

# Supplementary Material: Learning to Detect Human-Object Interactions

Yu-Wei Chao<sup>1</sup>, Yunfan Liu<sup>1</sup>, Xieyang Liu<sup>1</sup>, Huayi Zeng<sup>2</sup>, and Jia Deng<sup>1</sup>

<sup>1</sup>University of Michigan, Ann Arbor

{ywchao, yunfan, lxieyang, jiadeng}@umich.edu

<sup>2</sup>Washington University in St. Louis\*

{zengh}@wustl.edu

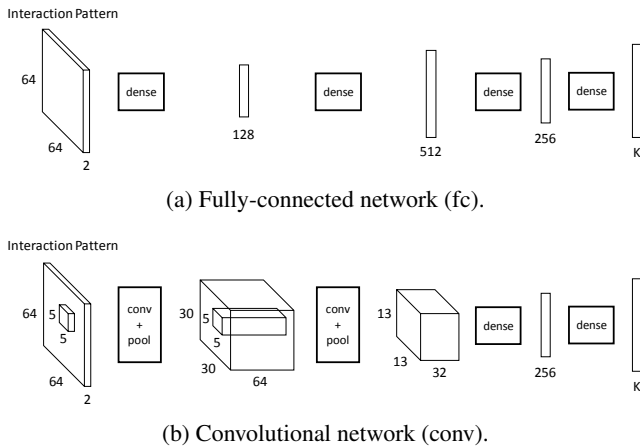


Figure 1: Two different architectures for the pairwise stream.

## 1. Pairwise Stream

In our HO-RCNN, we consider two different network architectures for the pairwise stream: a fully-connected network (fc) and a convolutional network (conv). The two architectures are illustrated in Fig. 1. Both architectures take as input an Interaction Pattern and produce as output a vector of classification scores on  $K$  HOI classes of interest. For fair comparison, both networks have approximately the same number of parameters, and are trained with identical schemes.

\*Work done at the University of Michigan as a visiting student.