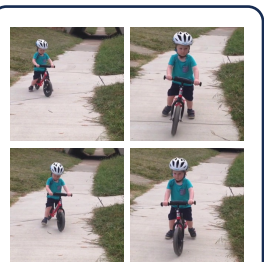


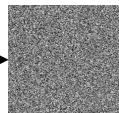
## Train:

### Stage 1: Instance Descriptor Learning



**Diffusion  
Encoder**

[v] person



### Diffusion Decoder

ResNet

2D-Attention

⋮

ResNet

2D-Attention

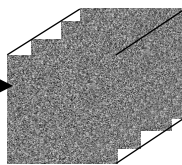


### Stage 2: Temporal Encoding



**Diffusion  
Encoder**

[v] person is  
riding a bike



### Diffusion Decoder

ResNet

2D-Attention

Temporal-  
Attention

⋮

ResNet

2D-Attention

Temporal-  
Attention



Freezing Layers

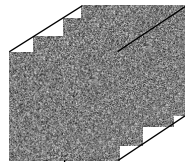


Optimising Layers



## Inference:

**DDIM  
Inversion**



[v] person  
is walking  
in the road

