

Summary

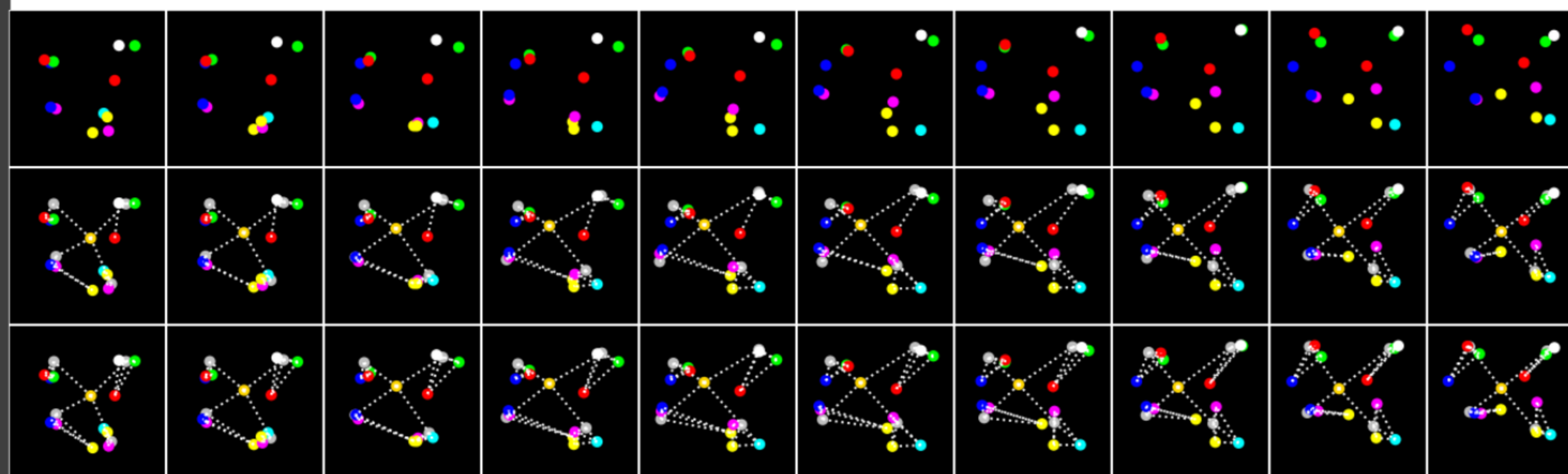
- Common-sense physical reasoning, requires learning about **objects**, their interactions and dynamics directly from raw visual input
- Real world objects vary greatly in terms of their properties, which complicates modelling their dynamics
- Many objects can be viewed as a hierarchy of parts that locally behave independently, but also act more globally as a single whole
- We learn a **hierarchical world model (HRI)** that explicitly distinguishes objects at multiple levels of abstraction and relations between them
- Use cases: interacting with complex objects (e.g. humanoid robot), or learning a hierarchically structured predictive model for MBRL

Experimental Questions

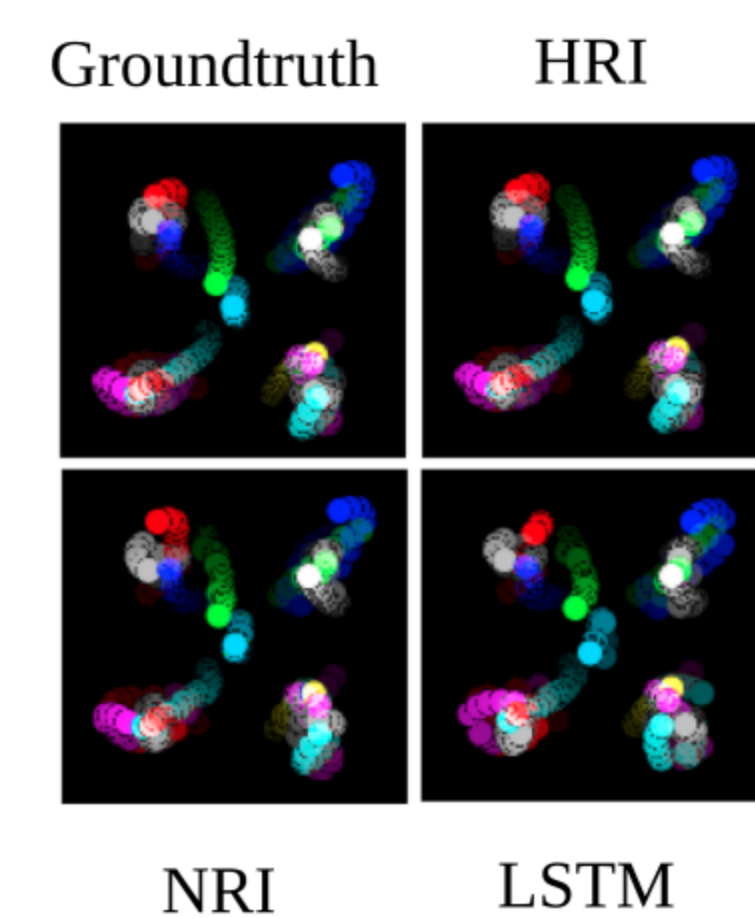
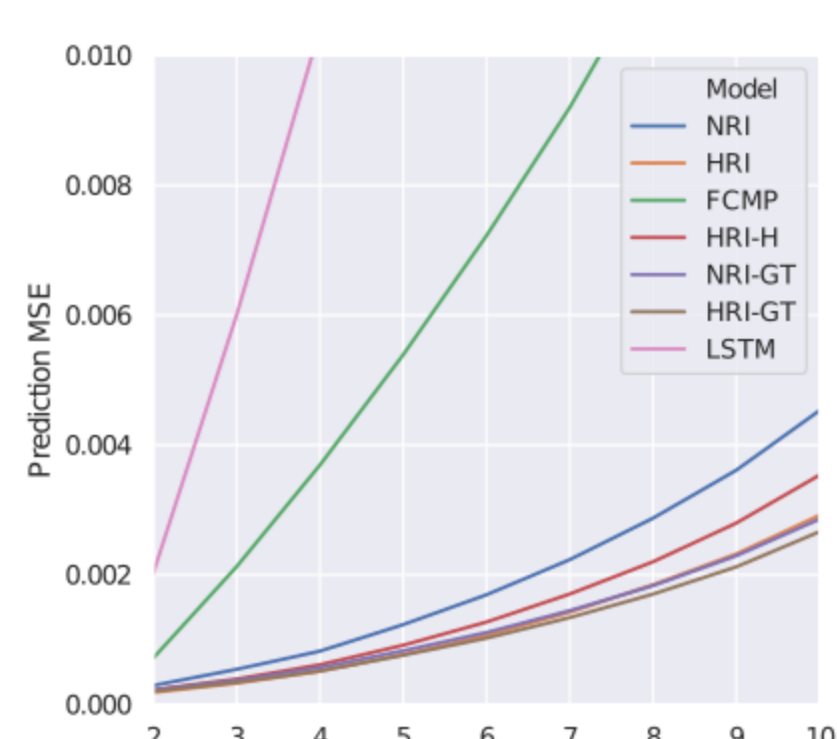
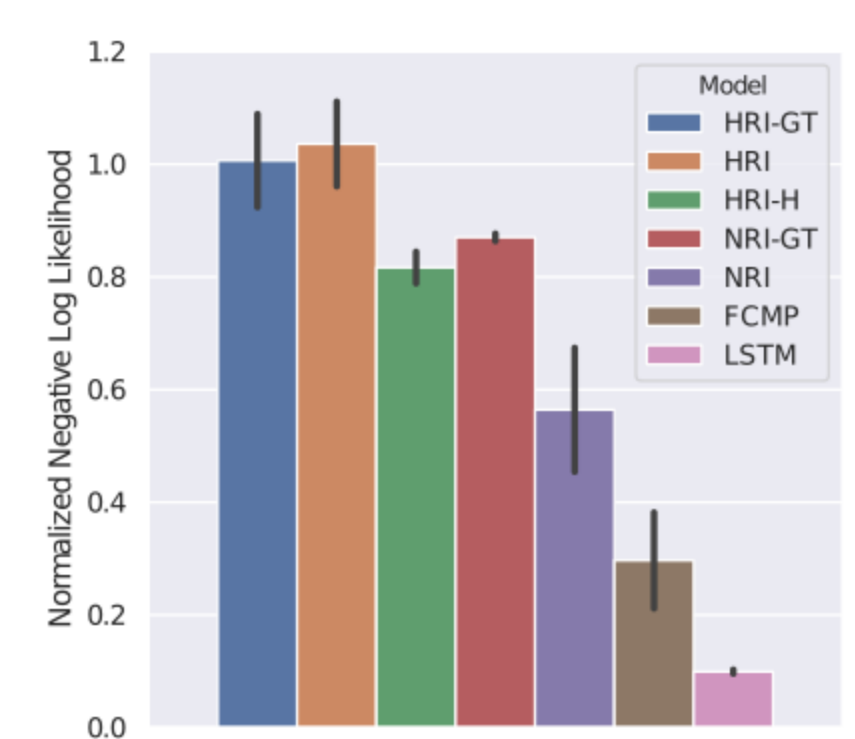
- Can parts, objects and relations be inferred from visual data?
- What is the effect of relational inference?
- Does hierarchical graph inference and reasoning help?
- Can HRI be scaled to real world visual scenarios?

State Springs

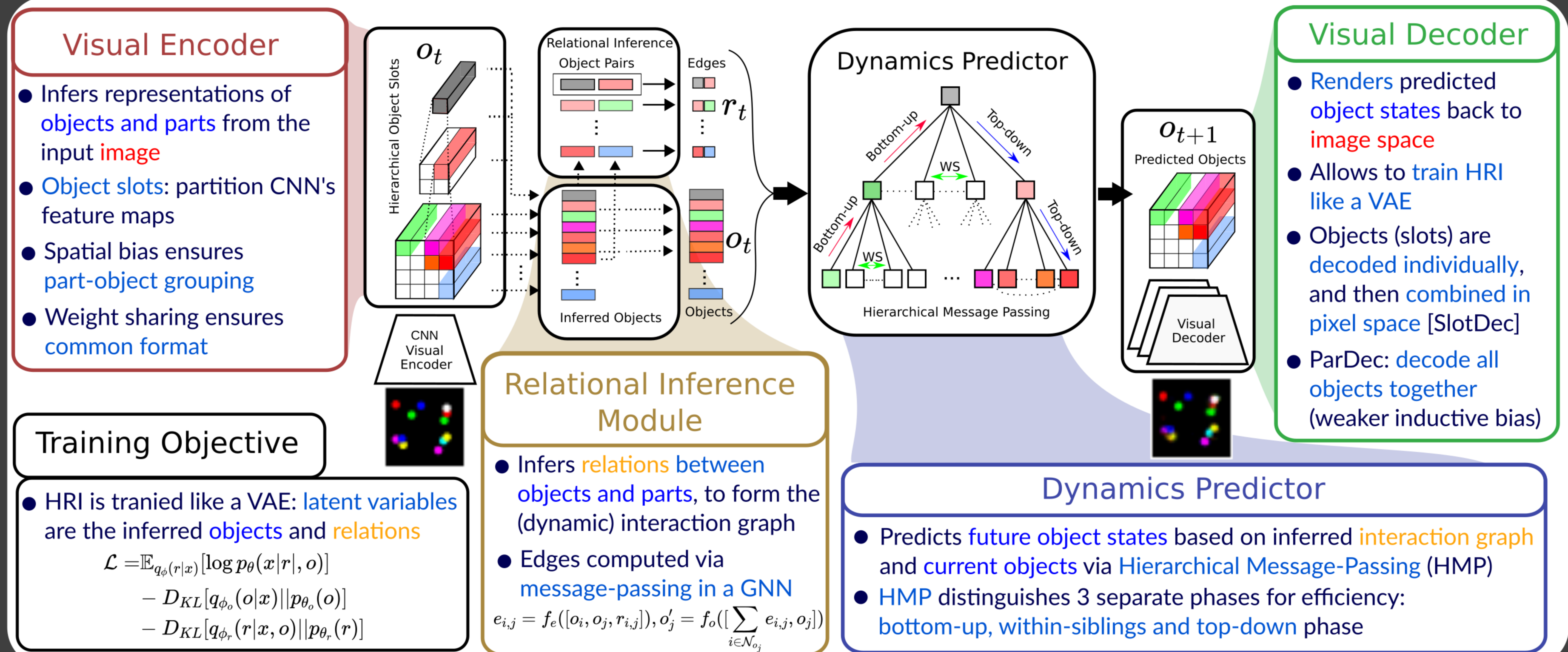
- Synthetic dataset of state trajectories of objects and parts connected via finite-length springs in a hierarchical structure
- HRI makes accurate predictions about the future state of the environment and correctly infers the underlying interaction graph



- Comparisons to baselines show benefits of relational inference (HRI vs LSTM) and hierarchical message passing (HRI vs NRI)

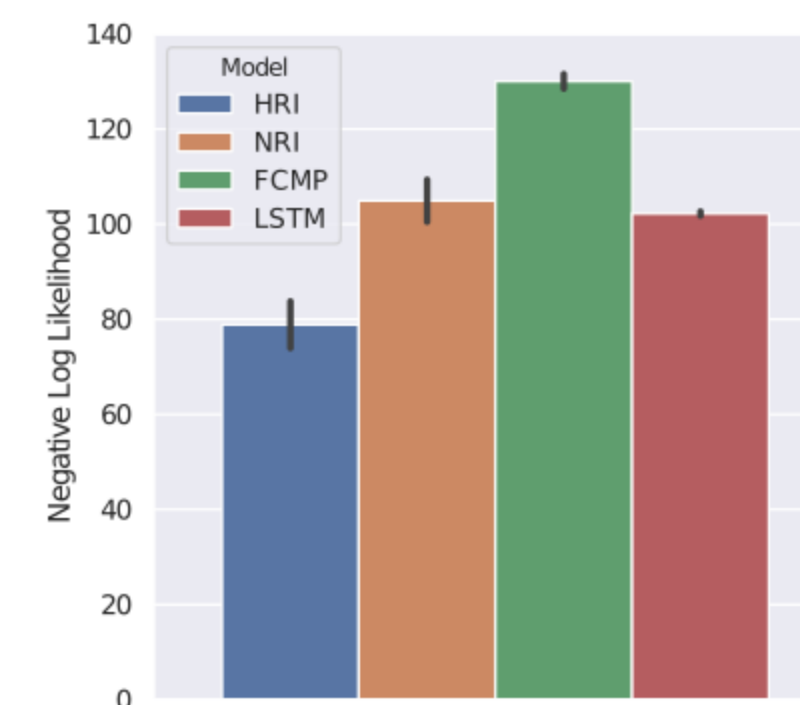
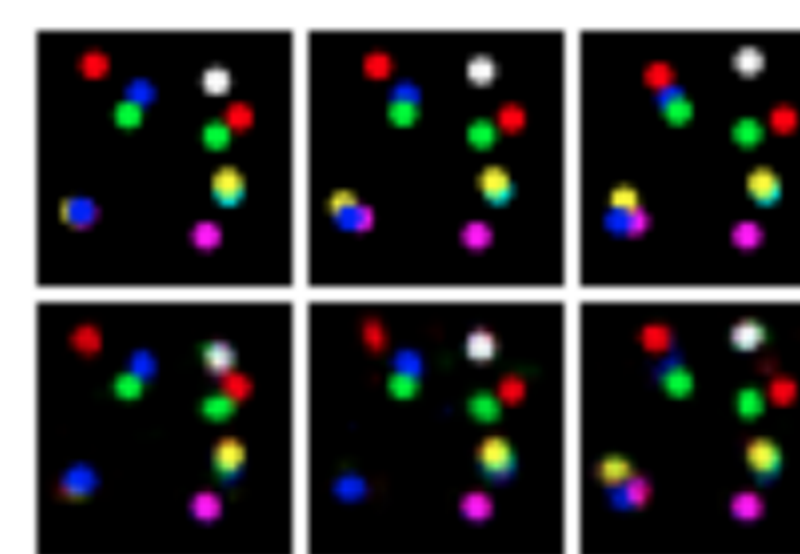


HRI model

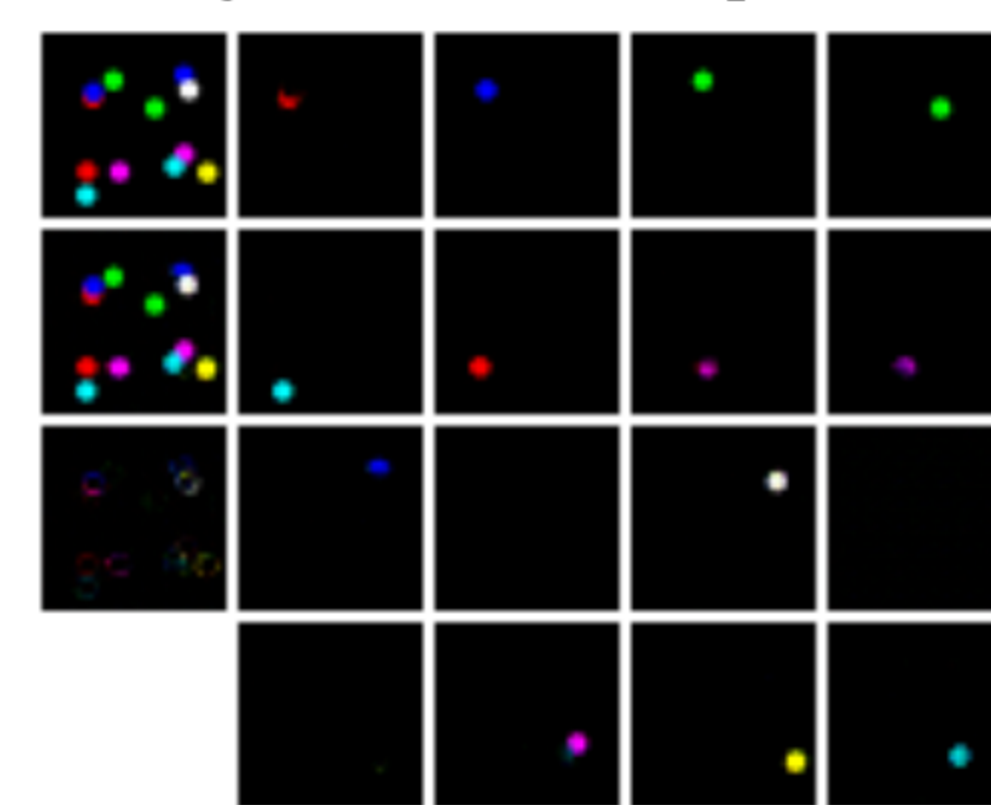


Visual Springs

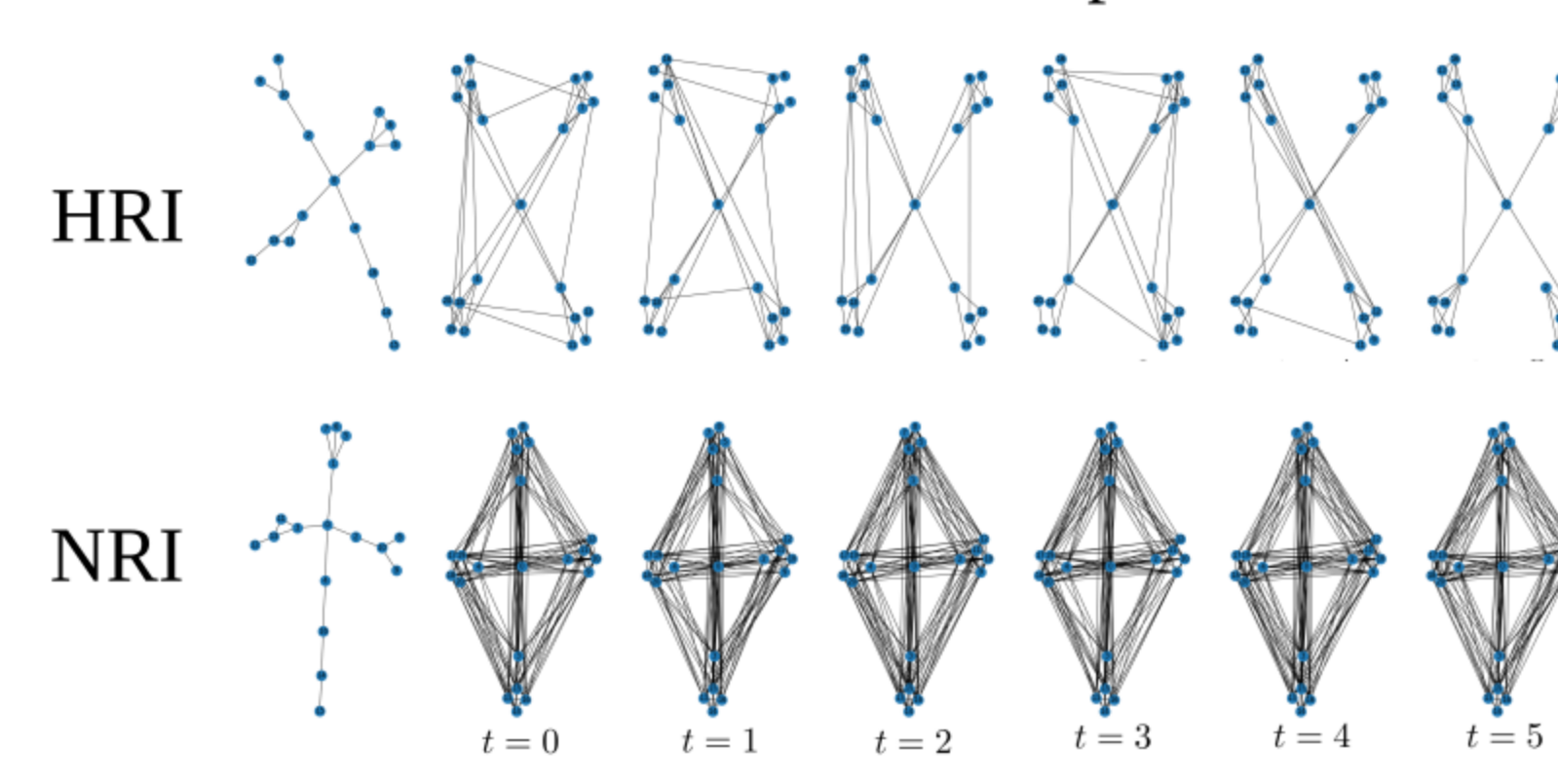
- Rendered version of the springs dataset only a part of the hierarchy is visible
- HRI **outperforms** both NRI and LSTM in terms of NLL for physical reasoning
- **Spatial slots** learned by the visual encoder capture **individual objects** and yield object representations
- HRI outperforms NRI at **inferring the interaction graph** in the visual domain



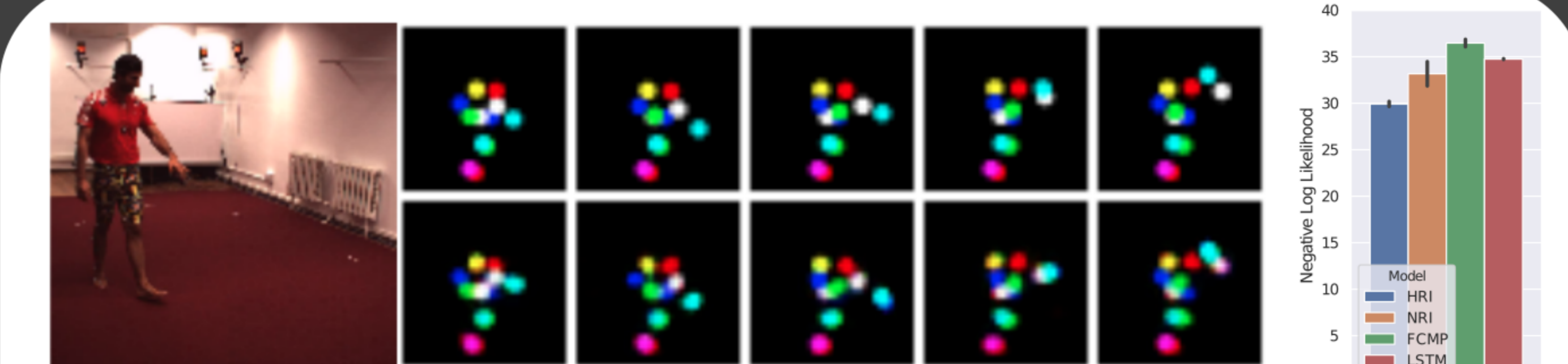
Object Slots Introspection



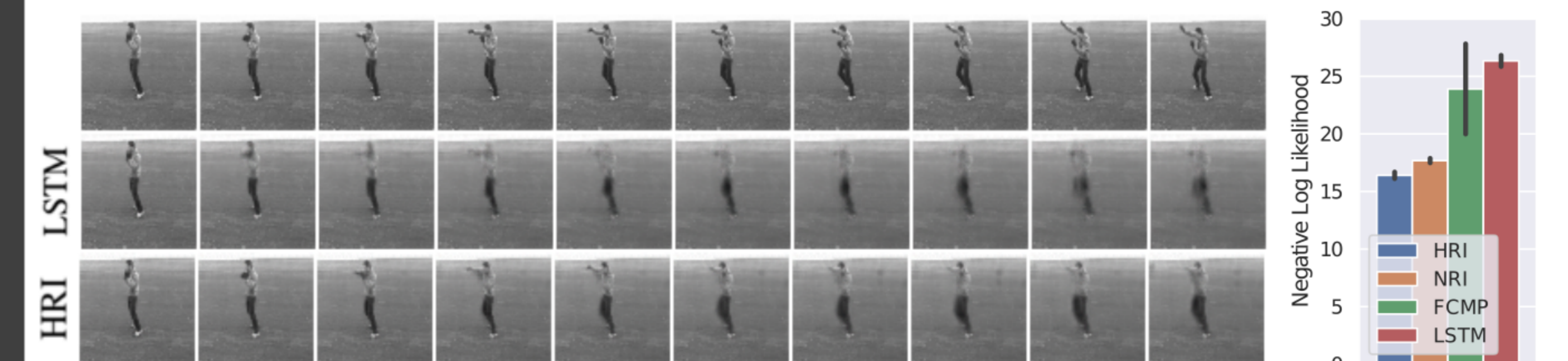
Inferred Graphs



Human 3.6M and KTH



- Render joints - a visual representation based on **keypoints**
- **Non-stationary dynamics** and **less hierarchical interactions**
- Still able to observe **improvements** using HRI



- Using **raw video frames** as input to the model
- Predictions made by HRI are **much more accurate** (limbs clearly depicted)

