

## Weakly-supervised Visual Grounding of Phrases with Linguistic Structures Fanyi Xiao, Leonid Sigal\*, and Yong Jae Lee **University of California, Davis** \*Disney Research

#### Produce visual groundings of linguistic Goal phrases at all levels (words, phrases, sentences), by training with weak supervision (img-cap pairs).

a man

sandwich







## **Previous works** erges from a vellow colapsable toy tunnel onto th

• Treat the caption as a sequence of word tokens [*Xu 2014, Rohrbach 2016, …*].

• Require phrase-region pairs or pretrained detectors [Karpathy 2015, Plummer 2015, ...].

## **Motivation & Key ideas**

• Hard, if not impossible, to annotate a large collection of phrase-segment pairs. • A sentence is not simply a sequence of token. We leverage structure in natural language for

weakly supervised visual grounding.

#### **Benefits of exploiting** linguistic structures

 Avoid grounding nonsensical tokens ("a", "the", etc.)

• Enrich training data, in a linguistically sound way (i.e., words, phrases, sentence) • Transfer linguistic structure to visual domain







# **Network architecture**

Loss functions









0.3	IOU@0.4	IOU@0.5	Avg mAP
2	0.199	0.110	0.203
7	0.213	0.118	0.219
.6	0.203	0.114	0.211
4	0.240	0.138	0.238
7	0.246	0.159	0.251