### Applying Methods in RL to Character Animation Unsupervised RL IW Seminar

# Topic Agenda (Approx 9 min.)

- **Motivation & Goal**
- **Problem Background**
- **Related Work**
- **Approach** 
  - AMP Mod 1
  - AMP Mod 2
  - **Custom Dataset**
- **Implementation** 
  - **Code Review**
- Results
- Conclusion
  - **Future Work**



Motivation & Goal



#### Motivation & Goal

- Generation of assets for animated movies / games takes a lot of work! Involves creation of:
  - o 3D Mesh
  - Textures
  - Rigging
  - Animation of Rigging
- Prior methods in using motion tracking to generate animation sequences have proven to
   be data-intensive and expensive
- Unsupervised RL is well suited to learning distributional characteristics of certain motions
- Personal Interest

#### Goal

"Can we reduce pose error for Adversarial Motion Planning (AMP) agents for specific tasks (walking) by using custom training data or applying modifications to the base algorithm?

# Problem & Background

### Diving into AMP

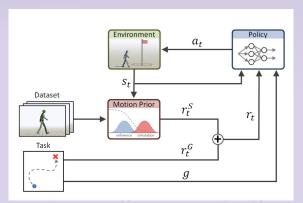
- Enables physics-based characters to learn natural movements from
   unstructured motion data without manual reward engineering or explicit
   motion planning
- Uses adversarial training to develop a motion prior that guides character
   behavior to match reference motion styles
- Combines task objectives with learned motion priors

Learns motion prior D(st, st+1) as discriminator trained to differentiate between real motion transitions from dataset vs those generated by policy  $\pi$ 

The discriminator is trained by solving a least-squares regression problem to predict a score of 1 for samples from the dataset and -1 for samples recorded from the policy. The reward function for training the policy is then given by

$$r(s_t, s_{t+1}) = \max \left[ 0, \ 1 - 0.25(D(s_t, s_{t+1}) - 1)^2 \right].$$
 (7)

The additional offset, scaling, and clipping are applied to bound the reward between [0, 1], as is common practice in previous RL frameworks [Peng et al. 2018a, 2016; Tassa et al. 2018]. Mao et al. [2017] showed that this least-squares objective minimizes the Pearson  $\chi^2$  divergence between  $d^{\mathcal{M}}(\mathbf{s},\mathbf{s}')$  and  $d^{\pi}(\mathbf{s},\mathbf{s}')$ .



### Diving into AMP

#### ALGORITHM 1: Training with AMP

```
1: input M: dataset of reference motions
 2: D \leftarrow initialize discriminator
 3: \pi \leftarrow initialize policy
 4: V \leftarrow initialize value function
 5: \mathcal{B} \leftarrow \emptyset initialize reply buffer
  6: while not done do
          for trajectory i = 1, ..., m do
              \tau^i \leftarrow \{(\mathbf{s}_t, \mathbf{a}_t, r_t^G)_{t=0}^{T-1}, \mathbf{s}_T^G, \mathbf{g}\} collect trajectory with \pi for time step t=0,...,T-1 do
                  d_t \leftarrow D(\Phi(\mathbf{s}_t), \Phi(\mathbf{s}_{t+1}))
                  r_t^S \leftarrow calculate style reward according to Equation 7 using d_t
                  r_t \leftarrow w^G r_t^G + w^S r_t^S
                  record r_t in \tau^i
               end for
              store \tau^i in \mathcal{B}
          end for
          for update step = 1, ..., n do
              b^{\mathcal{M}} \leftarrow \text{sample batch of } K \text{ transitions } \{(\mathbf{s}_j, \mathbf{s}_j')\}_{j=1}^K \text{ from } \mathcal{M}
              b^{\pi} \leftarrow \text{ sample batch of } K \text{ transitions } \{(\mathbf{s}_j, \mathbf{s}_j')\}_{i=1}^K \text{ from } \mathcal{B}
19:
              update D according to Equation 8 using b^{\mathcal{M}} and b^{\pi}
          end for
          update V and \pi using data from trajectories \{\tau^i\}_{i=1}^m
23: end while
```

and 7 illustrate examples of motions learned by the characters. Performance is evaluated using the average pose error, where the pose error  $e_t^{\text{pose}}$  at each time step t is computed between the pose of the simulated character and the reference motion using the relative positions of each joint with respect to the root (in units of meters),

$$e_t^{\text{pose}} = \frac{1}{N^{\text{joint}}} \sum_{j \in \text{joints}} \left\| (\mathbf{x}_t^j - \mathbf{x}_t^{\text{root}}) - (\hat{\mathbf{x}}_t^j - \hat{\mathbf{x}}_t^{\text{root}}) \right\|_2.$$
 (10)

 $\mathbf{x}_t^j$  and  $\hat{\mathbf{x}}_t^j$  denote the 3D Cartesian position of joint j from the simulated character and the reference motion, and  $N^{\mathrm{joint}}$  is the total number of joints in the character's body. This method of evaluating

### Related Work

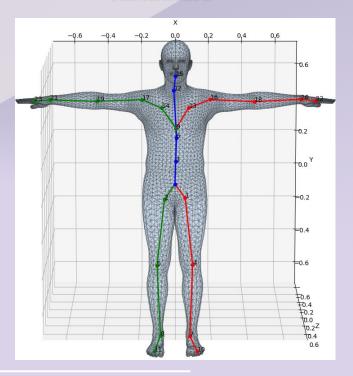
#### Related Work

#### *IsaacSim*





#### HumanML3D



#### Related Work

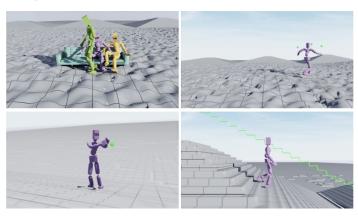
#### **ProtoMotions: Physics-based Character Animation**

"Primitive or fundamental types of movement that serve as a basis for more complex motions."

- · What is this?
- Installation guide
- Training built-in agents
- Evaluating your agent
- · Building your own agent/environment

#### What is this?

This codebase contains our efforts in building interactive physically-simulated virtual agents. It supports both IsaacGym and IsaacSim.



# Approach



#### Approach

#### Algorithm 1 Training with Adaptive AMP

```
Input: \mathcal{M}: dataset of reference motions
\alpha \leftarrow 0.5 {initialize adaptive weight}
D \leftarrow \text{initialize discriminator}
\pi \leftarrow \text{initialize policy}
V \leftarrow initialize value function
\mathcal{B} \leftarrow \emptyset {initialize reply buffer}
while not done do
    \begin{array}{ll} \textbf{for} \ \text{trajectory} \ i = 1, \ldots, m \ \textbf{do} \\ \tau^i \leftarrow \{(s_t, a_t, r_t^G)_{t=0}^{T-1}, s_T^G, g\} \ \{\text{collect trajectory with} \ \pi\} \end{array} 
      for time step t = 0, \dots, T-1 do
          d_t \leftarrow D(\Phi(s_t), \Phi(s_{t+1}))
          r_t^S \leftarrow calculate style reward according to Equation 7 using d_t
          r_t \leftarrow (1-\alpha)w^G r_t^G + \alpha w^S r_t^S {adaptive weighting}
          record r_t in \tau^i
       end for
       store \tau^i in \mathcal{B}
   end for
   for update step = 1, \ldots, n do
       b^{\mathcal{M}} \leftarrow \text{sample batch of } K \text{ transitions } \{(s_j, s_j')\}_{j=1}^K \text{ from } \mathcal{M}
      b^{\pi} \leftarrow \text{sample batch of } K \text{ transitions } \{(s_j, s'_i)\}_{i=1}^{K} \text{ from } \mathcal{B}
      update D according to Equation 8 using b^{\mathcal{M}} and b^{\pi}
   end for
   update V and \pi using data from trajectories \{\tau^i\}_{i=1}^m
   \alpha \leftarrow update adaptive weight based on relative losses
end while
```

#### Algorithm 2 Training with Hierarchical AMP

```
Input: \mathcal{M}: dataset of reference motions
D_L, D_H \leftarrow initialize local and global discriminators
\pi \leftarrow \text{initialize policy}
V \leftarrow initialize value function
\mathcal{B} \leftarrow \emptyset {initialize reply buffer}
while not done do
   for trajectory i = 1, \ldots, m do
       \tau^i \leftarrow \{(s_t, a_t, r_t^G)_{t=0}^{T-1}, s_T^G, g\}  {collect trajectory with \pi}
       for time step t = 0, \dots, T-1 do
          d_t^L \leftarrow D_L(\Phi_L(s_t), \Phi_L(s_{t+1})) \text{ {local features}} 
d_t^H \leftarrow D_H(\Phi_H(s_{t:t+k})) \text{ {global features}}
          r_t^S \leftarrow \text{calculate combined style reward using } d_t^L, d_t^H
          r_t \leftarrow w^G r_t^G + w^S r_t^S
          record r_{t} in \tau^{i}
       end for
       store \tau^i in \mathcal{B}
    end for
    for update step = 1, \ldots, n do
       b^{\mathcal{M}} \leftarrow \text{sample batch of } K \text{ transitions } \{(s_i, s_i')\}_{i=1}^K \text{ from } \mathcal{M}
       b^{\pi} \leftarrow \text{sample batch of } K \text{ transitions } \{(s_j, s_i')\}_{i=1}^K \text{ from } \mathcal{B}
       update D_L, D_H according to Equation 8 using b^{\mathcal{M}} and b^{\pi}
    end for
   update V and \pi using data from trajectories \{\tau^i\}_{i=1}^m
end while
```

#### Approach

- AMASS dataset is comprised of the following subsets:
  - ACCAD, BMLhandball, BMLmovi, BMLrub, CMU, CNRS, DanceDB, DFaust, EKUT,
     EyesJapanDataset, GRAB, HDM05, HUMAN4D, HumanEva, KIT, MoSh, PosePrior, SFU, SOMA,
     SSM, TCDHands, TotalCapture, Transitions, WEIZMANN
- Selected Subset: HumanEva Dataset
  - o <a href="http://humaneva.is.tue.mpg.de/">http://humaneva.is.tue.mpg.de/</a>
  - Focuses on basic human actions
  - Walking, jogging, gesturing
  - Reasoning: try to avoid overfitting for the simple task of walking humanoid



## Implementation

### Implementation (Algo 1)

New

```
def calculate_extra_reward(self):
    rew = super().calculate extra reward()
    if self.disable discriminator:
        return rew
   discriminator_obs = self.experience_buffer.discriminator_obs
   disc r = self.calculate discriminator reward(
        discriminator_obs.view(self.num_envs * self.num_steps, -1)
    ).view(self.num steps, self.num envs)
    self.experience buffer.batch update data("discriminator rewards", disc r)
    extra reward = disc r + rew
    return extra_reward
```

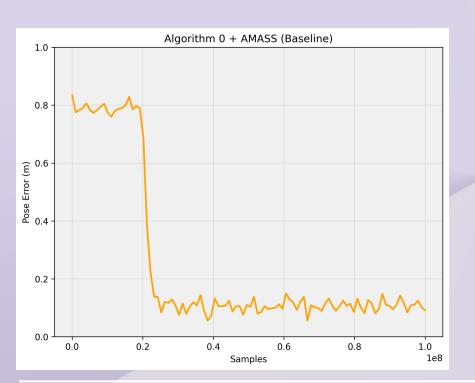
Old

#### Implementation (Algorithm 2)

Old

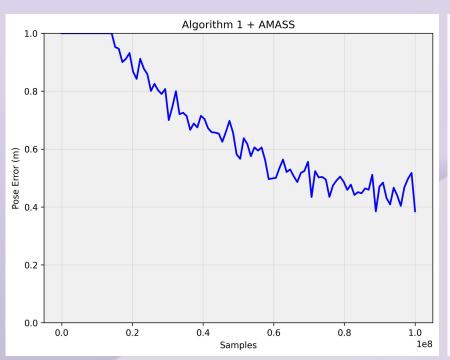
```
def discriminator_forward(self, obs: Tensor, return_norm_obs=False) -> Tensor:
    args = {"obs": obs}
    return self.discriminator(args, return_norm_obs=return_norm_obs)
```

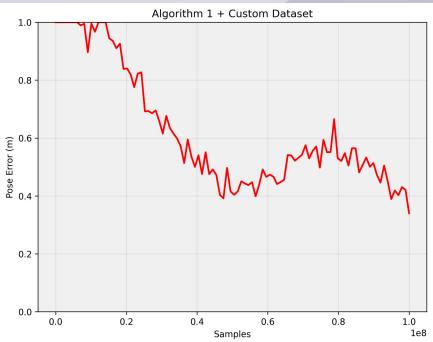
New

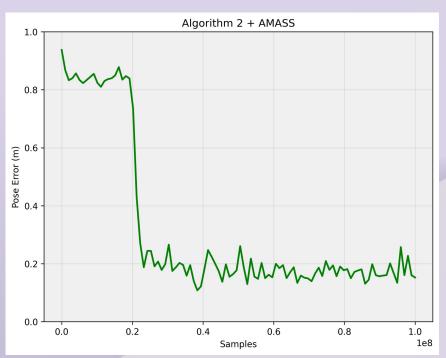


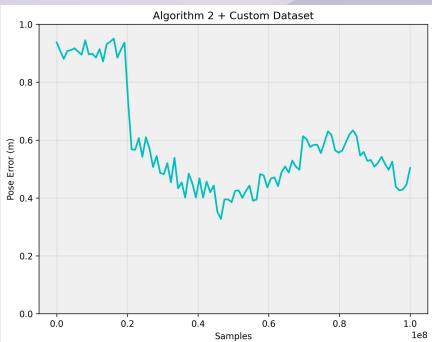
$$e_t^{\mathrm{pose}} = \frac{1}{N^{\mathrm{joint}}} \sum_{j \in \mathrm{joints}} \left\| (\mathbf{x}_t^j - \mathbf{x}_t^{\mathrm{root}}) - (\hat{\mathbf{x}}_t^j - \hat{\mathbf{x}}_t^{\mathrm{root}}) \right\|_2.$$

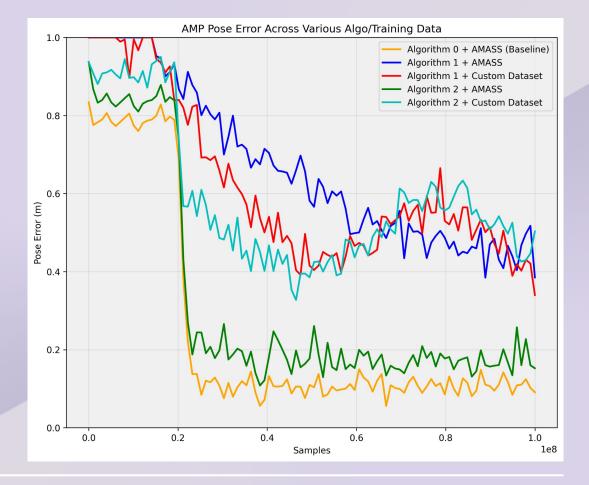
(Pose Error compared to motion file provided in ProtoMotions repo for basic humanoid walking)











### Conclusion



#### Conclusion

- Future Work
  - Try different subsets of the AMASS dataset
  - Experiment with different task focuses or try a generalized implementation
  - Experiment with different adjustments to AMP base algorithm
- Importance
  - Contribute to reducing expense/labor of animation for movies / games
  - Train agents to understand characteristics of motion from motion capture data