

SEEING WHAT YOU'RE TOLD: SENTENCE-GUIDED ACTIVITY RECOGNITION IN VIDEO

N. Siddharth, Andrei Barbu, & Jeffrey Mark Siskind

Stanford University, MIT, Purdue University

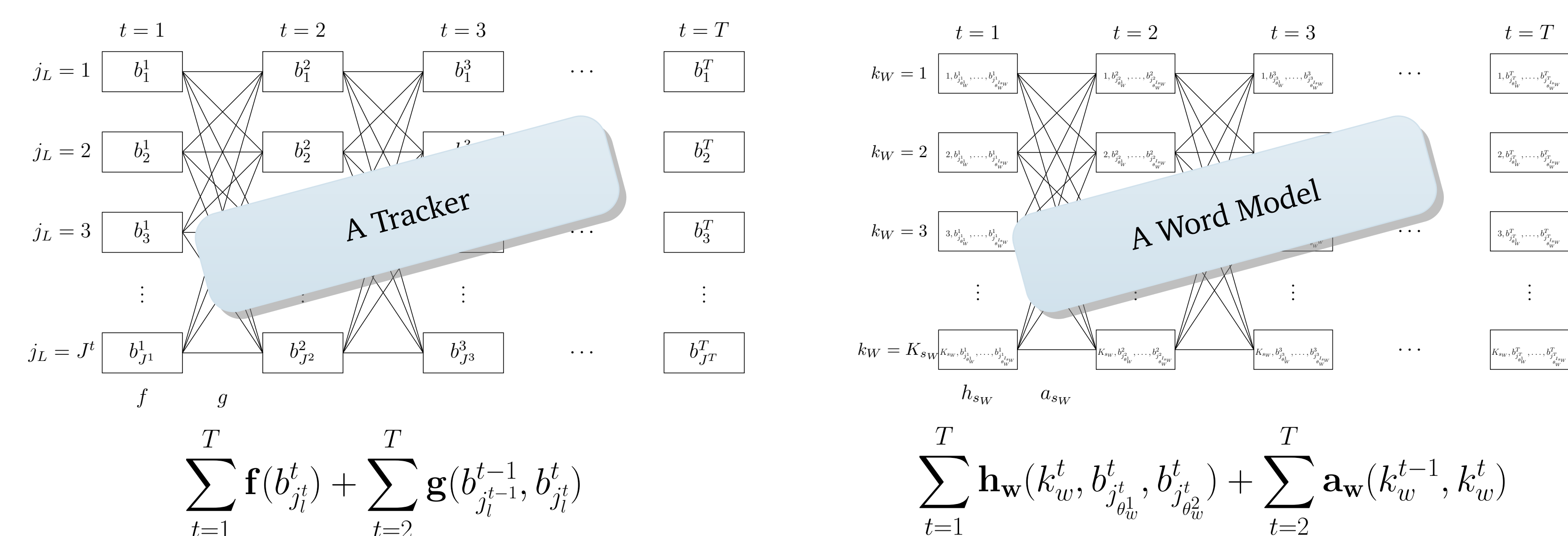
What?



The person to the left of the chair put down the blue object.

Recognize

How?

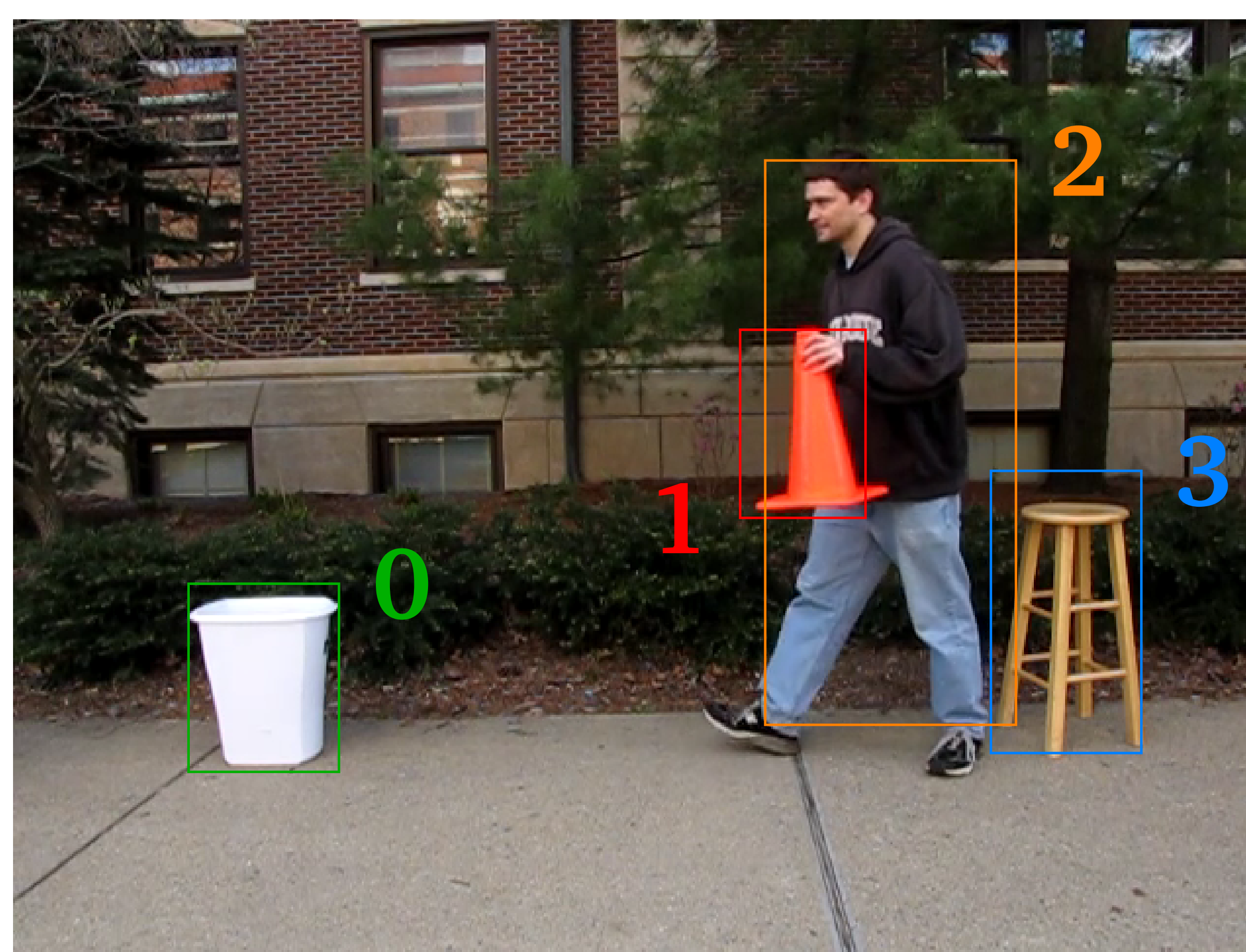
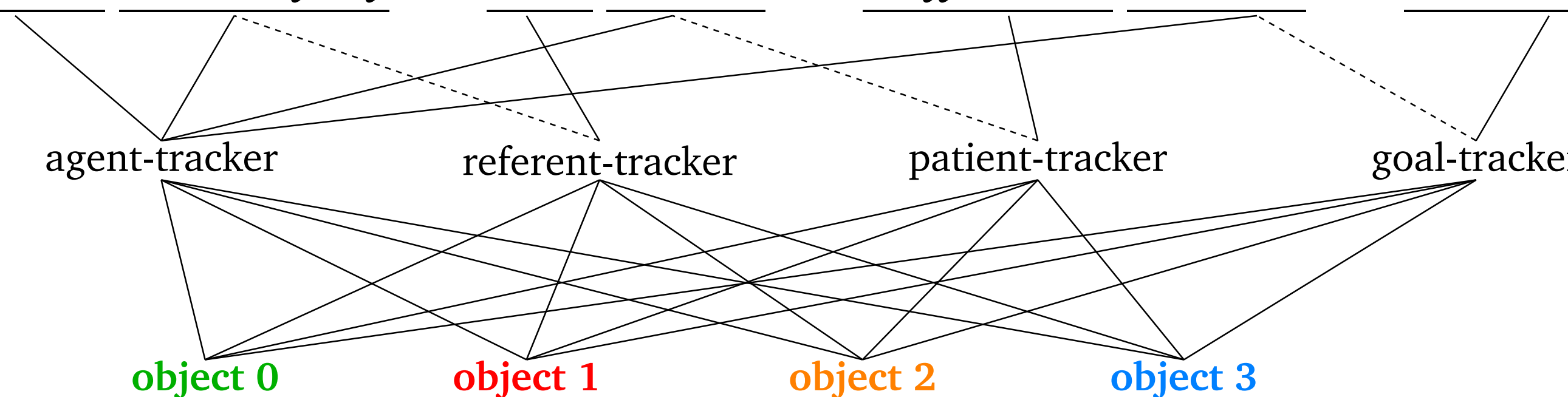


Joint

$$\max_{j_l^1, \dots, j_l^T, k_w^1, \dots, k_w^T} \sum_{l=1}^L \sum_{t=1}^T f(b_{j_l^t}^t) + \sum_{t=2}^T g(b_{j_l^{t-1}}^{t-1}, b_{j_l^t}^t) + \sum_{w=1}^W \sum_{t=1}^T h_w(k_w^t, b_{j_{a_l^t}}^t, b_{j_{p_l^t}}^t) + \sum_{t=2}^T a_w(k_w^{t-1}, k_w^t)$$

L tracks W words

The person to the left of the stool carried the traffic-cone towards the trash-can.



$$\mathcal{S} : (\mathbf{B}, \mathbf{s}, \Lambda) \rightarrow (\tau, \mathbf{J})$$

\mathbf{B} – detections, \mathbf{s} – sentence, Λ – lexicon

τ – score, \mathbf{J} – tracks

Look!

Performing disparate tasks simply by leveraging the framework differently

Sentential Focus of Attention

Same video | Different sentences → Different tracks



The person carried an object towards the trash can. The person carried an object away from the trash can.

Semantic Video Retrieval

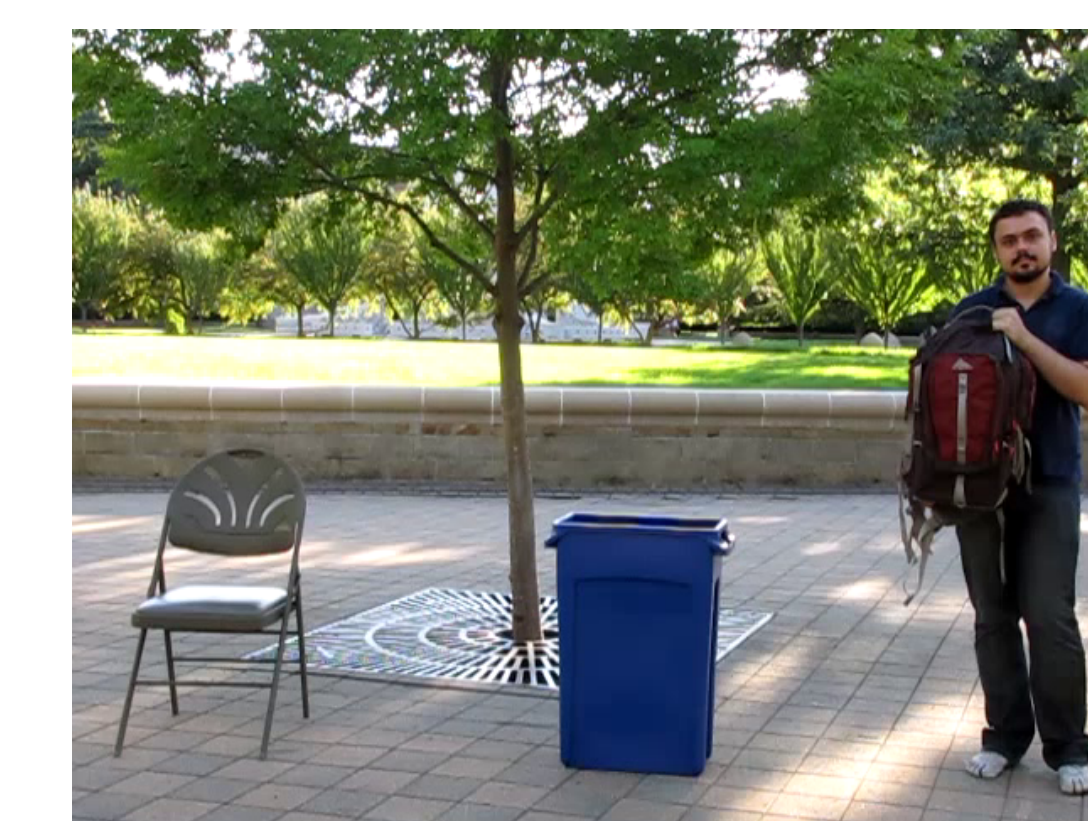
Many videos | A sentence → That video



The person carried the backpack away from the trash can.

Generation of Sentential Descriptions

Many sentences[†] | A video → That sentence

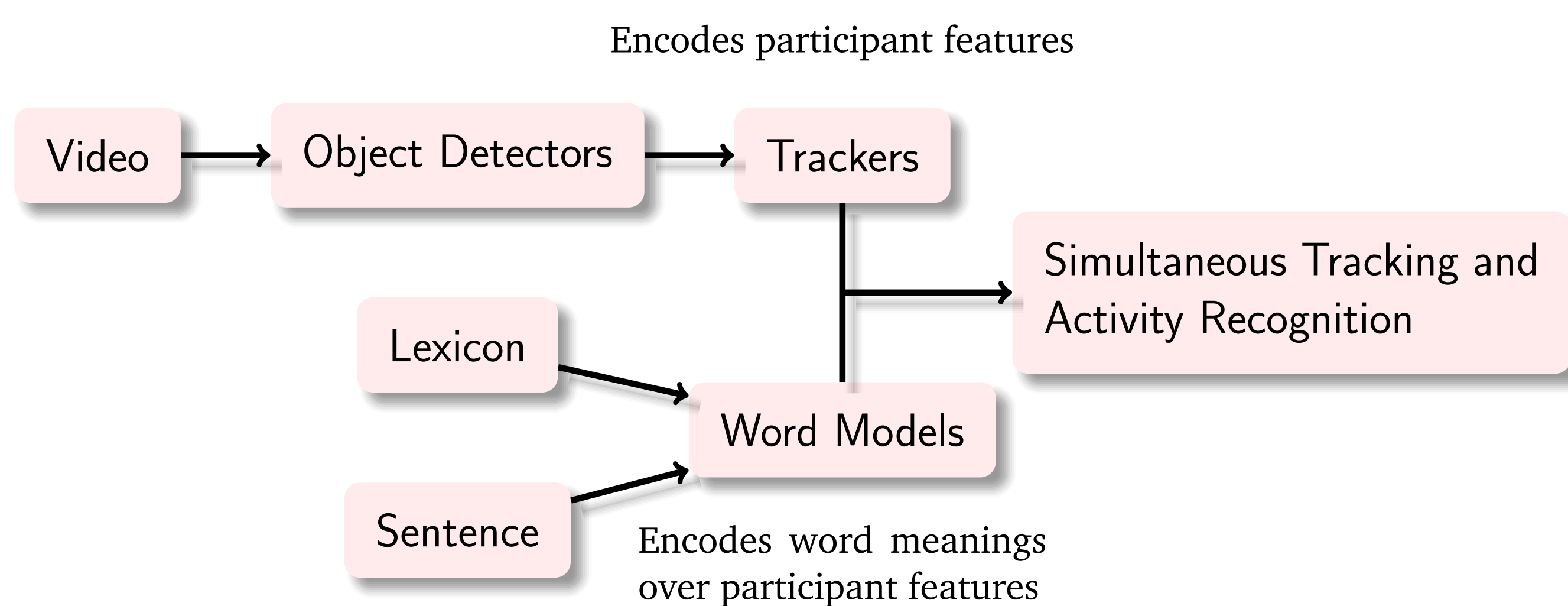


The person to the right of the trash-can picked up the backpack.

†

- S → NP VP
- NP → D [A] N [PP]
- D → an | the
- A → blue | red
- N → person | backpack | trash-can | chair | object
- PP → P NP
- P → to the left of | to the right of
- VP → V NP [ADV] [PPM]
- V → picked up | put down | carried | approached
- ADV → quickly | slowly
- PPM → PM NP
- PM → towards | away from

Characteristics



- Joint evaluation of both Tracking and Activity Recognition
 - A Tracker for each participant in the activity
 - A Word Model for each lexical entry in the lexicon

- Integration of top-down sentential information and bottom-up tracker information

- Recognize different parts of speech

| | |
|-------------|--------------|
| Verbs | Adverbs |
| Nouns | Adjectives |
| Determiners | Prepositions |

- Sensitive to sentence structure

The person approached the object. \neq The object approached the person.

This research was supported, in part, by ARL, under Cooperative Agreement Number W911NF-10-2-0060, and the Center for Brains, Minds and Machines, funded by NSF STC award CCF-1231216. The views and conclusions contained in this document are those of the authors and do not represent the official policies, either express or implied, of ARL or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for Government purposes, notwithstanding any copyright notation herein.