# Long-term Human Motion Prediction with Scene Context

- Zhe Cao, Hang Gao , Karttikeya Mangalam , Qi-Zhi Cai, Minh Vo, and Jitendra Malik

- Prakash K Naikade,
prna00001@stud.uni-saarland.de

# Problem

- **Given:** $N$-step 2D human pose history $\mathbf{X}_{1:N}$ and Scene Image $\mathbf{I}$ ($N^{th}$ video frame)
- **Predict:** next $T$-step 3D human poses together with their locations ($\mathbf{Y}_{N+1:N+T}$)
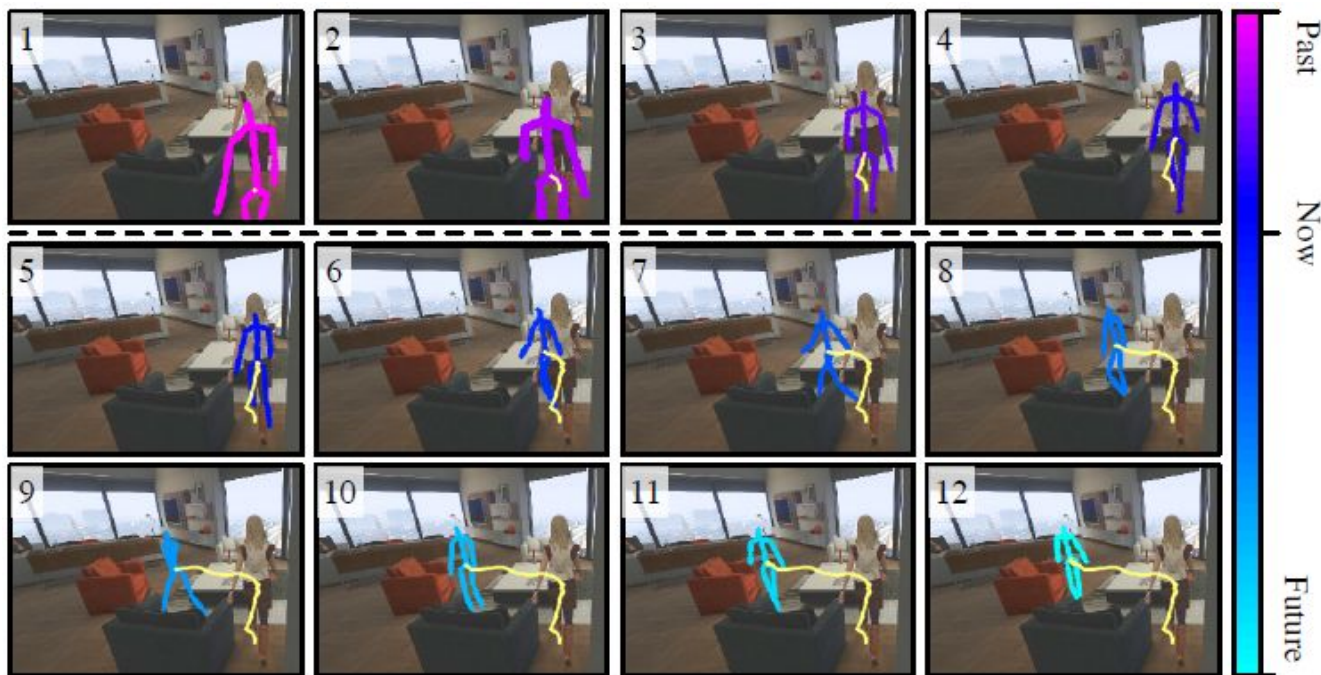


Fig 1: Problem description - Long-term 3D human motion prediction

# Motivation

- Human movement is,
    1. Goal directed
    2. Constrained by environment
    3. Multimodal future



Fig 2: Predict long-term human motion with scene context

# Motivation



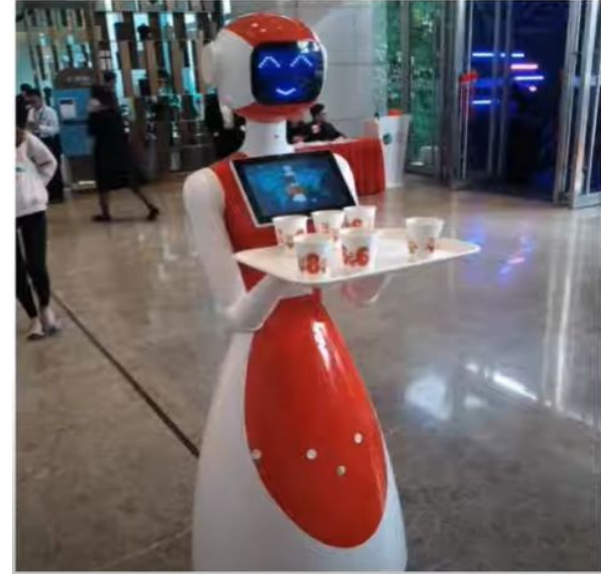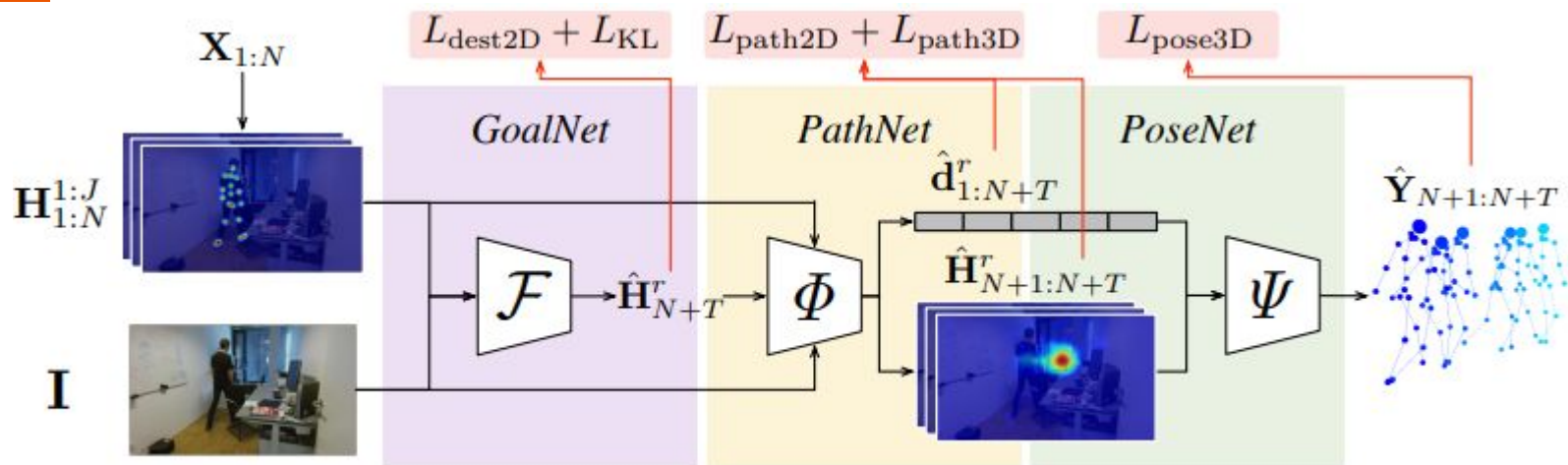Fig 3: Smart glass for vision impaired people



Fig 4: Home robot serves nearby passengers

# Key Contribution

1. **Formulated a new task**: Long-term 3D human motion prediction with scene context in terms of 3D poses and 3D locations.

2. **GTA-IM Dataset:** Created new synthetic dataset with diverse recordings of human-scene interaction and clean annotations.

   - Renderer Scripting - To generate one million RGBD frames of 1920 × 1080 resolution
   - Labels generated automatically:
     - RGBD Video
     - 3D human pose
     - Camera pose
     - Global coordinates of paths
     - Action labels
     - Human Segmentation

3. **Developed a novel three-stage computational framework:** Framework utilizes scene context for goal-oriented motion prediction.
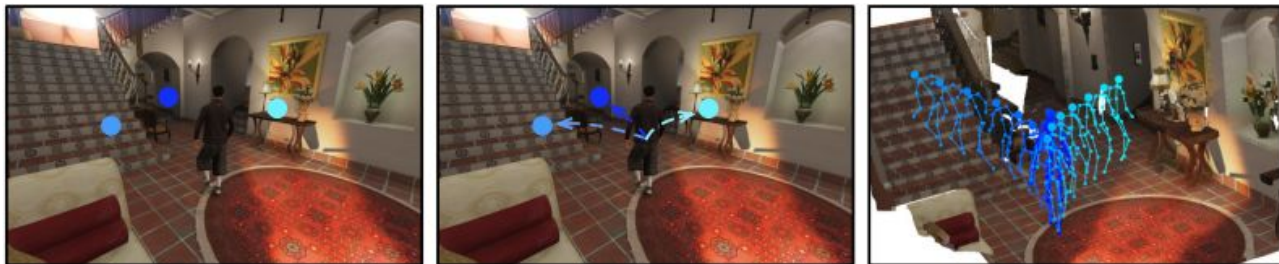
# Proposed Solution



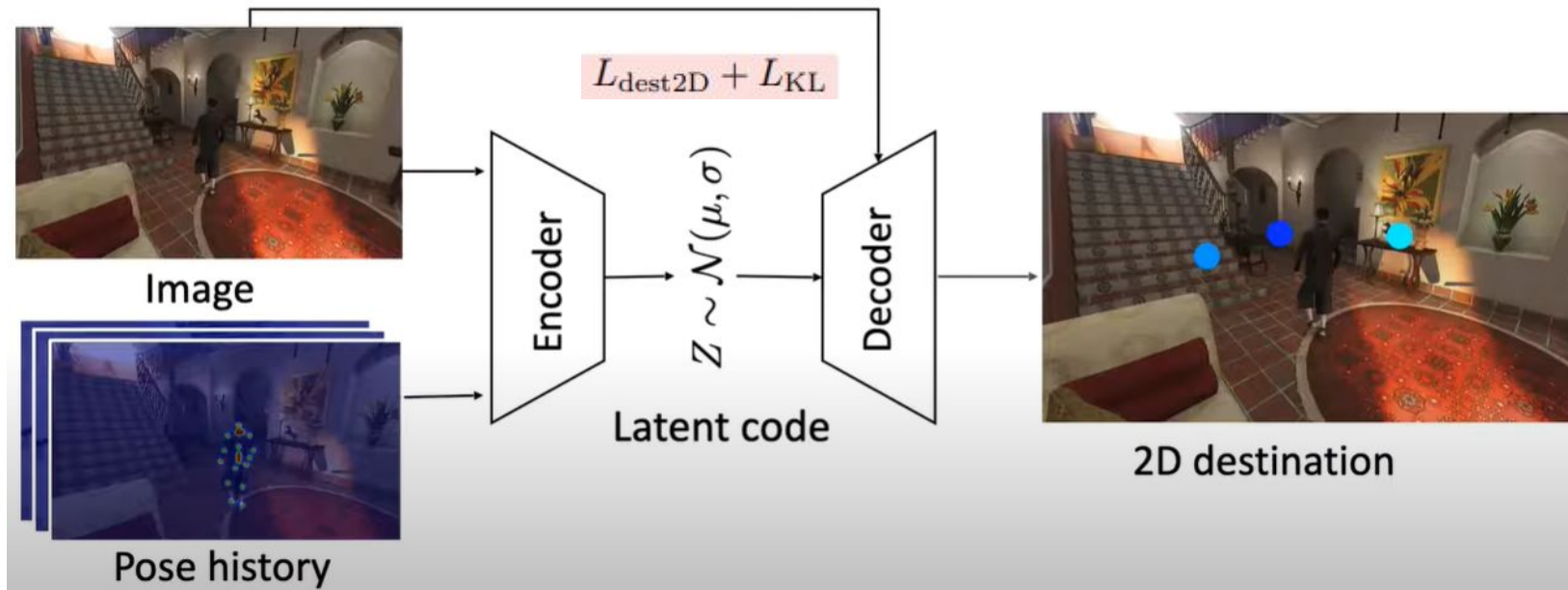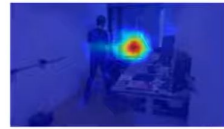Fig 4: Proposed Pipeline and Network architecture

# GoalNet



Fig 5: GoalNet - Predicting 2D Movement Destination
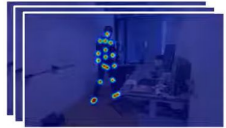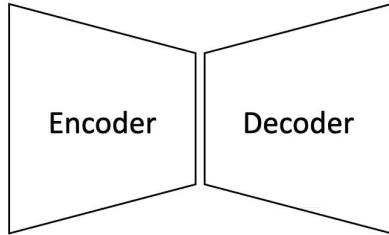
# PathNet



Destination

Image

Pose history

$$L_{\text{path2D}} + L_{\text{path3D}}$$

Encoder   Decoder

Depth vector

2D path

3D human path represented as 2D path
and depth values of human center

Fig 6: PathNet - Predicting 3D path towards each destination

8

# PoseNet

Final 3D pose
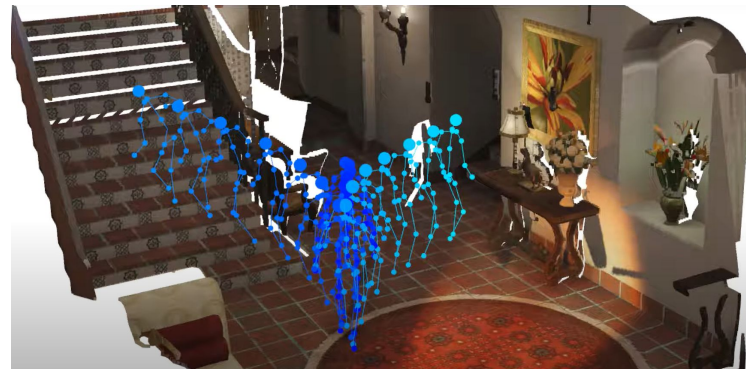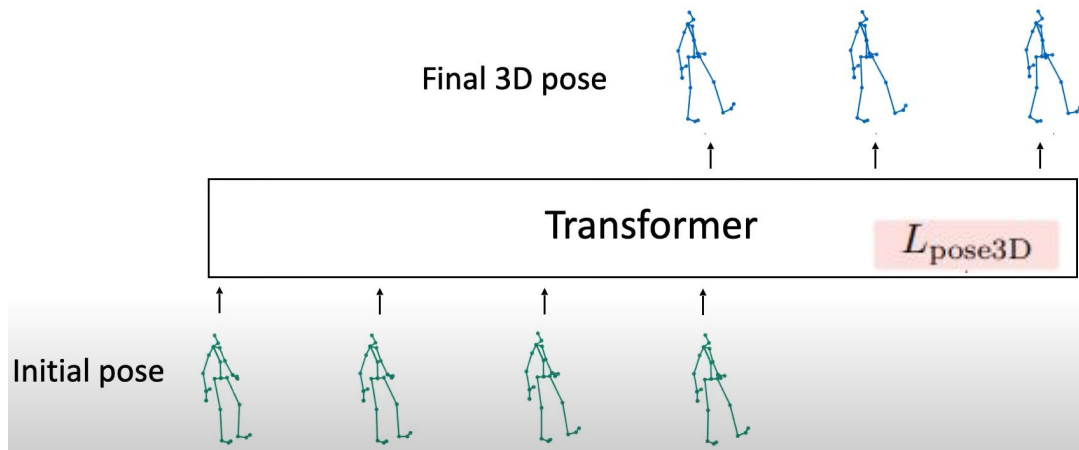
Transformer $L_{pose3D}$

Initial pose



Fig 7: PathNet - Generating 3D pose following the path

# Evaluation

1. **Average 3d distance between the two second long prediction and the ground truth**

# Limitations

1. **Resulting 3D poses may not strictly meet all physical constraints:** Use multi-view/temporal images.

2. **Dynamic objects and multiple moving people**

3. **Naturalness and feasibility of the stochastic human motion predictions**

4. **Domain gap between synthetic and realistic image dataset**

**Thank You !**

**Questions Please !**