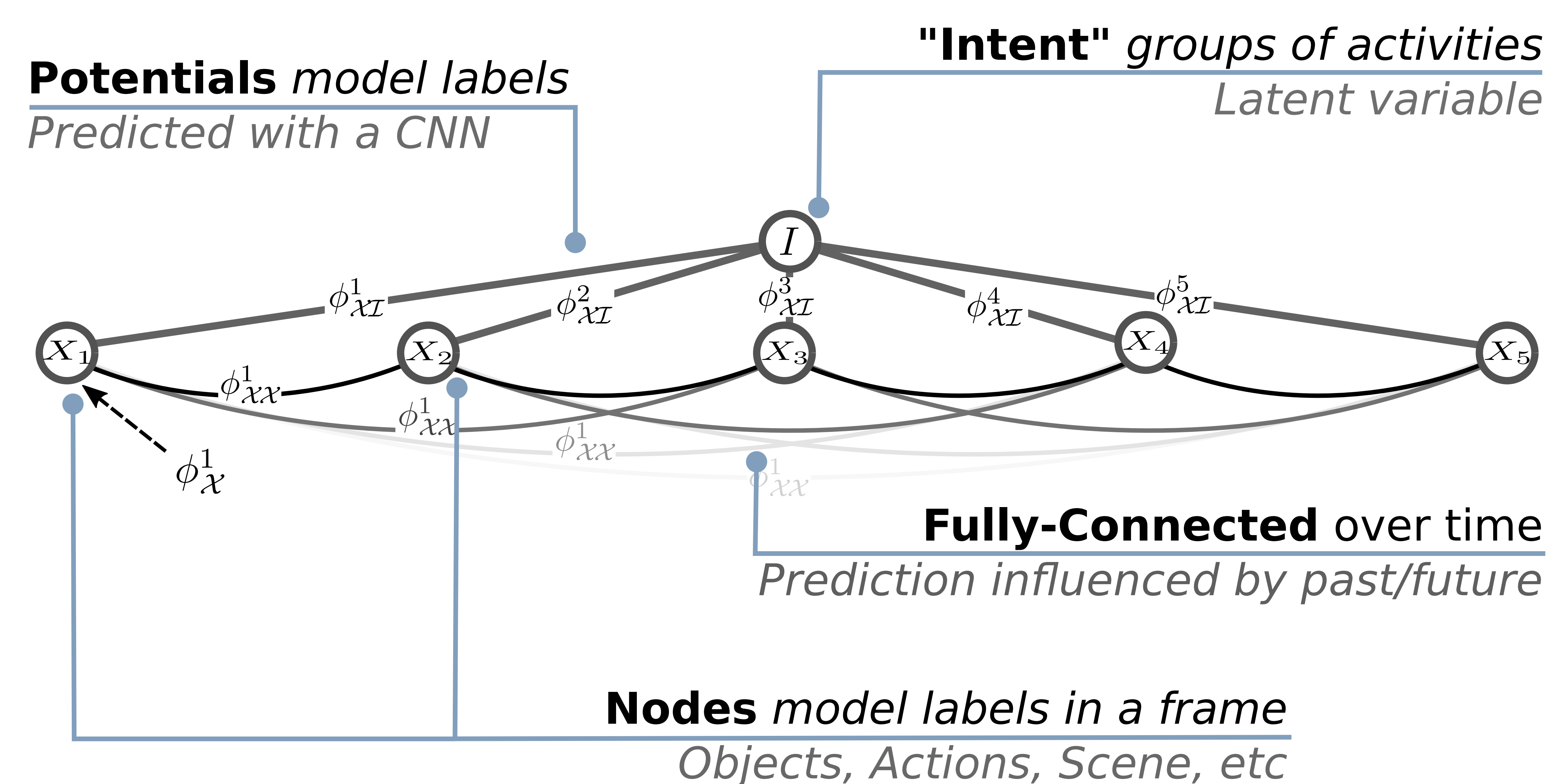


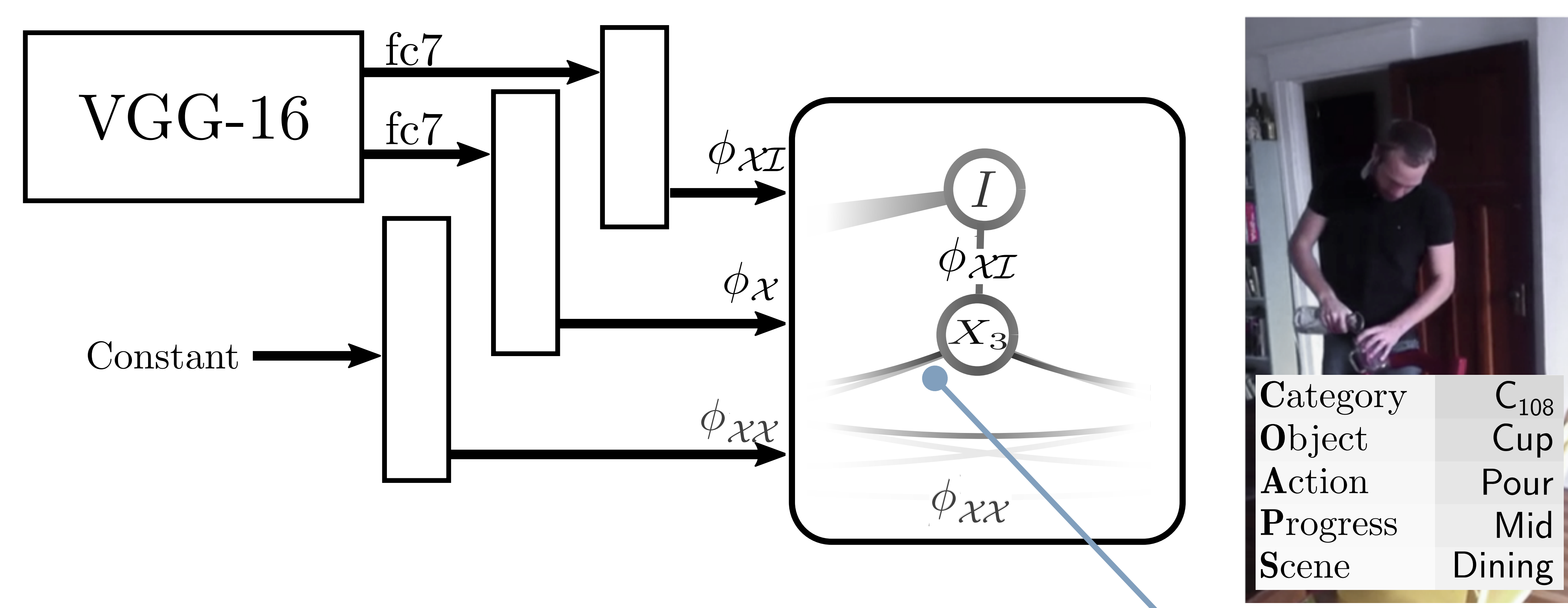
Model

- We model activities as a fully-connected CRF over time
- Nodes for each frame in the video, plus "intent"

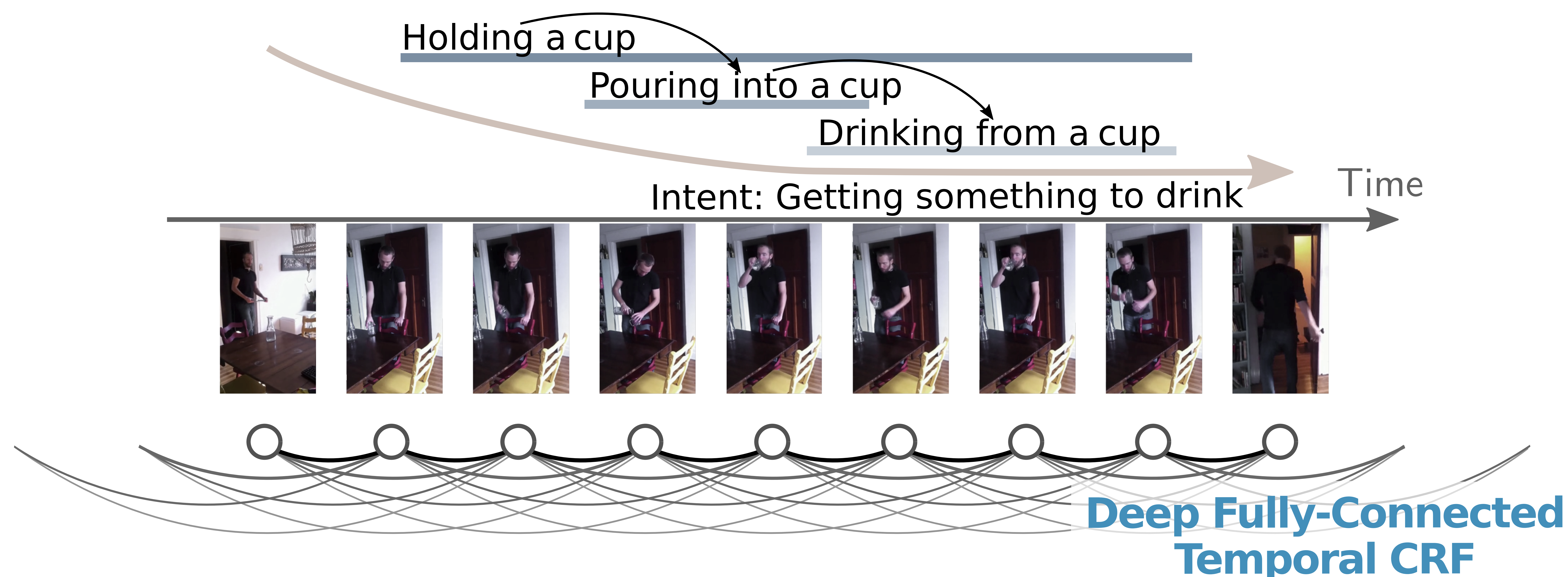


Zooming in...

- A Two-Stream network predicts the potentials



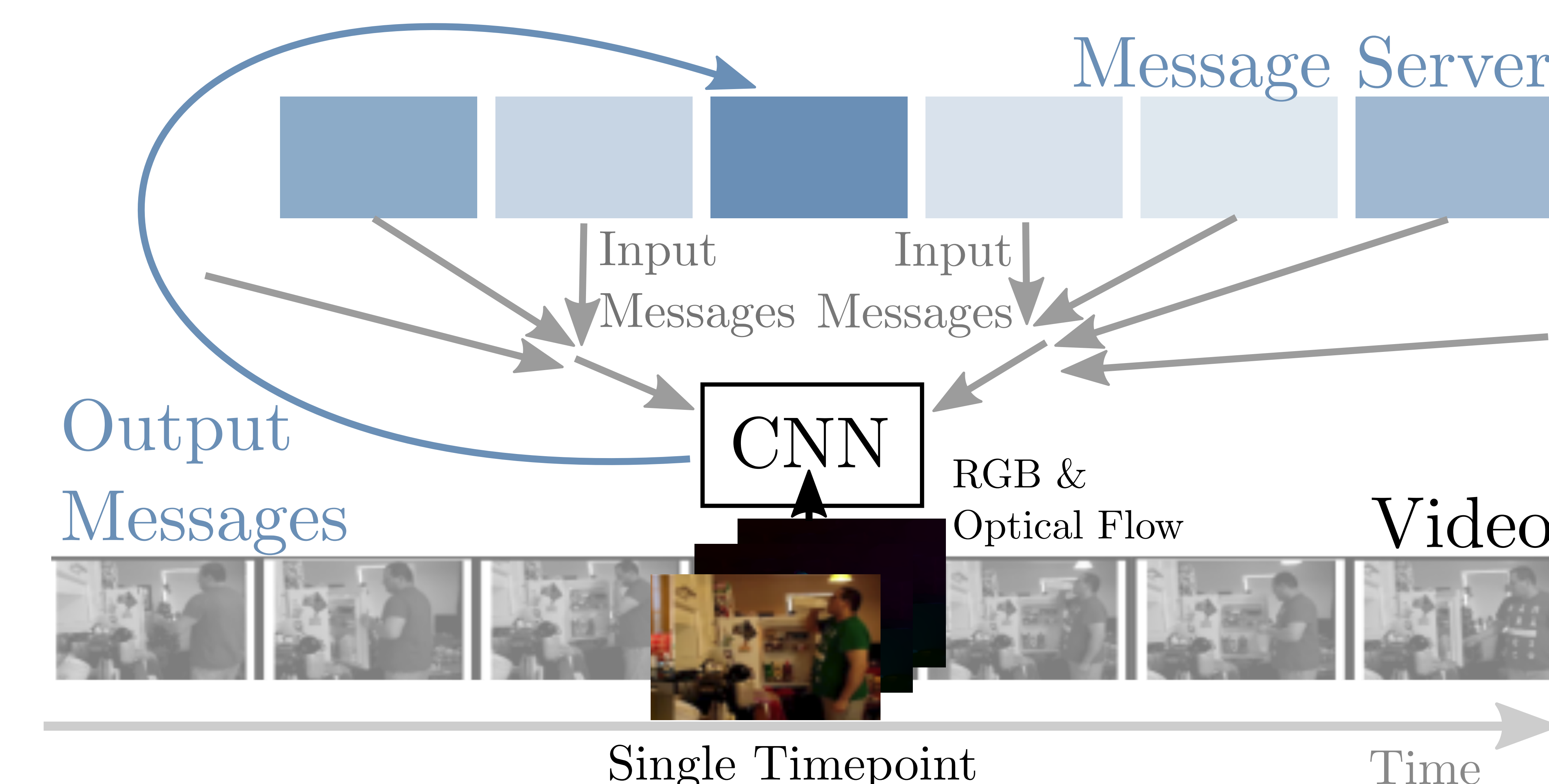
Understanding Sequences of Activities



High-order Temporal Model Trained From Randomly Sampled Individual Images

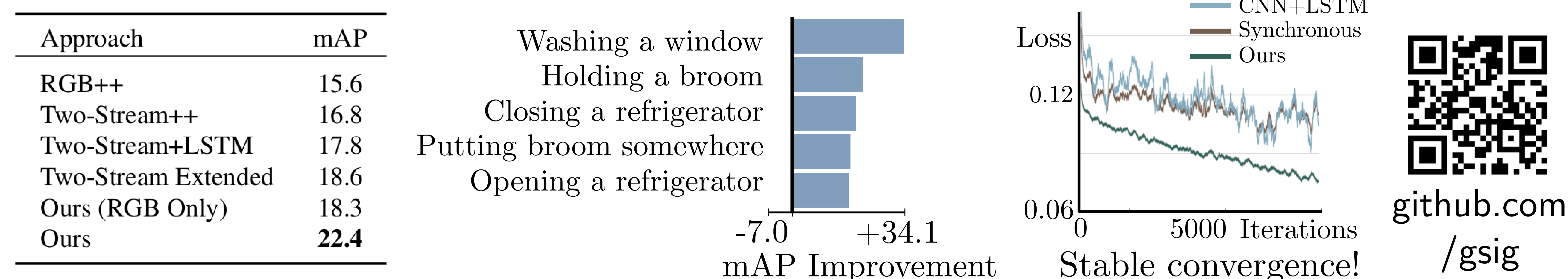
Training Method

- With *Mean-Field*, training is simple:
 1. Sample any frame
 2. Get stored "messages"
 3. Backprop using just this frame



Results

- Evaluation on the Charades dataset
157 classes, 66k instances, object, verb, scene, captions, etc



Learned "Intent"

- We model a latent variable for groups of activities
- Learns higher-level activity concepts:

