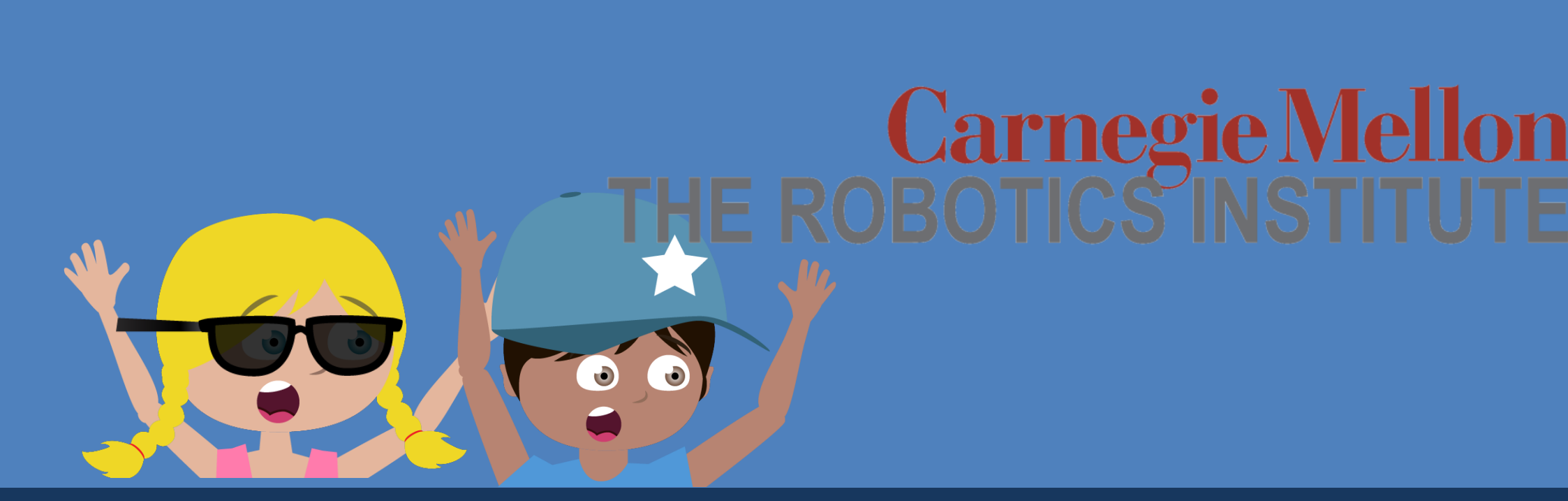




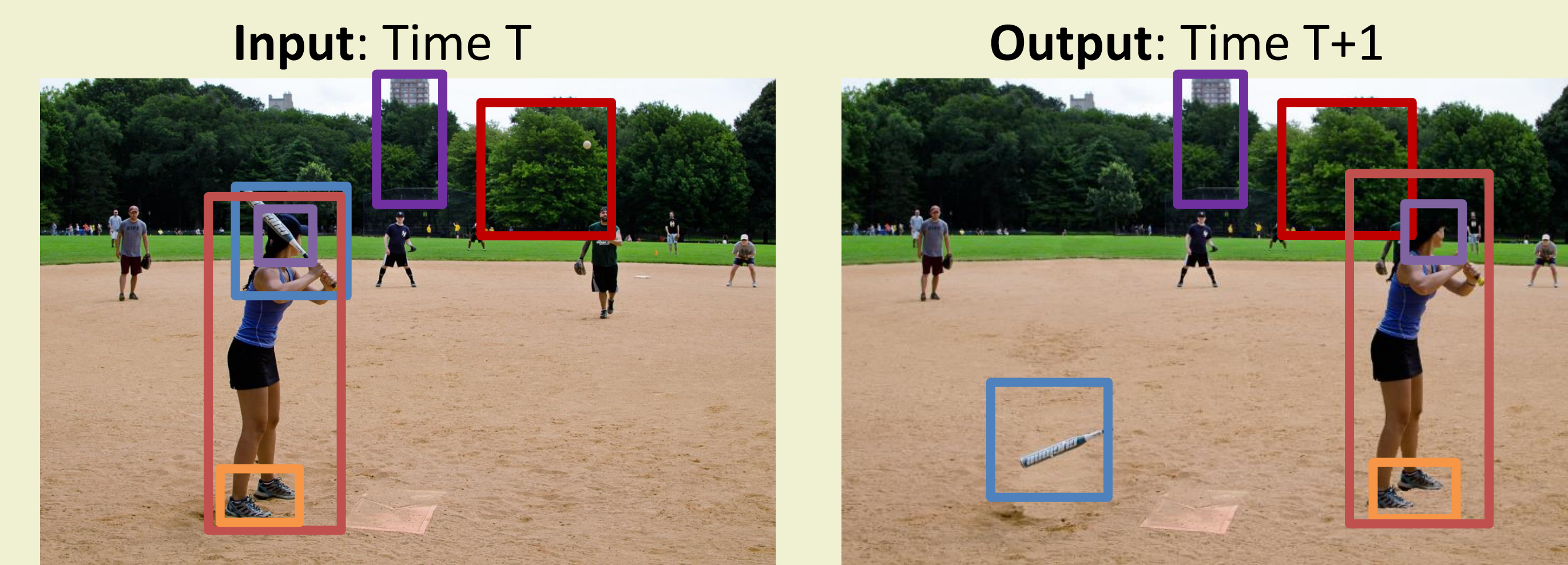
# Predicting Object Dynamics in Scenes

David F. Fouhey<sup>1</sup>, Larry Zitnick<sup>2</sup>

<sup>1</sup>Robotics Institute, Carnegie Mellon University, <sup>2</sup>Microsoft Research



## Overview



**Goal:** Learn mapping  
**Input:** bounding boxes at time T  
**Output:** bounding boxes at time T+1

**How can we learn spatiotemporal common-sense to predict that:**

- The woman will probably move right?
- The woman's hat will probably go with her?
- The trees will stay still?

**Problems:**

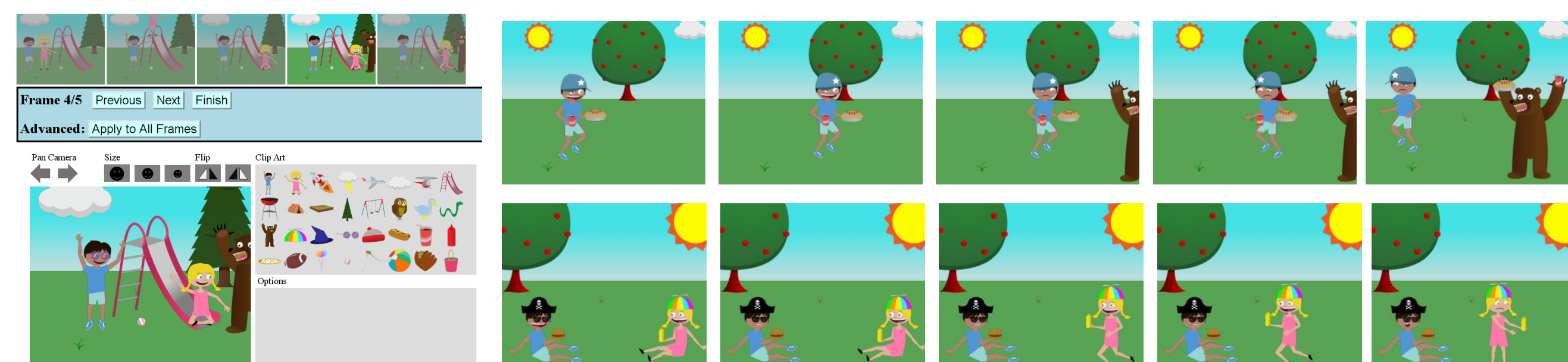
- How can we model and satisfy these constraints?
- How can we do this if we cannot detect people and hats reliably?

## Datasets

Prediction models were trained on Abstract Scenes Data and applied to both abstract and natural scenes.

### Abstract Scenes Data:

5,000 sequences of 5 scenes  
 Gathered via Amazon Mechanical Turk (AMT)



Interface

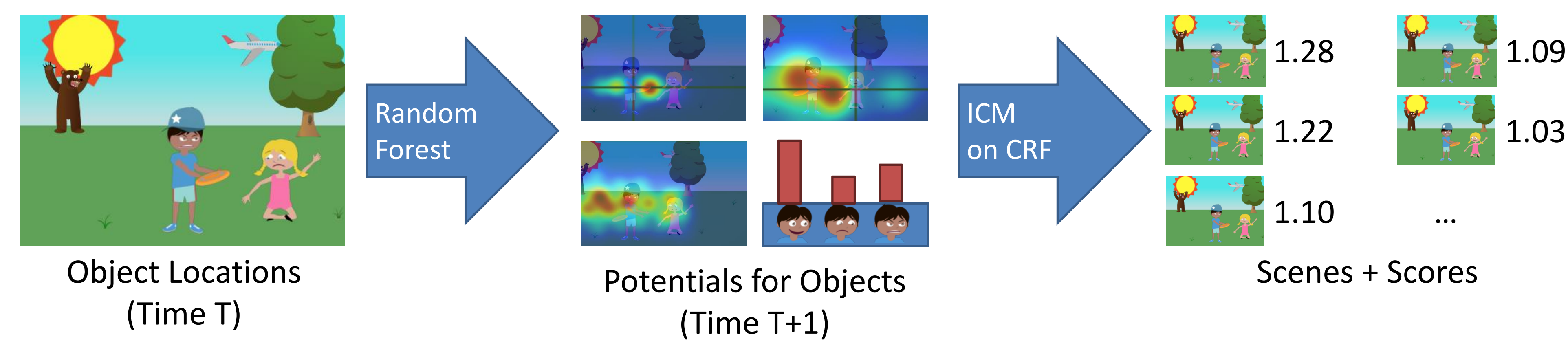
Sample Stories

### Natural Images:

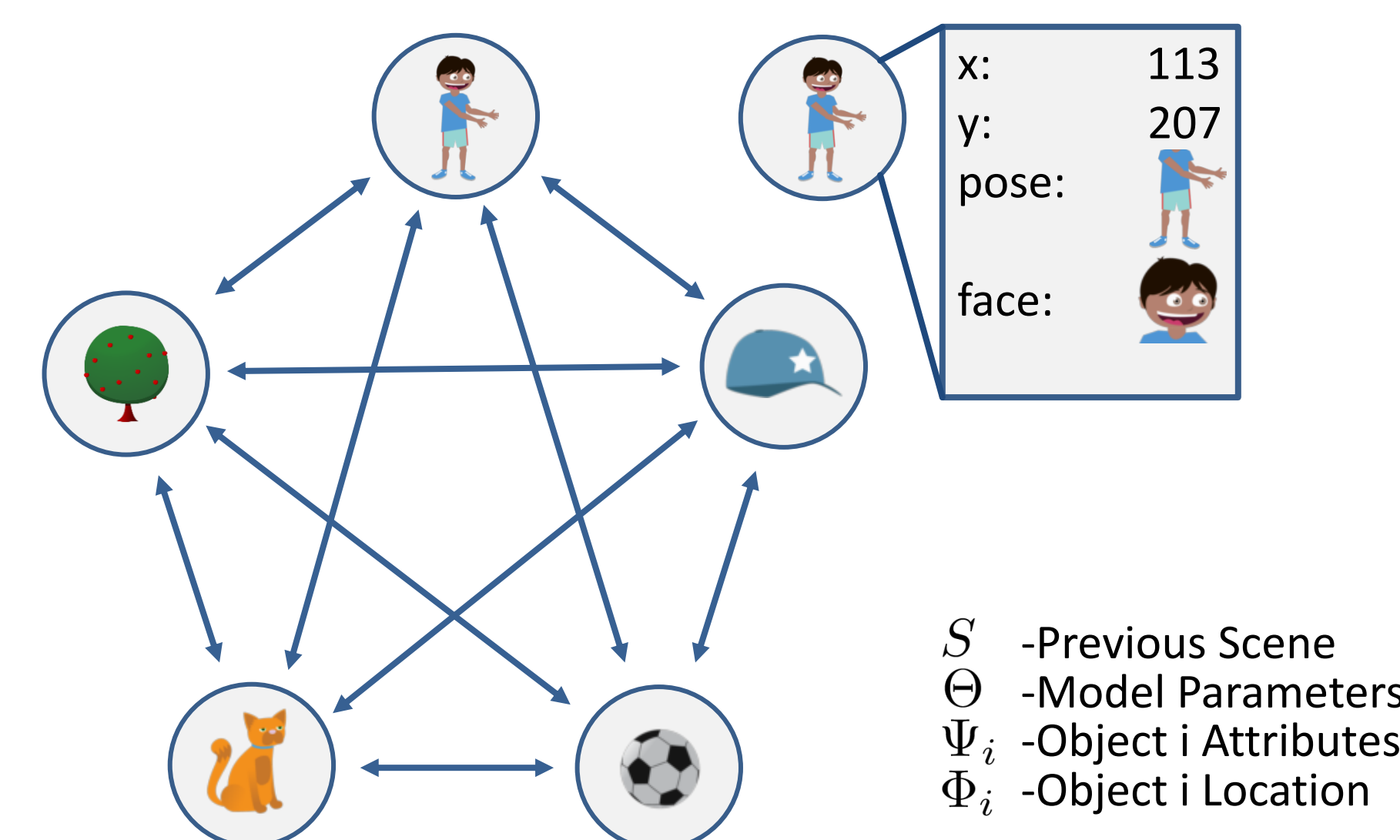
225 natural images from flickr.com with bounding-boxes labeled  
 Only one image; ground-truth predictions labeled by AMT



## Method

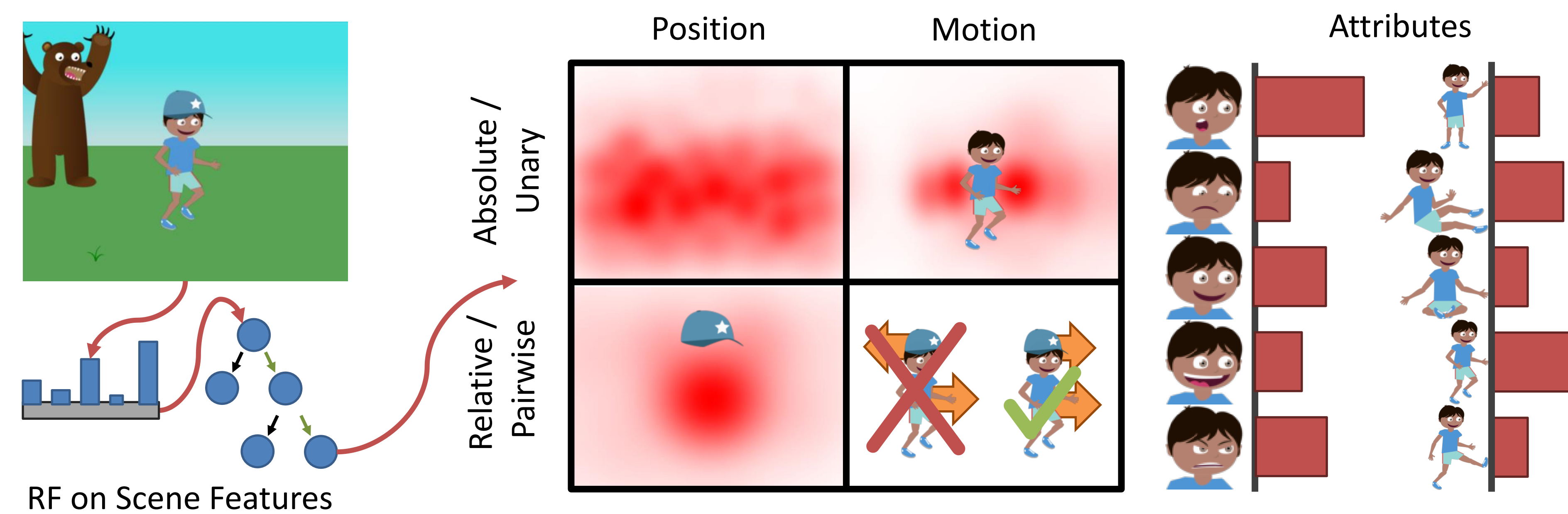


### CRF Scene Model



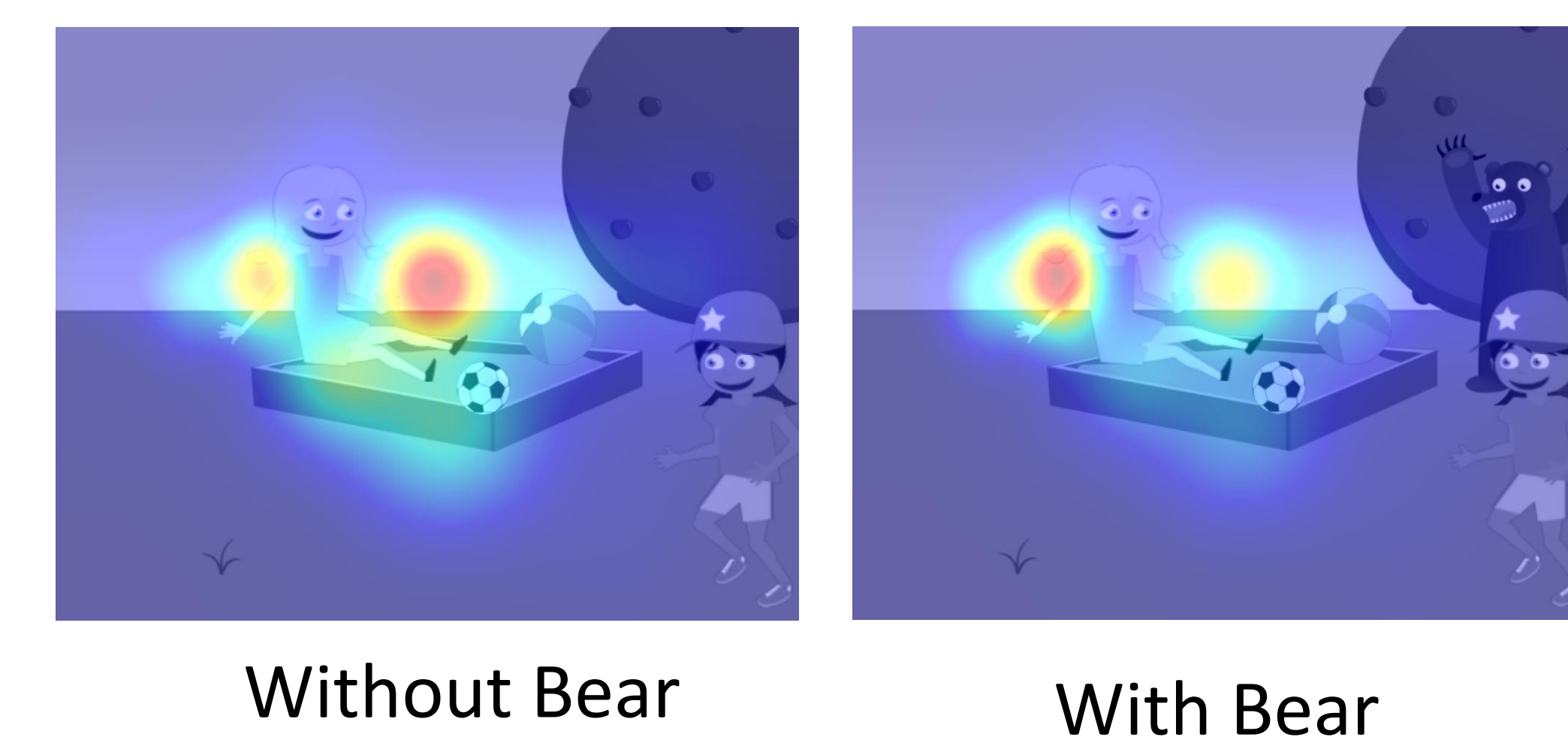
$$\log P(\Phi, \Psi | S, \Theta) = \sum_i \left( \underbrace{\lambda(\Phi_i; \theta_\lambda)}_{abs. location} + \underbrace{\omega_i(\Phi_i, S; \theta_\omega)}_{abs. motion} + \underbrace{\pi_i(\Psi_i, S; \theta_\pi)}_{attributes} \right) + \sum_{i,j} \left( \underbrace{\phi_{i,j}(\Phi_i, \Phi_j; \theta_\phi)}_{rel. location} + \underbrace{\varphi_{i,j}(\Phi_i, \Phi_j, S; \theta_\varphi)}_{rel. motion} \right) + \underbrace{\alpha(\Phi, S; \Theta_\alpha)}_{motion prior} + Z(S, \Theta)$$

### CRF Potentials learned by Random Forest

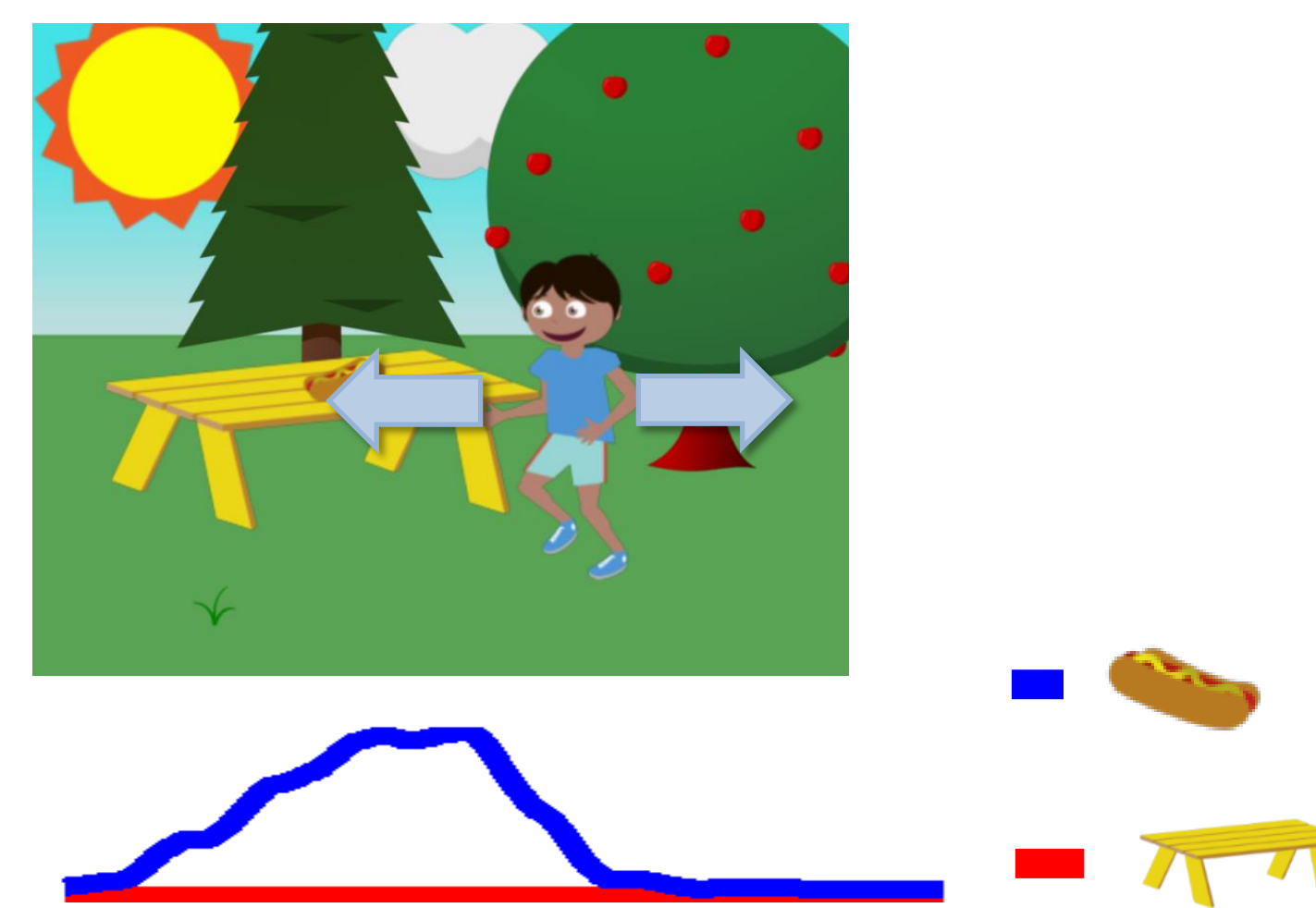


### Learned Motion Potentials

Distribution on Jenny's motion vector in the same scene with and without a bear present.



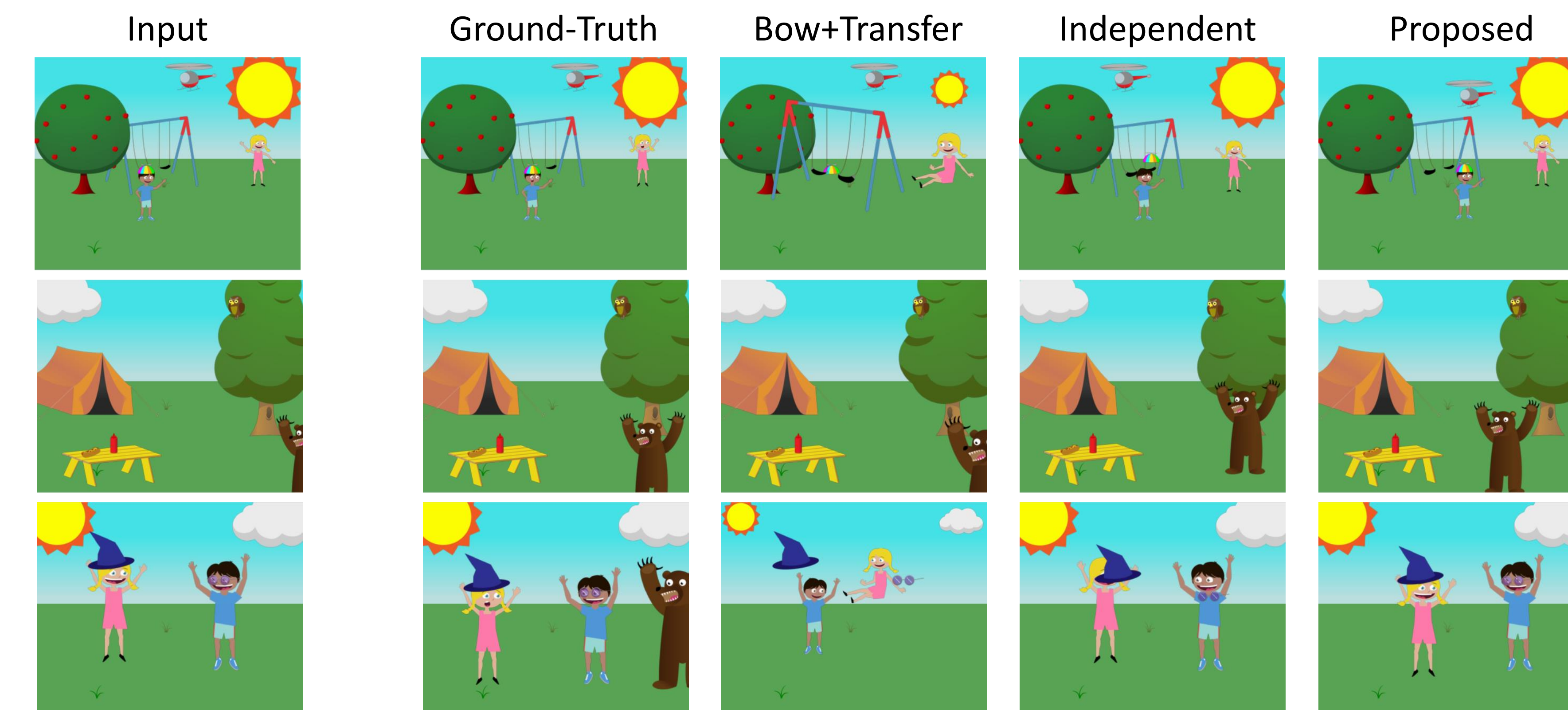
P(motion) as a function of Mike's position



## Results

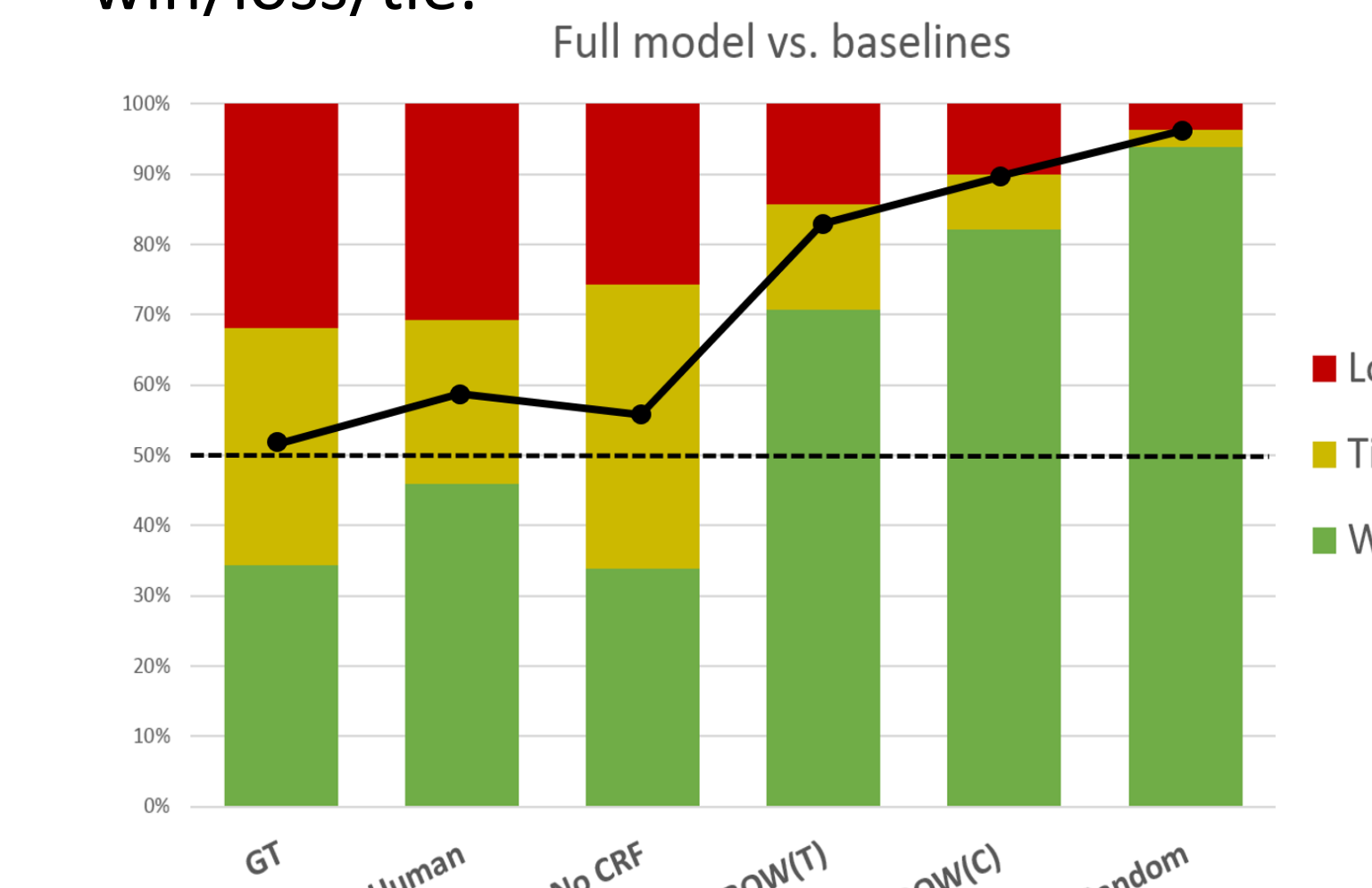
Compare with:

**Global Prediction:** BoW + Motion Vector Transfer;  
**Independent Prediction:** Use RF on same features;  
**Humans:** Ask another Turker to complete.



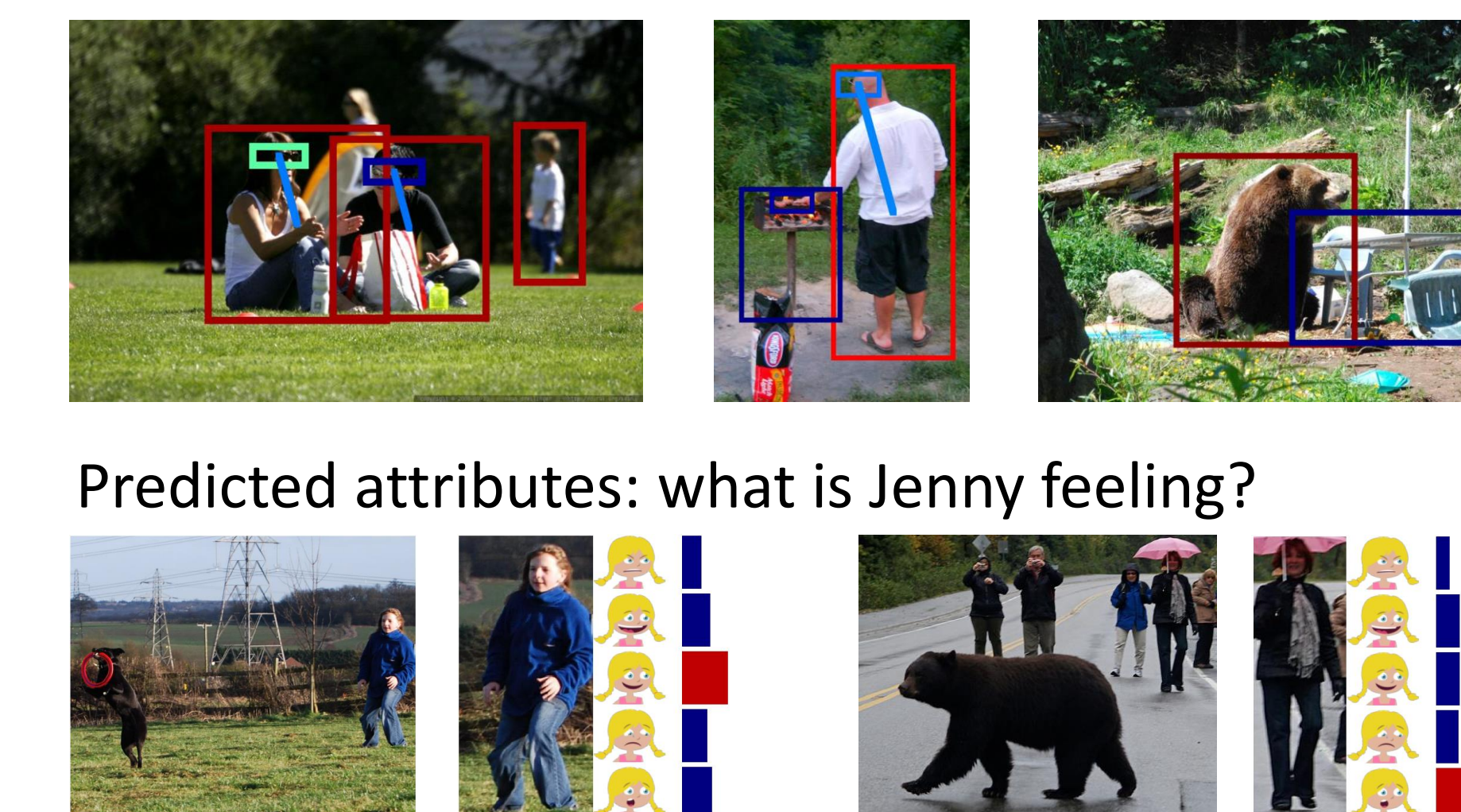
### Human Evaluation

Another human rates pairs of predictions win/loss/tie.



### Natural Image Results

Who moves and how?  
 Hot color: likely to move. Line: moves together.



### Quantitative Evaluation

Ask yes/no question about scene; compare original sequence and prediction.

F1 Score:	Abstract Scenes Dataset				Natural Scenes Dataset		
	Does X move?	Do X and Y move together?	Does X move Left or Right?	What are X's Attributes?	Does X move?	Do X and Y move together?	
Proposed	<b>49.3</b>	<b>42.9</b>	66.6	<b>61.9</b>	Proposed	91.2	45.2
Indep.	49.0	11.9	<b>75.1</b>	61.5	Indep.	87.5	14.4
Humans	48.8	31.8	70.0	52.3	Humans	<b>97.8</b>	<b>75.3</b>
Global	39.6	16.2	61.9	29.1	Global	69.3	5.7

## Conclusions

- Neither global nor independent prediction models produce accurate joint motion: joint motion must be modeled explicitly.
- Our models are short-term common-sense only; long-term prediction with narratives is an interesting future direction.