

# Supplementary Material: Joint Inference of Groups, Events and Human Roles in Aerial Videos

Tianmin Shu<sup>1</sup>, Dan Xie<sup>1</sup>, Brandon Rothrock<sup>2</sup>, Sinisa Todorovic<sup>3</sup> and Song-Chun Zhu<sup>1</sup>

<sup>1</sup> Center for Vision, Cognition, Learning and Art, University of California, Los Angeles

{stm512, xiedan}@g.ucla.edu sczhu@stat.ucla.edu

<sup>2</sup>Jet Propulsion Laboratory, California Institute of Technology

brandon.rothrock@jpl.nasa.gov

<sup>3</sup>School of Electrical Engineering and Computer Science, Oregon State University

sinisa@onid.orst.edu

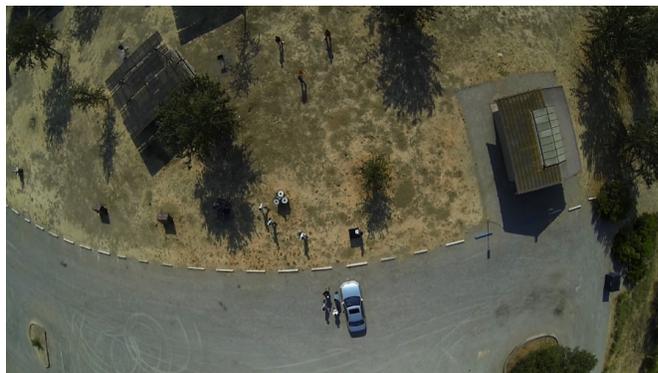
## 1. Dataset



Figure 1: Image of our hex-rotor in the air with a GoPro camera.

We assembled a new low-cost hex-rotor with a GoPro camera shown in Fig. 1, which is able to eliminate the high frequency vibration of the camera and hold in air autonomously through a GPS and a barometer. It can also fly 20 ~ 90m above the ground and stays 5 minutes in air. We use this hex-rotor to take a set of videos with some plots at a park where the terrain is interesting: hiking routes, parking lots, camping sites, picnic areas with shelters, restrooms, tables, trash bins and BBQ ovens. By detecting/tracking humans and objects in the videos, we can recognize events such as BBQ, queuing, exchanging objects, loading/unloading, etc.

We have collected some events with scripts involving the interactions between humans and objects. Fig. 2 shows two frames captured from the original videos. The original videos are pre-processed, including camera calibration and frame registration. After pre-processing, there are totally 27 videos in the dataset, the length of which ranges from 2 minutes to 5 minutes. We annotate the hierarchical semantic information of objects, roles, events and groups in the



(a)



(b)

Figure 2: Two frames of the original aerial videos from two different sites.

videos. Tab. 1 is a summary of events, roles, objects and the number of instances in our dataset.

ID	Event	Objects	Roles	# of instances
1	<i>Exchange Box</i>	<i>Box, Car</i>	<i>Deliverer, Receiver</i>	11
2	<i>Play Frisbee</i>	<i>Frisbee</i>	<i>Player</i>	13
3	<i>Info Consult</i>	<i>Desk, Info Booth</i>	<i>Consultant, Visitor</i>	11
4	<i>Pick Up</i>	<i>Car</i>	<i>Driver, Passenger</i>	9
5	<i>Queue for Vending Machine</i>	<i>Vending Machine</i>	<i>Queuing Person</i>	9
6	<i>Group Tour</i>	N/A	<i>Guide, Tourist</i>	6
7	<i>Throw Trash</i>	<i>Trash Bin</i>	<i>Thrower</i>	12
8	<i>Sit on Table</i>	<i>Table &amp; Seat</i>	<i>Customer</i>	10
9	<i>Pinic</i>	<i>Blanket</i>	<i>Picnic Person</i>	4
10	<i>Serve Table</i>	<i>Table &amp; Seat</i>	<i>Waiter, Customer</i>	6
11	<i>Sell BBQ</i>	<i>BBQ Oven</i>	<i>Chef, Buyer</i>	6

Table 1: Summary of the dataset.

Object type	Precision	Recall
Buildings	100.0%	95.65%
Cars	10.38%	30.64%
Small static objects	16.16%	53.33%

Table 2: Object detection accuracy. Small static objects include *Table & Seat*, *Info Booth*, *Desk*, *BBQ Oven*, and *Trash Bin*. We have accurate building detection results while the detection accuracy for small static objects and cars is not very ideal, which affects the recognition of events that involve these objects.

Precision	Recall	Rate of broken tracklets
31.36%	39.05%	4.44

Table 3: Tracking accuracy. Tracking generates a large number of false alarms and fails when humans/objects have little motions, which can be indicated by the low precision and recall. The rate of broken tracklets shows that ground truth trajectories are broken into multiple fragments in the tracking results. The poor tracking results greatly increase the difficulty of our inference.

## 2. Detection and Tracking Evaluation

The accuracy of the detection of buildings, cars and small static objects is shown in Tab. 2.

Tab. 3 shows the accuracy of tracking.