

## RGB observations



$t_{0-2}$

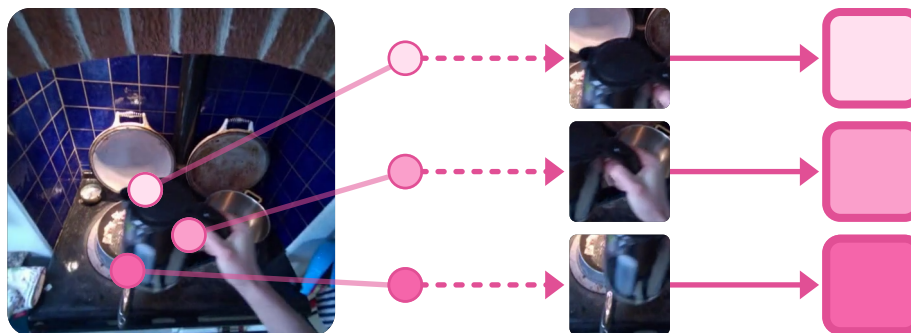
$t_{0-1}$

$t_0$

## Action

*pour the water into the pot*

## 2D query point features



$t_0$

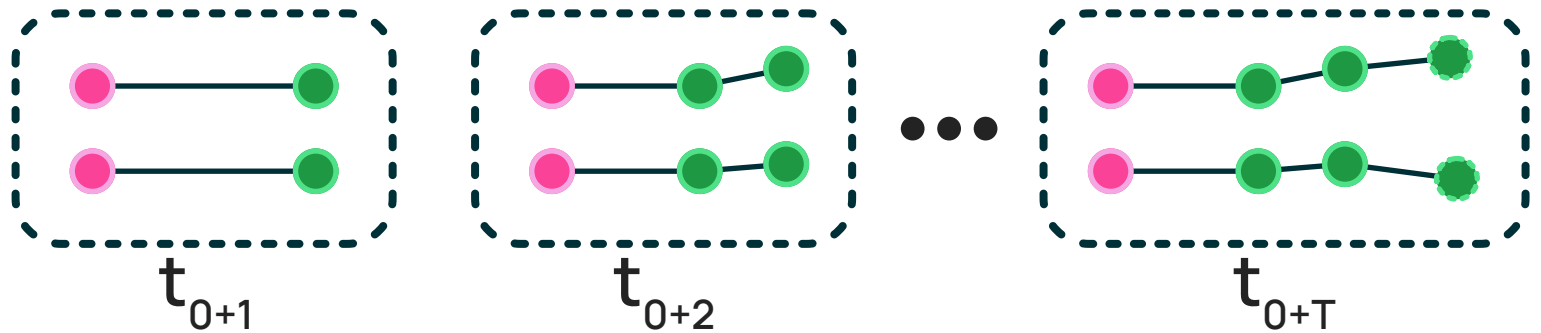


OR

## Autoregressive prediction



`<track coord="t0 p1 x y z  
p2 x y z"></track>` → `<track coord="t1 p1 x y z  
p2 x y z; t2 ..."></track>`



DIT head

## Flow-matching prediction

