

One-Step Time-Dependent Future Video Frame Prediction with a Convolutional Encoder-Decoder Neural Network

Vedran Vukotić^{1,2,3}, **Silvia-Laura Pintea**¹, **Christian Raymond**^{2,3}, **Guillaume Gravier**^{3,4}, **Jan van Gemert**¹
 vedran.vukotic@irisa.fr s.l.pintea@tudelft.nl christian.raymond@irisa.fr guillaume.gravier@irisa.fr j.c.vangemert@tudelft.nl

Problem

• given the current video frame at time t_0 and an arbitrary temporal displacement t predict the frame at $t_0 + t$

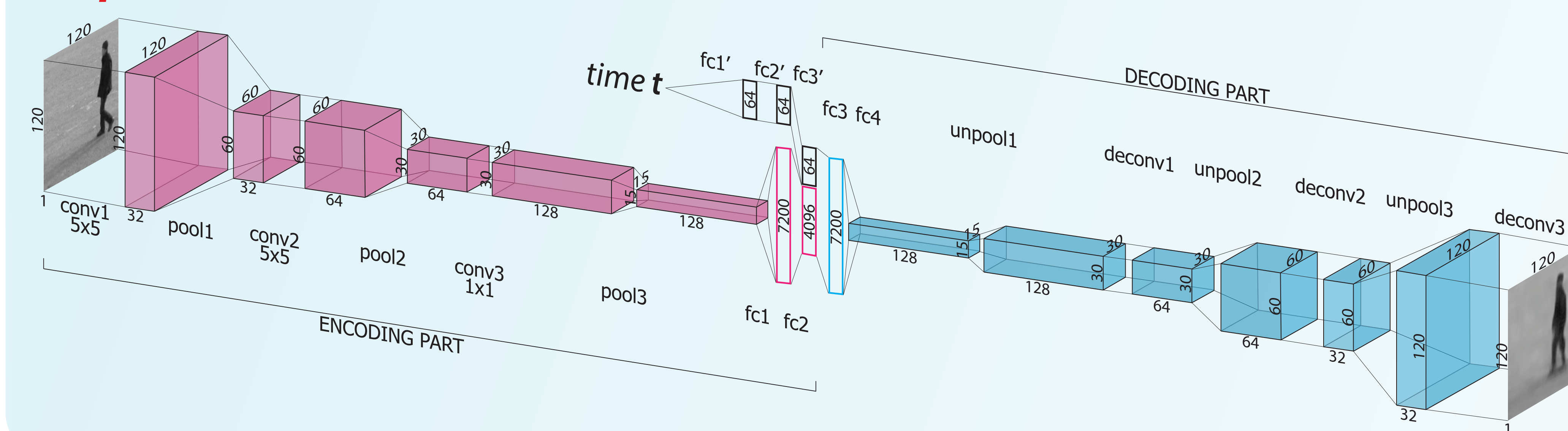
Goal:

• anticipate future motion-induced appearance change

Means:

- creating representation that encodes appearance changes over time
- embedding the input image and a continuous time variable
- translating back to the image space to visualize the anticipated video frame

Proposed Architecture



Idea

- **encoding network**
 - **image encoding branch**
 - **time encoding branch** (time modeled as a continuous variable t)
- **decoding network**
- inspired by the architecture in [1]

Existing Work

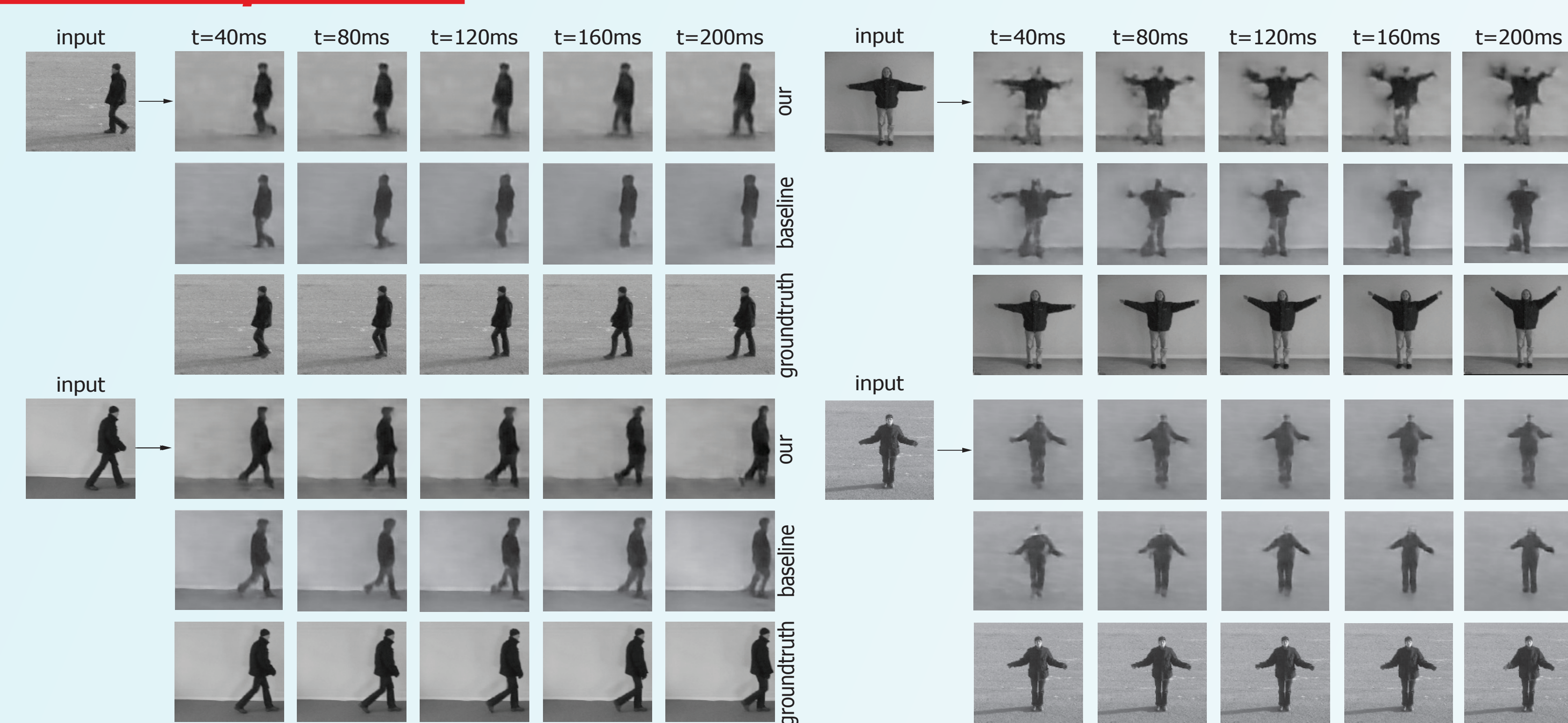
Predicting Future Motion:

- given an image, predict optical flow at the next timestep
- given an image predict motion trajectories

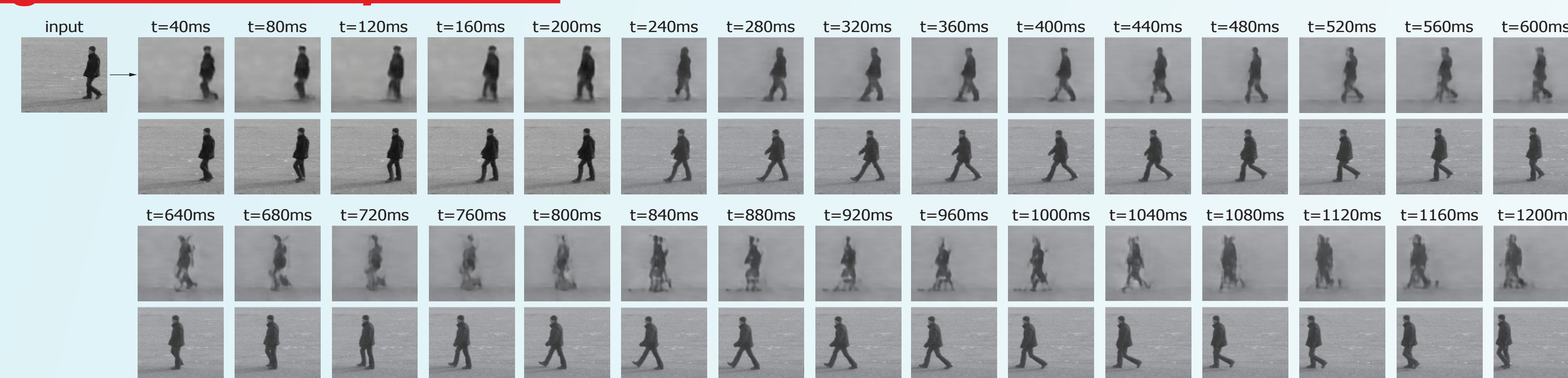
Predicting Future Appearance:

- hallucinating possible images (conditioned GANs)
- predicting future pixels from previous pixels (Pixel Networks)
- **autoencoding methods** - predicting the future image at the next timestep

Example Anticipations

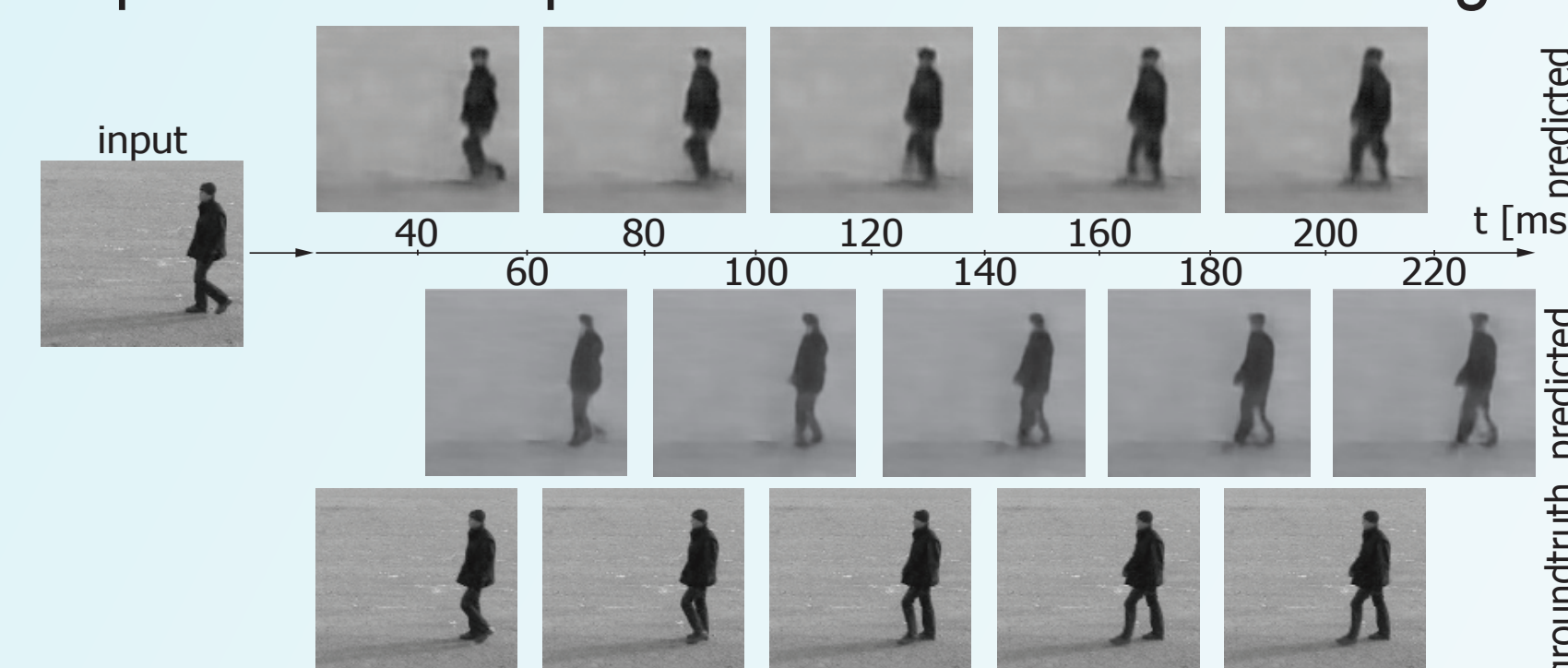


Long-Term Anticipations



Unseen Time Displacements

• anticipations at temporal distances not seen during training:



Conclusion

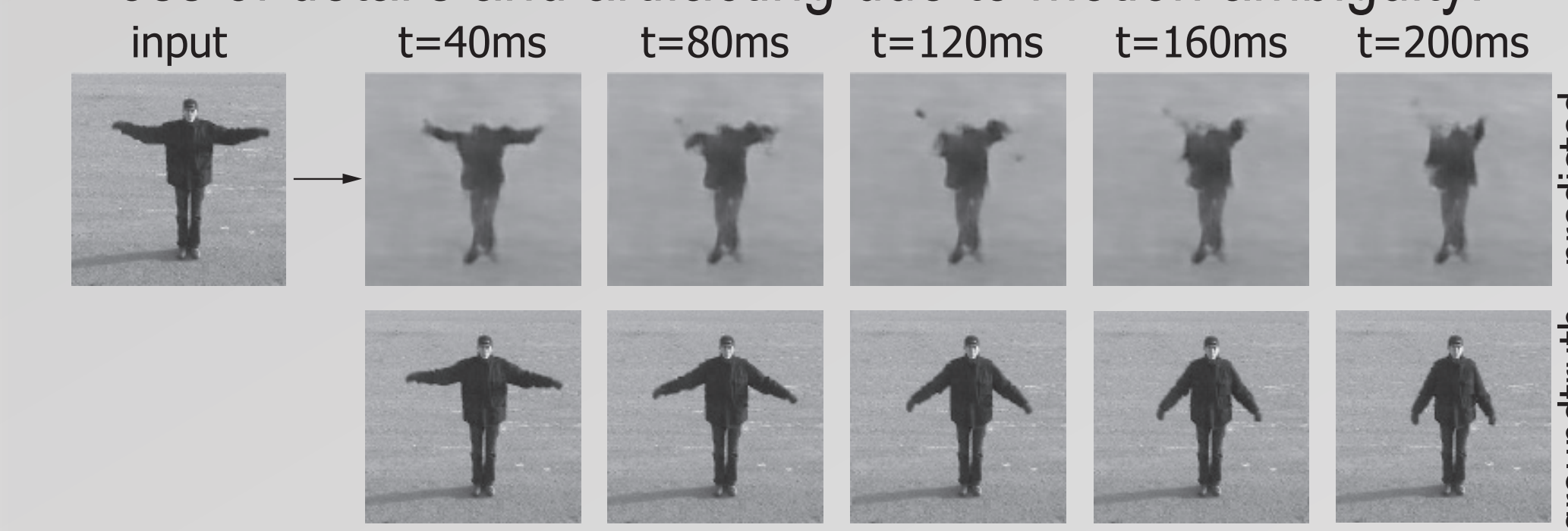
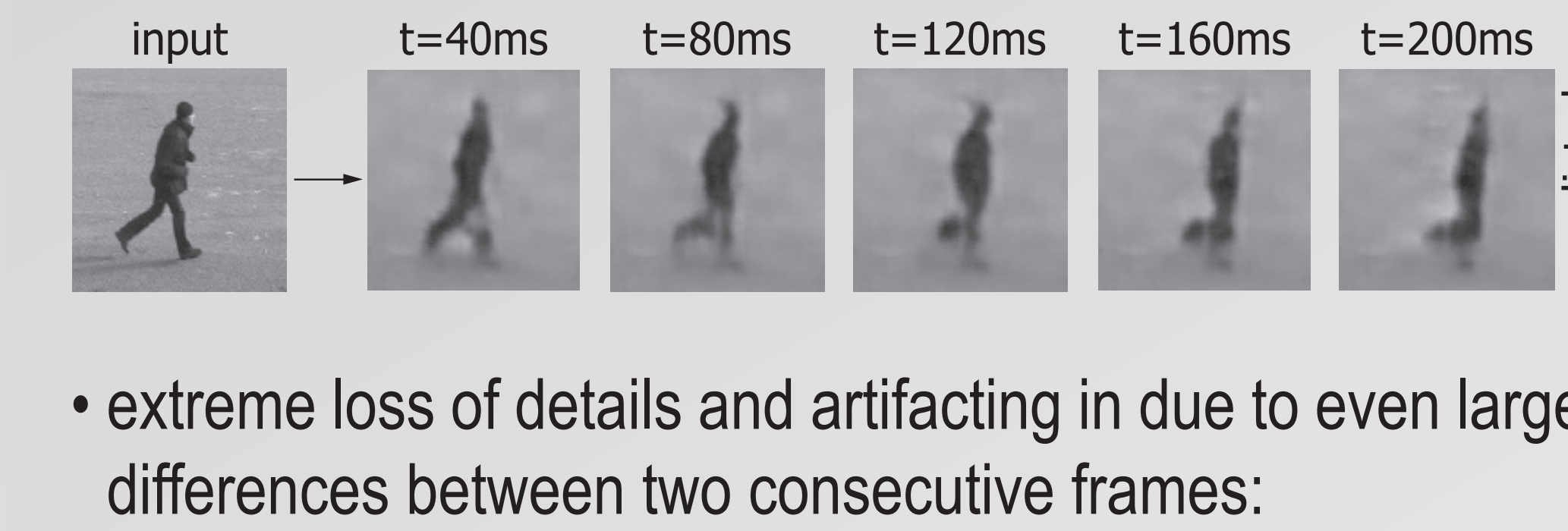
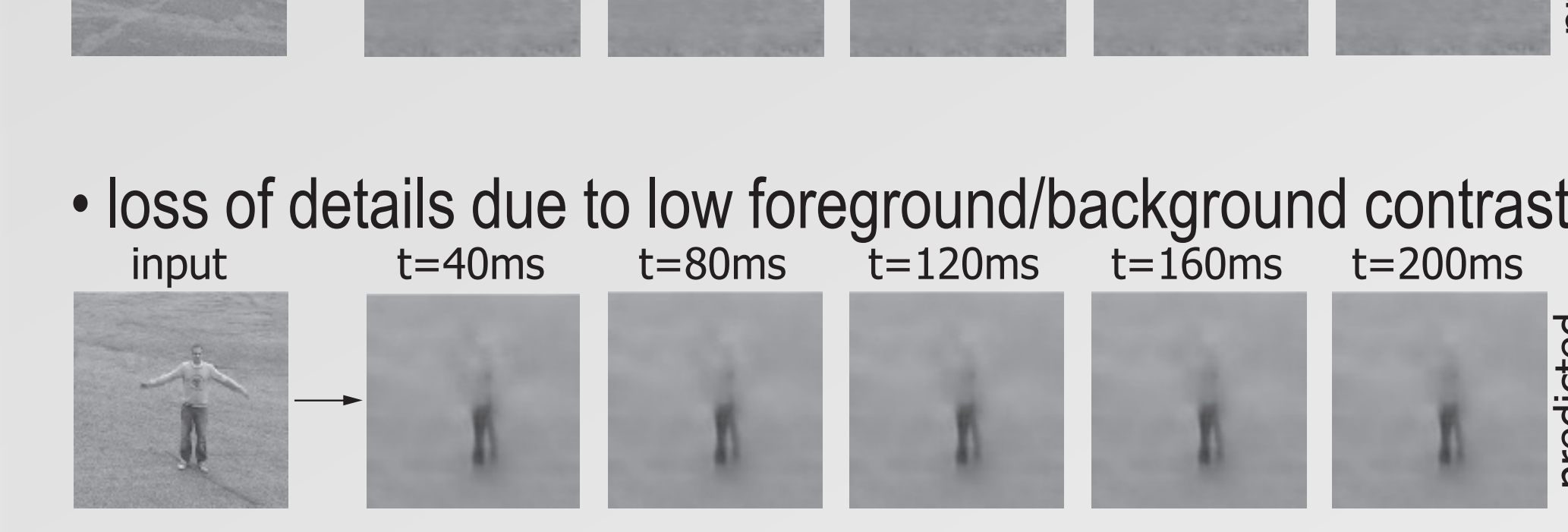
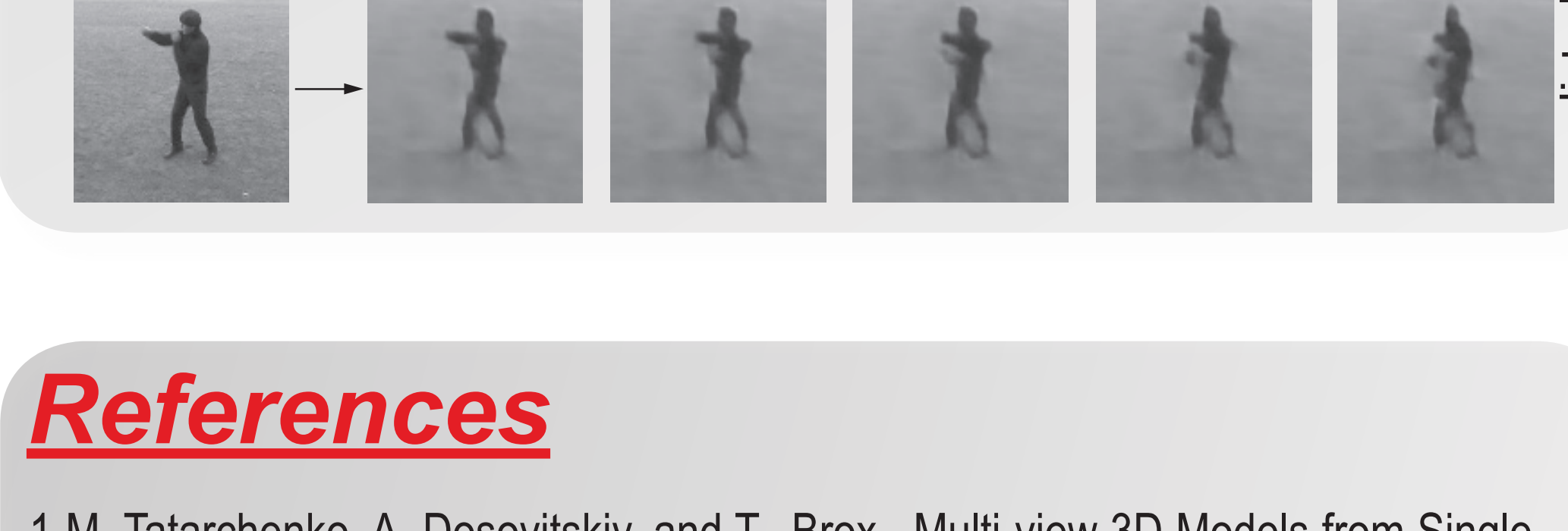

Successes:

- successfully predicts future frames at **arbitrary** temporal displacements, including temporal displacements **never seen during training**
- predictions are done directly, in one step

Downsides:

- unable to tackle ambiguities; artifacting and loss of details caused by addressable issues in videos

Downsides

- loss of details and artifacting due to motion ambiguity:
 
- loss of details due to larger differences between two consecutive frames:
 
- extreme loss of details and artifacting in due to even larger differences between two consecutive frames:
 
- loss of details due to low foreground/background contrast:
 
- artifacting due to small and sporadic movements:
 

References

1.M. Tatarchenko, A. Dosovitskiy, and T. Brox. Multi-view 3D Models from Single Images with a Convolutional Network. In ECCV 2016, Amsterdam, The Netherlands, 2016.