be obtained from <http://deepmind.com/kinetics>.

2. Data Collection Process

The data collection process evolved from Kinetics-400 to Kinetics-600. The overall pipeline was the same: 1) action class sourcing, 2) candidate video matching, 3) candidate clip selection, 4) human verification, 5) quality analysis and filtering. In words, a list of class names is created, then a list of candidate YouTube URLs is obtained for each class name, and candidate 10s clips are sampled from the videos. These clips are sent to humans in Mechanical Turk who decide whether those clips contain the action class that they are supposed to. Finally, there is an overall curation process including clip de-duplication, and selecting the higher quality classes and clips. Full details can be found in the original publication [6].

The main differences in the data collection process between Kinetics-400 and 600 were in the first two steps: how

¹<http://activity-net.org/challenges/2018/evaluation.html>

Version	Train	Valid.	Test	Held-out Test	Total Train	Total	Classes
Kinetics-400 [6]	250–1000	50	100	0	246,245	306,245	400
Kinetics-600	450–1000	50	100	around 50	392,622	495,547	600

Table 1: Kinetics Dataset Statistics. The number of clips for each class in the various splits (left), and the totals (right). With Kinetics-600 we have released the ground truth test set labels, and also created an additional held-out test set for the purpose of the Activity-Net Challenge.

action classes were sourced, and how candidate YouTube videos were matched with classes.

2.1. Action class sourcing

For Kinetics-400, class names were first sourced from existing datasets, then from the everyday experience of the authors, and finally by asking the humans in Mechanical Turk what classes they were seeing in videos that did not contain the classes being tested. For Kinetics-600 we sourced many classes from Google’s Knowledge Graph, in particular from the hobby list. We also obtained class ideas from YouTube’s search box auto-complete, for example by typing an object or verb, then following up on promising auto-completion suggestions and checking if there were many videos containing the same action.

2.2. Candidate video matching

In Kinetics-400 we matched YouTube videos with each class by searching for videos having some of the class name words in the title, while allowing for variation in stemming. There was no separation between the class name and the query text, which turned out to be a limiting factor: in many cases we exhausted the pool of candidates, or had impractically low yields. We tried matching directly these queries to not just the title but also other metadata but this proved of little use (in particular the video descriptions seemed to have plenty of spam). We tried two variations that worked out much better:

Multiple queries. In order to get better and larger pools of candidates we found it useful to manually create sets of queries for each class and did so in two different languages: English and Portuguese. These are two out of six languages with the most native speakers in the world², have large YouTube communities (especially in the USA and Brazil), and were also natively spoken by this paper’s authors. As an example the queries for folding paper were: “folding paper” (en), “origami” (en) and “dobrar papel” (pt). We found also that translating action descriptions was not always easy, and sometimes required observing the videos returned by puta-

tive translated queries on YouTube and tuning them through some trial and error.

Having multiple languages had the positive side effect of also promoting greater dataset diversity by incorporating a more well-rounded range of cultures, ethnicities and geographies.

Weighted ngram matching. Rather than matching directly using textual queries we found it beneficial to use weighted ngram representations of the combination of the metadata of each video and the titles of related ones. Importantly, these representations were compatible with multiple languages. We combined this with standard title matching to get a robust similarity score between a query and all YouTube videos, which, unlike the binary matching we used before, meant we never ran out of candidates, although the post-mechanical-turk yield of the selected candidates became lower for smaller similarity values.

3. From Kinetics-400 to Kinetics-600

Kinetics-600 is an approximate superset of Kinetics-400 – overall, 368 of the original 400 classes are exactly the same in Kinetics-600 (except they have more examples). For the other 32 classes, we renamed a few (e.g. “dying hair” became “dyeing hair”), split or removed others that were too strongly overlapping with other classes, such as “drinking”. We split some classes: “hugging” became “hugging baby” and “hugging (not baby)”, while “opening bottle” became “opening wine bottle” and “opening bottle (not wine)”.

A few video clips from 30 classes of the Kinetics-400 validation set became part of the Kinetics-600 test set, and some from the training set became part of the new validation set. It is therefore not ideal to evaluate models on Kinetics-600 that were pre-trained on Kinetics-400, although it should make almost no difference in practice. The full list of new classes in Kinetics-600 is given in the appendix.

4. Benchmark Performance

As a baseline model we used I3D [2], with standard RGB videos as input (no optical flow). We trained the model from scratch on the Kinetics-600 training set, picked hyper-

²According to <https://www.babbel.com/en/magazine/the-10-most-spoken-languages-in-the-world/>

Acc. type	Valid	Test	Test + HeldOut Test
Top-1	71.9	71.7	69.7
Top-5	90.1	90.4	89.1
100.0 – $avg(\text{Top-1}, \text{Top-5})$	19.0	19.0	20.6

Table 2: Performance of an I3D model with RGB inputs on the Kinetics-600 dataset, without any test time augmentation (processing a center crop of each video convolutionally in time). The first two rows show accuracy in percentage, the last one shows the metric used at the Kinetics challenge hosted by the ActivityNet workshop.

parameters on validation, and report performance on validation, test set and the combination of the test and held-out test sets. We used 32 P100 GPUs, batch size 5 videos, 64 frame clips for training and 251 frames for testing. We trained using SGD with momentum, starting with a learning rate of 0.1, decreasing it by a factor of 10 when the loss saturates. Results are shown in table 2.

The top-1 accuracy on the test set was 71.7, whereas on Test+Held-out was 69.7, which shows that the held-out test set is harder than the regular test set. On Kinetics-400 the corresponding result was 68.4, hence the task overall seems to have become slightly easier. There are several factors that may help explain this: even though Kinetics-600 has 50% extra classes, it also has around 50% extra training examples; and also, some of the ambiguities in Kinetics-400 have been removed in Kinetics-600. We also used fewer GPUs (32 instead 64), which resulted in half the batch size.

Kinetics challenge. There was a first Kinetics challenge at the ActivityNet workshop in CVPR 2017, using Kinetics-400. The second challenge occurred at the ActivityNet workshop in CVPR 2018, this time using Kinetics-600. The performance criterion used in the challenge is the average of Top-1 and Top-5 error. There was an improvement between the winning systems of the two challenges, with error going down from 12.4% (in 2017) to 11.0% (in 2018) [1, 4].

5. Conclusion

We have described the new Kinetics-600 dataset, which is 50% larger than the original Kinetics-400 dataset. It represents another step towards our goal of producing an action classification dataset with 1000 classes and 1000 video clips for each class. We explained the differences in the data collection process between the initial version of the dataset made available in 2017 and the new one. This publication coincides with the release of the test set annotations for both Kinetics-400 and Kinetics-600; we hope these will facilitate research as it will no longer be necessary to submit results to an external evaluation server.

Acknowledgements:

The collection of this dataset was funded by DeepMind. The authors would like to thank Sandra Portugues for help-

ing to translate queries from English to Portuguese, and Aditya Zisserman and Radhika Desikan for data clean up.

References

- [1] Y. Bian, C. Gan, X. Liu, F. Li, X. Long, Y. Li, H. Qi, J. Zhou, S. Wen, and Y. Lin. Revisiting the effectiveness of off-the-shelf temporal modeling approaches for large-scale video classification. *arXiv preprint arXiv:1708.03805*, 2017.
- [2] 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.
- [3] C. Gu, C. Sun, D. A. Ross, C. Vondrick, C. Pantofaru, Y. Li, S. Vijayanarasimhan, G. Toderici, S. Ricco, R. Sukthankar, et al. AVA: A video dataset of spatio-temporally localized atomic visual actions. *CoRR, abs/1705.08421*, 4, 2017.
- [4] D. He, F. Li, Q. Zhao, X. Long, Y. Fu, and S. Wen. Exploiting spatial-temporal modelling and multi-modal fusion for human action recognition. *arXiv preprint arXiv:1806.10319*, 2018.
- [5] Y. Jiang, J. Liu, A. R. Zamir, G. Toderici, I. Laptev, M. Shah, and R. Sukthankar. Thumos challenge: Action recognition with a large number of classes, 2014.
- [6] W. Kay, J. Carreira, K. Simonyan, B. Zhang, C. Hillier, S. Vijayanarasimhan, F. Viola, T. Green, T. Back, P. Natsev, M. Suleyman, and A. Zisserman. The kinetics human action video dataset. *arXiv preprint arXiv:1705.06950*, 2017.
- [7] 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.
- [8] G. A. Sigurdsson, G. Varol, X. Wang, A. Farhadi, I. Laptev, and A. Gupta. Hollywood in homes: Crowdsourcing data collection for activity understanding. In *European Conference on Computer Vision*, pages 510–526. Springer, 2016.
- [9] 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.

A. List of New Human Action Classes in Kinetics-600

This is the list of classes in Kinetics-600 that were not in Kinetics-400, or that have been renamed.

1. acting in play
2. adjusting glasses

3. alligator wrestling
4. archaeological excavation
5. arguing
6. assembling bicycle
7. attending conference
8. backflip (human)
9. base jumping
10. bathing dog
11. battle rope training
12. blowdrying hair
13. blowing bubble gum
14. bodysurfing
15. bottling
16. bouncing on bouncy castle
17. breaking boards
18. breathing fire
19. building lego
20. building sandcastle
21. bull fighting
22. bulldozing
23. burping
24. calculating
25. calligraphy
26. capsizing
27. card stacking
28. card throwing
29. carving ice
30. casting fishing line
31. changing gear in car
32. changing wheel (not on bike)
33. chewing gum
34. chiseling stone
35. chiseling wood
36. chopping meat
37. chopping vegetables
38. clam digging
39. coloring in
40. combing hair
41. contorting
42. cooking sausages (not on barbeque)

43. cooking scallops
44. cosplaying
45. cracking back
46. cracking knuckles
47. crossing eyes
48. cumbia
49. curling (sport)
50. cutting apple
51. cutting orange
52. delivering mail
53. directing traffic
54. docking boat
55. doing jigsaw puzzle
56. drooling
57. dumpster diving
58. dyeing eyebrows
59. dyeing hair
60. embroidering
61. falling off bike
62. falling off chair
63. fencing (sport)
64. fidgeting
65. fixing bicycle
66. flint knapping
67. fly tying
68. geocaching
69. getting a piercing
70. gold panning
71. gospel singing in church
72. hand washing clothes
73. head stand
74. historical reenactment
75. home roasting coffee
76. huddling
77. hugging (not baby)
78. hugging baby
79. ice swimming
80. inflating balloons
81. installing carpet
82. ironing hair

83. jaywalking
84. jumping bicycle
85. jumping jacks
86. karaoke
87. land sailing
88. lawn mower racing
89. laying concrete
90. laying stone
91. laying tiles
92. leatherworking
93. licking
94. lifting hat
95. lighting fire
96. lock picking
97. longboarding
98. looking at phone
99. luge
100. making balloon shapes
101. making bubbles
102. making cheese
103. making horseshoes
104. making paper aeroplanes
105. making the bed
106. marriage proposal
107. massaging neck
108. moon walking
109. mosh pit dancing
110. mountain climber (exercise)
111. mushroom foraging
112. needle felting
113. opening bottle (not wine)
114. opening door
115. opening refrigerator
116. opening wine bottle
117. packing
118. passing american football (not in game)
119. passing soccer ball
120. person collecting garbage
121. photobombing
122. photocopying
123. pillow fight
124. pinching
125. pirouetting
126. planing wood
127. playing beer pong
128. playing blackjack
129. playing darts
130. playing dominoes
131. playing field hockey
132. playing gong
133. playing hand clapping games
134. playing laser tag
135. playing lute
136. playing maracas
137. playing marbles
138. playing netball
139. playing ocarina
140. playing pan pipes
141. playing pinball
142. playing ping pong
143. playing polo
144. playing rubiks cube
145. playing scrabble
146. playing with trains
147. poking bellybutton
148. polishing metal
149. popping balloons
150. pouring beer
151. preparing salad
152. pushing wheelbarrow
153. putting in contact lenses
154. putting on eyeliner
155. putting on foundation
156. putting on lipstick
157. putting on mascara
158. putting on sari
159. putting on shoes
160. raising eyebrows
161. repairing puncture
162. riding snow blower

163. roasting marshmallows
164. roasting pig
165. rolling pastry
166. rope pushdown
167. sausage making
168. sawing wood
169. scrapbooking
170. scrubbing face
171. separating eggs
172. sewing
173. shaping bread dough
174. shining flashlight
175. shopping
176. shucking oysters
177. shuffling feet
178. sipping cup
179. skiing mono
180. skipping stone
181. sleeping
182. smashing
183. smelling feet
184. smoking pipe
185. spelunking
186. square dancing
187. standing on hands
188. staring
189. steer roping
190. sucking lolly
191. swimming front crawl
192. swinging baseball bat
193. sword swallowing
194. tackling
195. tagging graffiti
196. talking on cell phone
197. tasting wine
198. threading needle
199. throwing ball (not baseball or American football)
200. throwing knife
201. throwing snowballs
202. throwing tantrum
203. throwing water balloon
204. tie dying
205. tightrope walking
206. tiptoeing
207. trimming shrubs
208. twiddling fingers
209. tying necktie
210. tying shoe laces
211. using a microscope
212. using a paint roller
213. using a power drill
214. using a sledge hammer
215. using a wrench
216. using atm
217. using bagging machine
218. using circular saw
219. using inhaler
220. using puppets
221. vacuuming floor
222. visiting the zoo
223. wading through mud
224. wading through water
225. waking up
226. walking through snow
227. watching tv
228. waving hand
229. weaving fabric
230. winking
231. wood burning (art)
232. yarn spinning