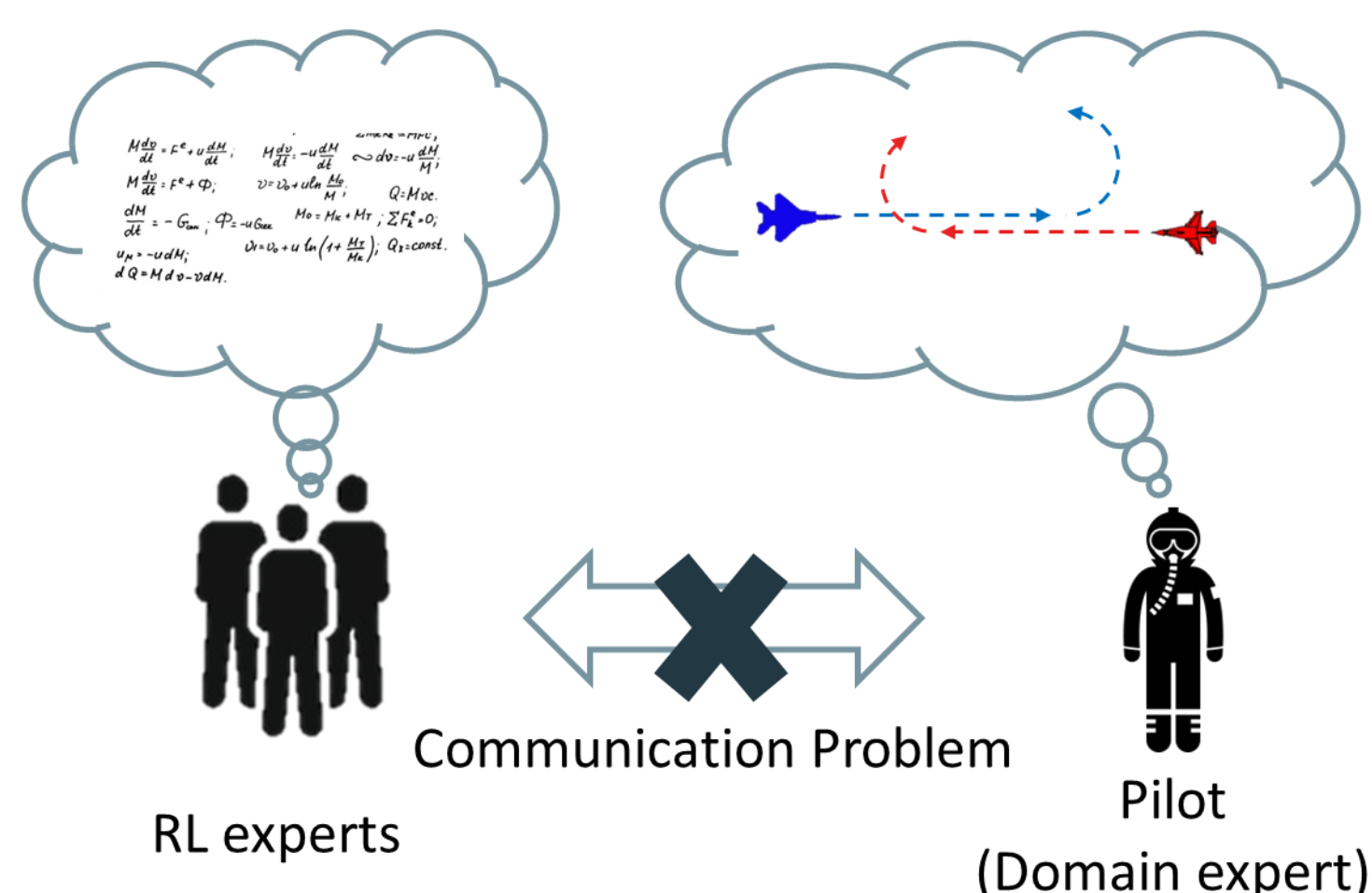




Motivation

- Create a visual representation of RL agents' behavior understandable to domain experts without ML expertise

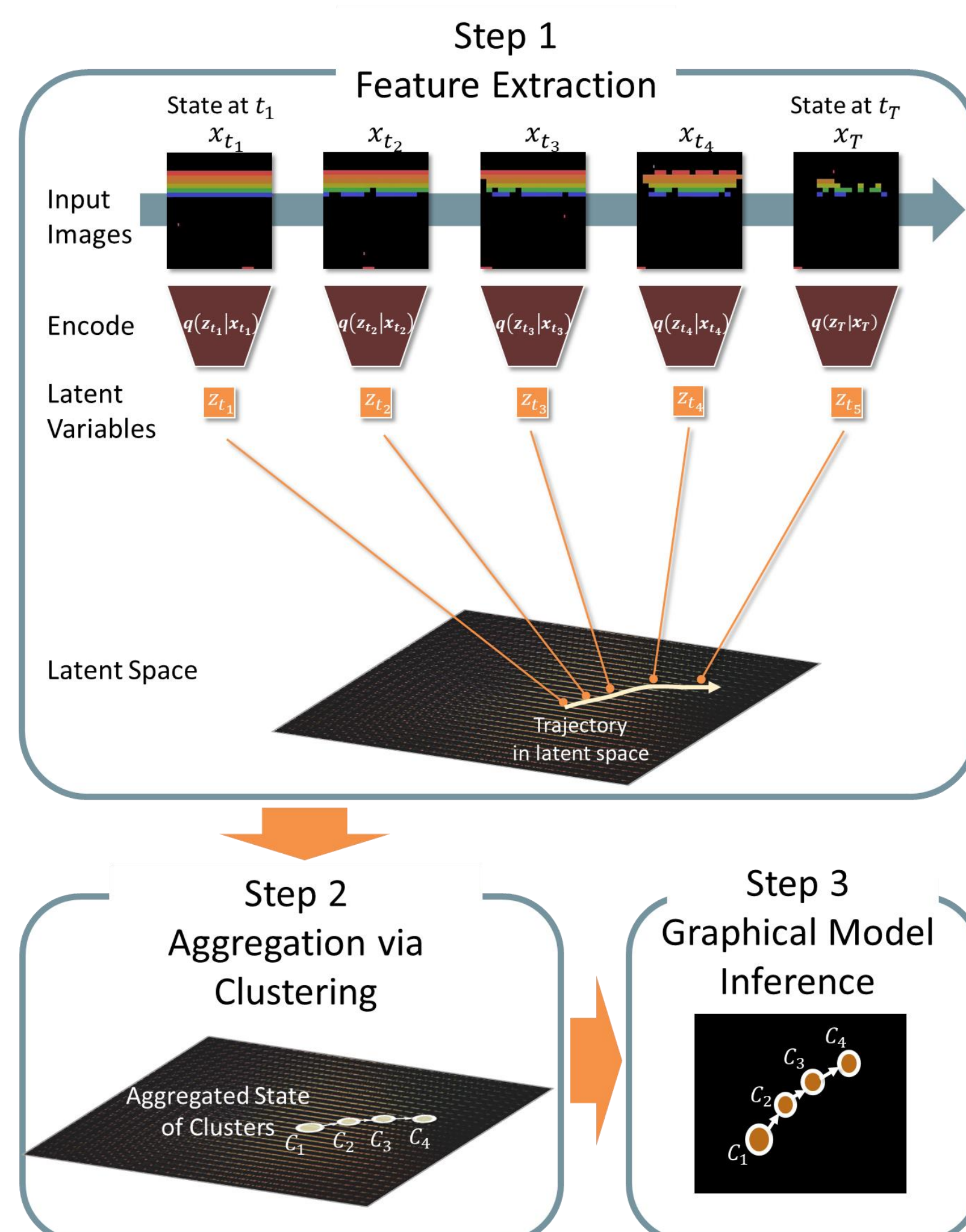


Approach

- Use replay to provide non descriptive knowledge about an agent
- Abstract & visualize it as a trajectory

Methodology

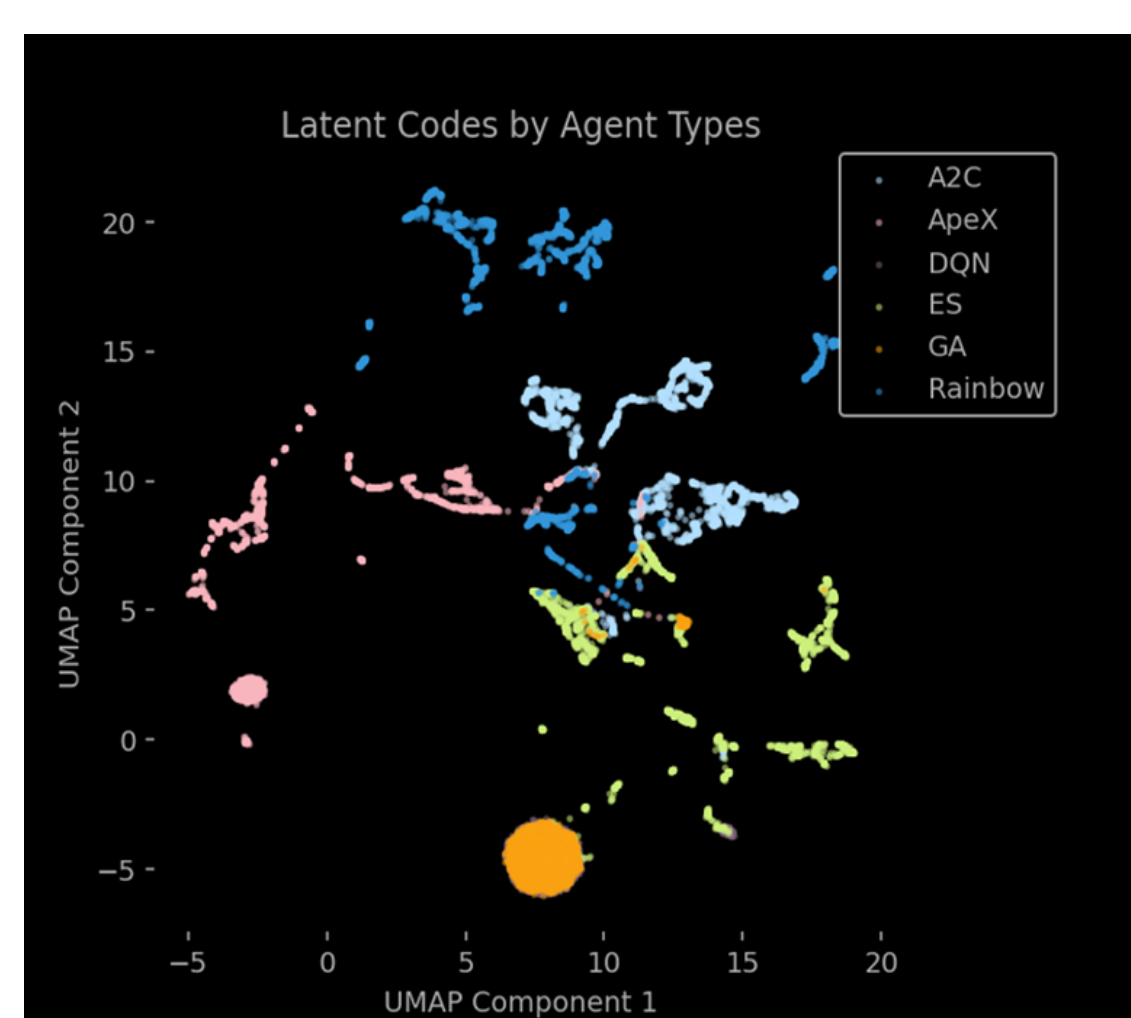
- Used β -VAE & ST-DBSCAN
- Tested to various Atari applications



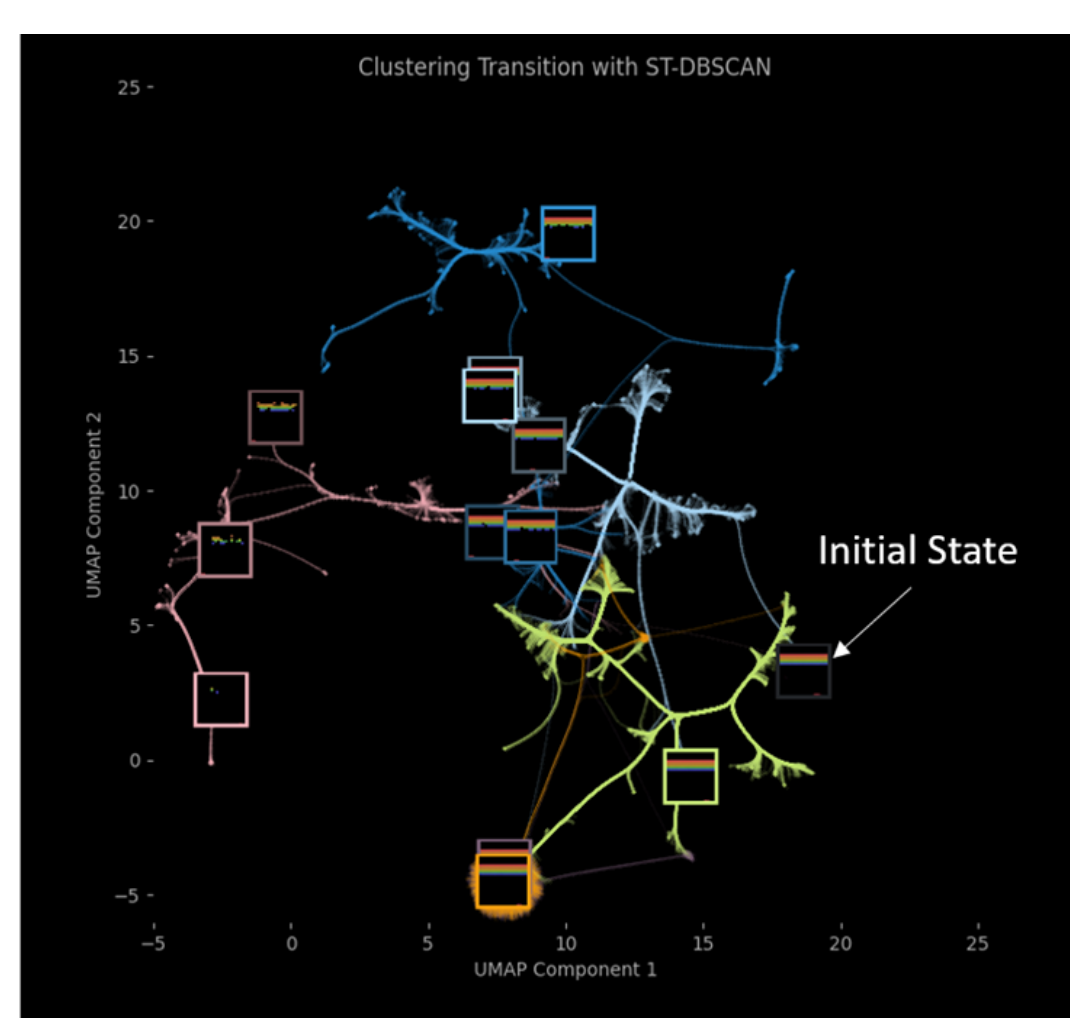
Result

- Extracted trajectories of several pre-trained agents with β -VAE
- Explored visualization ideas of abstracted trajectories

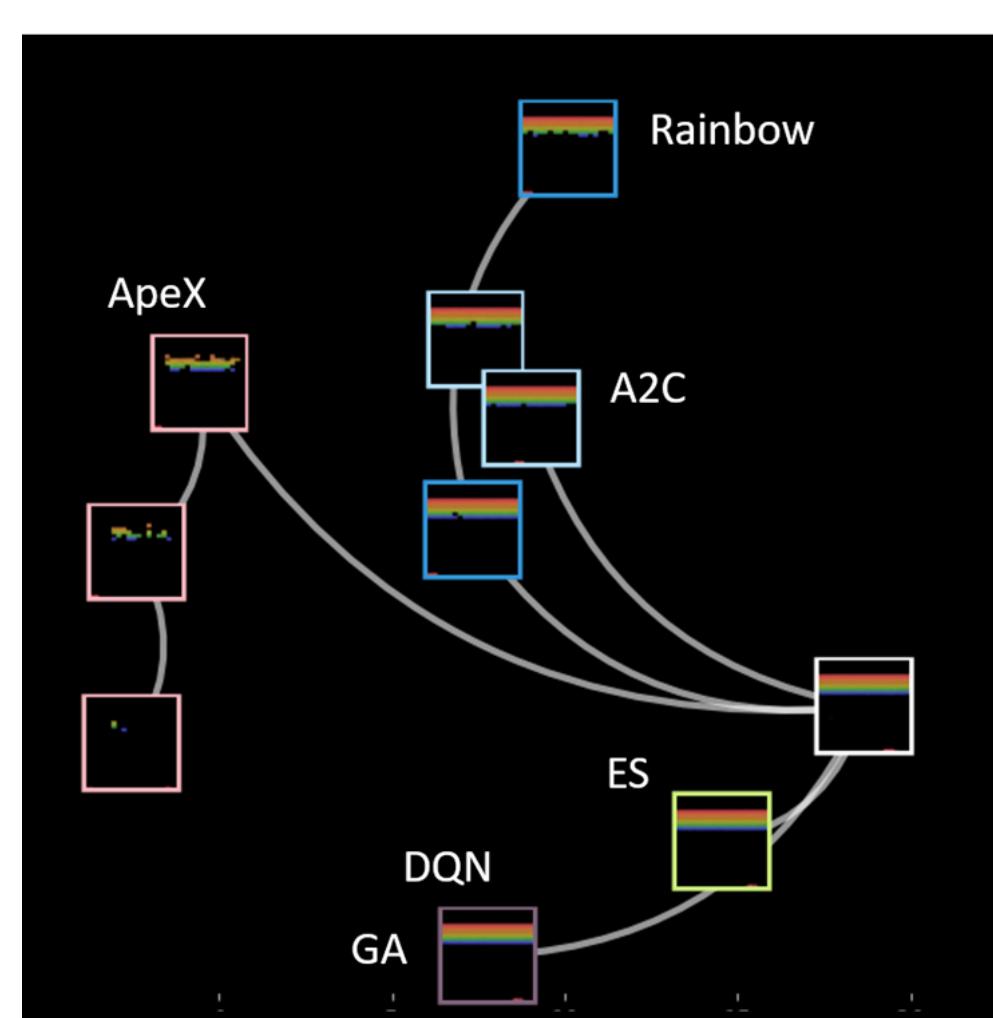
Raw trajectory



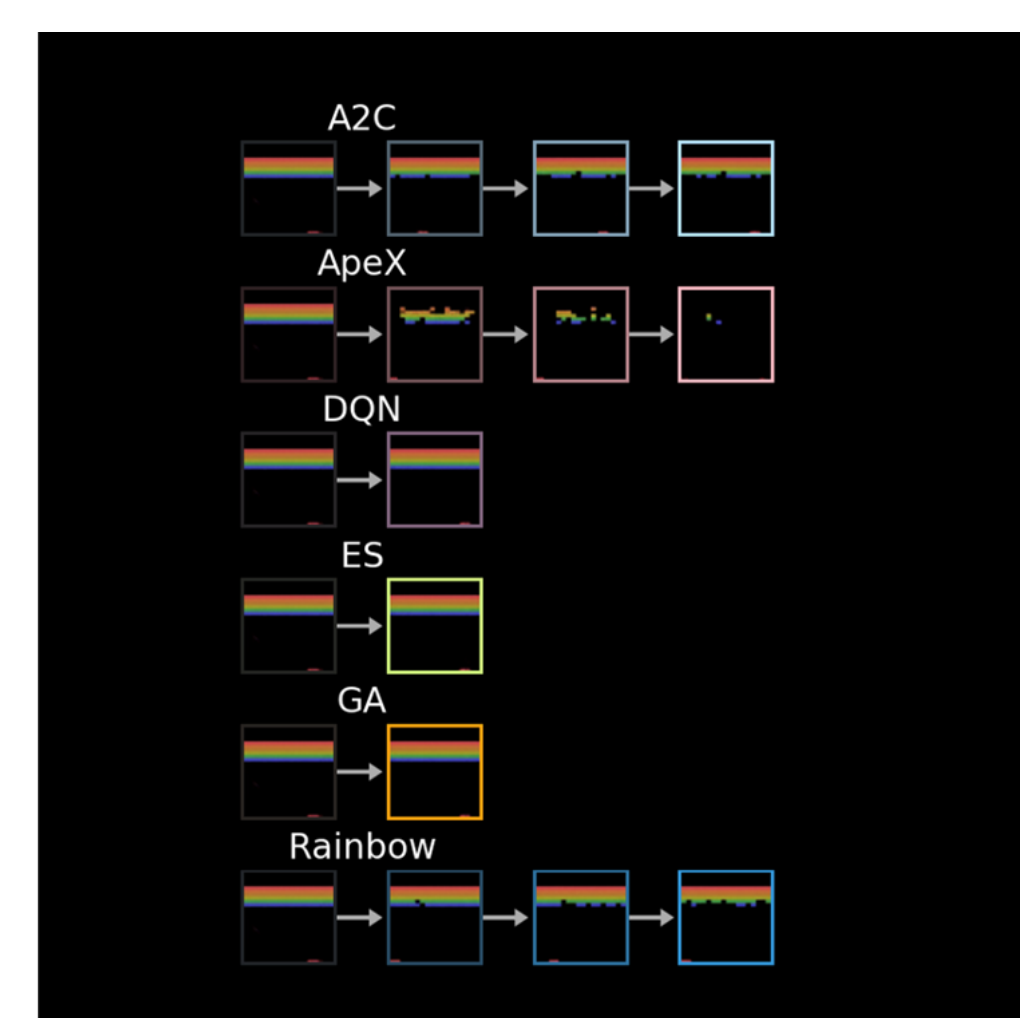
Clustering & Edge bundling



Abstracted trajectory #1



Abstracted trajectory #2



User Study Plan

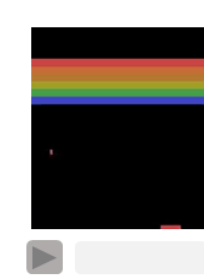
- Evaluate how well a user's mental model obtained from the proposed abstract trajectory agrees with agents' complete trajectory

Evaluation

- ✓ Accuracy
- ✓ Confidence
- ✓ Response time
- ✓ Preference

Complete Trajectory

Instruction:
The following movie shows a player playing a game. After watching the movie, please answer the following question.



Question:
Which player do you think the following game movie is played by?

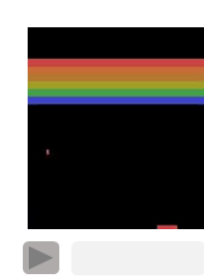
Player #1

Player #2

Player #3

Abstracted Trajectory

Instruction:
The following movie shows a player playing a game. After watching the movie, please answer the following question.



Question:
Each color-coded branch in the flowchart below is the result of a cutout of a different player playing a game. Which player do you think the following game movie is played by?

Player #1

Player #2

Player #3