

The Machine Learning of Time: Past & Future

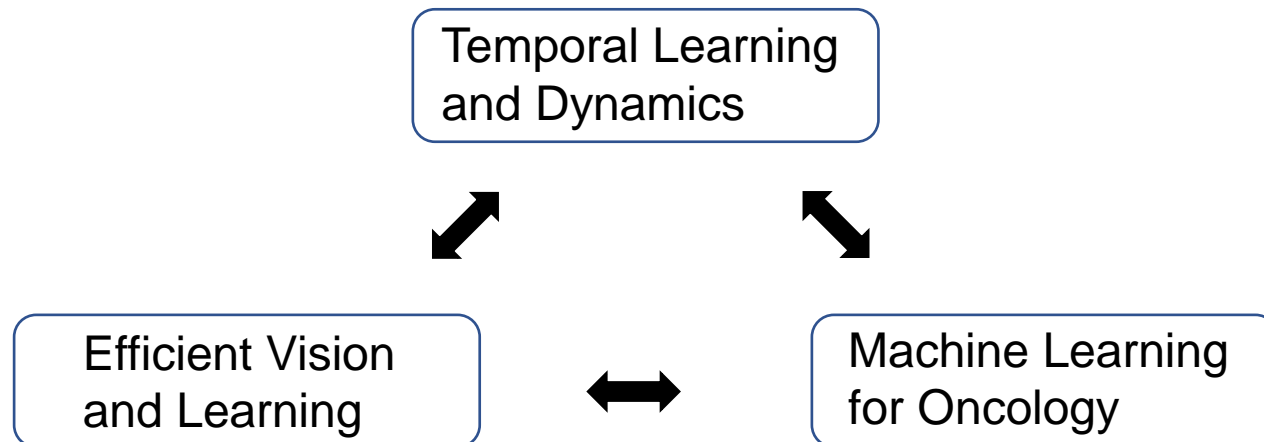
Efstratios Gavves, University of Amsterdam, egavves@uva.nl



www.i-aida.org

About Me

- Associate Professor at the University of Amsterdam
 - Director of QUVA and POP-AART Lab (we will be hiring!)
- Co-founder of Ellogon.AI
 - AI to personalize to immunotherapy in oncology
- ELLIS scholar



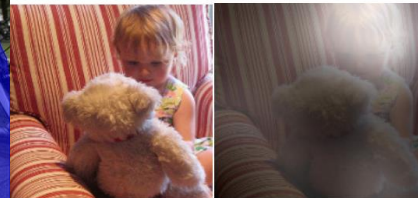
The Golden Age of Learning Algorithms



A woman is throwing a frisbee in a park.



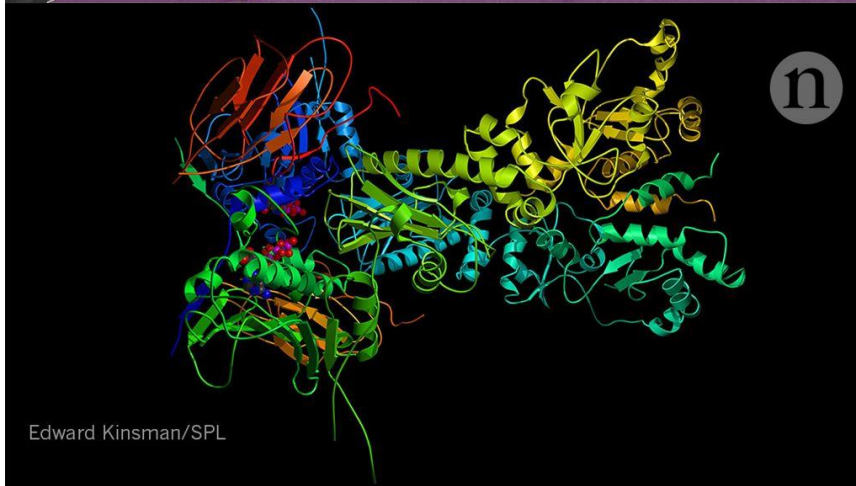
A dog is standing on a hardwood floor.



A little girl sitting on a bed with a teddy bear.



A group of people sitting on a boat in the water.



An Urgent Paradox

Apple falling: videos reversed, shuffled or normal \Rightarrow no difference^{1,2}

Normal video: 83.1%



Reversed frames: 82.9%



State-of-the-art spatiotemporal models ignore time

An Urgent Paradox

Apple falling: videos reversed, shuffled or normal \Rightarrow no difference^{1,2}

Normal video: 83.1%



Reversed frames: 82.9%



Urgent for forecasting



State-of-the-art spatiotemporal models ignore time

An Urgent Paradox

Apple falling: videos reversed, shuffled or normal \Rightarrow no difference^{1,2}

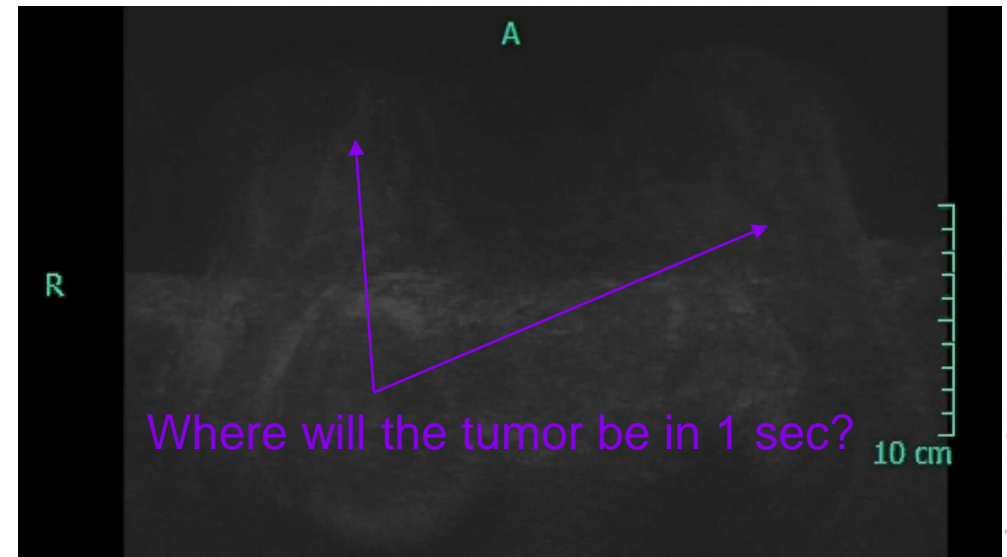
Normal video: 83.1%



Reversed frames: 82.9%



Urgent for future planning



State-of-the-art spatiotemporal models ignore time

An Urgent Paradox

Apple falling: videos reversed, shuffled or normal \Rightarrow no difference^{1,2}

Normal video: 83.1%



Reversed frames: 82.9%



Urgent for autonomous driving

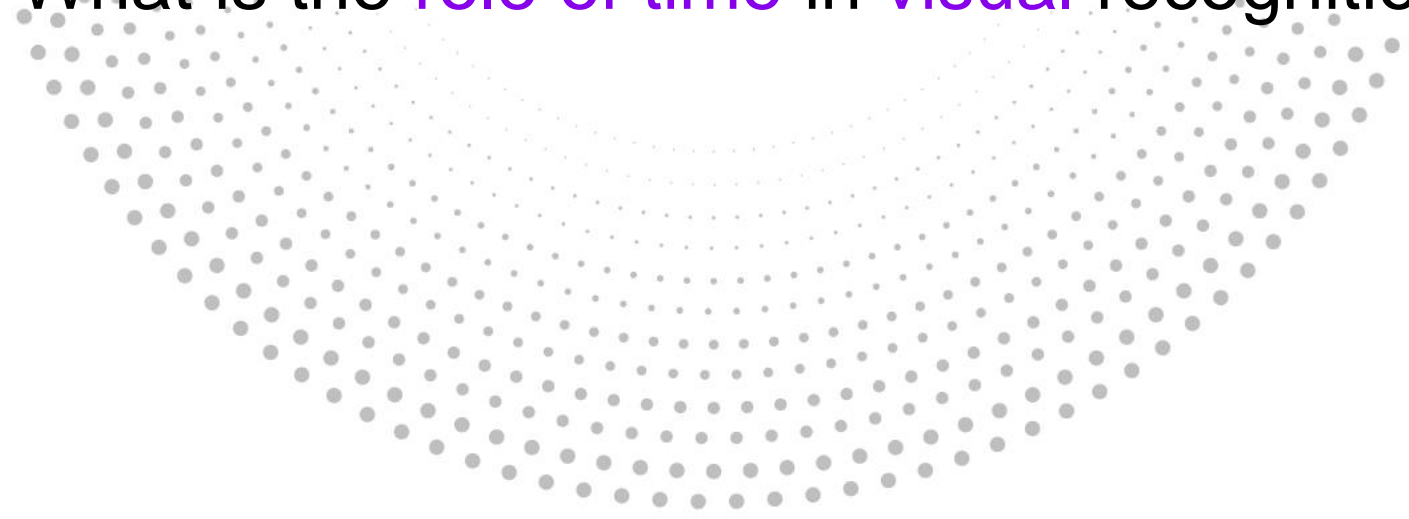


State-of-the-art spatiotemporal models **ignore time**



Central Question

What is the **role of time** in **visual** recognition?





The Vision

Models that **learn temporality** in **entangled** spatiotemporal sequences

The Vision

Models that **learn temporality** in entangled spatiotemporal sequences



+8.3

The Vision

Models that **learn temporality** in entangled spatiotemporal sequences



+1.1



-2.3



The Vision

Models that **learn temporality** in entangled spatiotemporal sequences



+1.1



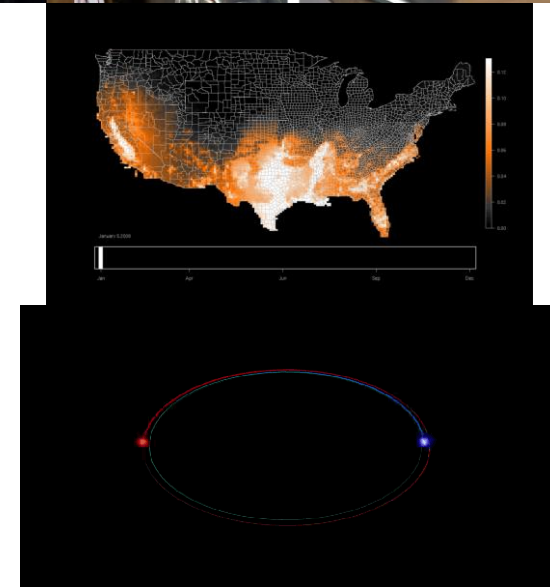
-5.1



+2.3

Entangled spatiotemporal data

- Data in thousands of dimensions confounded space and time
- Example 1: Long & complex videos
- Example 2: Migration patterns
- Example 3: Particles through time



What's the challenge?

- Thousands of frames → A lot of correlations and dynamics

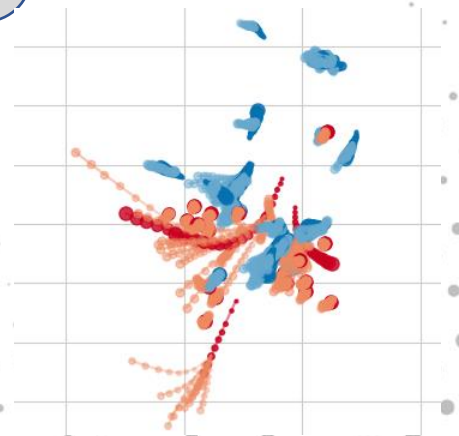
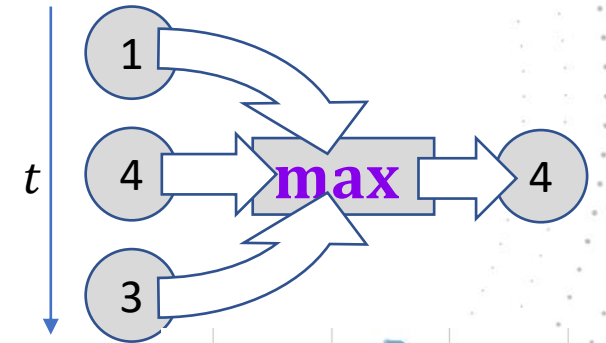


Challenge #1: State-of-the-art discards time by aggregating with set operations

Challenge #2: Hard to annotate manually → supervised learning is debatable

Challenge #3: A sequence is one of myriad possibilities → generative modelling critical

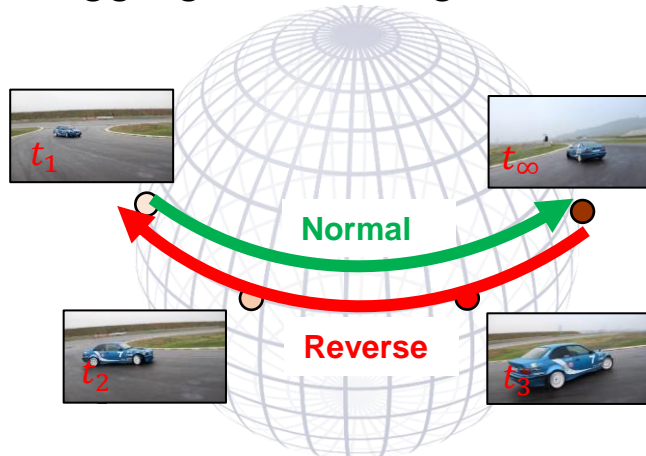
Challenge #4: Lack of standardization ← data is huge, algorithms very complex



Addressing the challenges

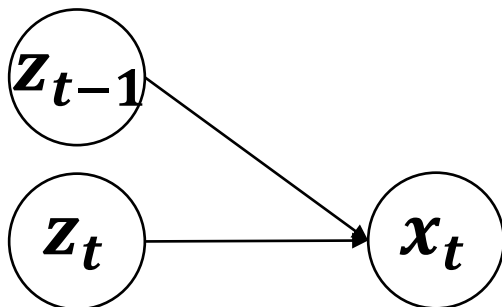
Time Geometry

- Learn spatiotemporal geometric manifolds
- Aggregate over a geodesic time path on manifold



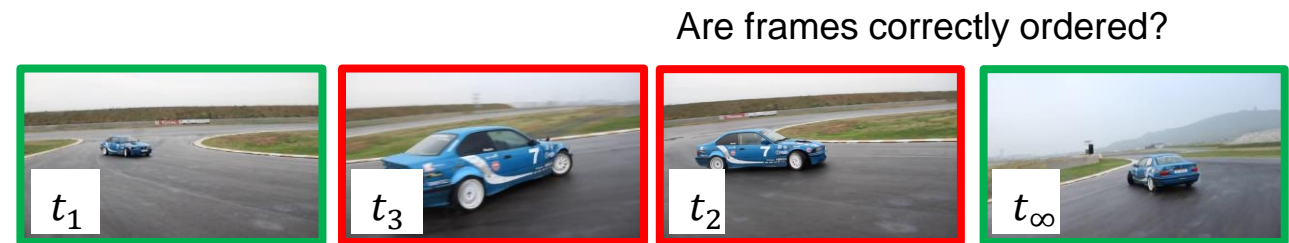
Time Generation

- Models that imagine all possible futures
- Spatiotemporal generative/bayesian models



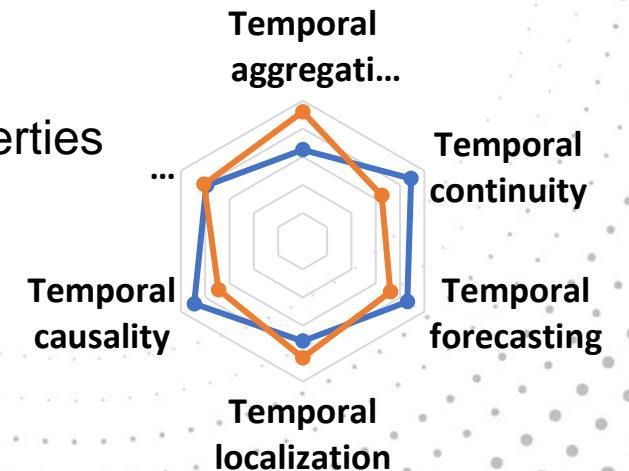
Time Supervision

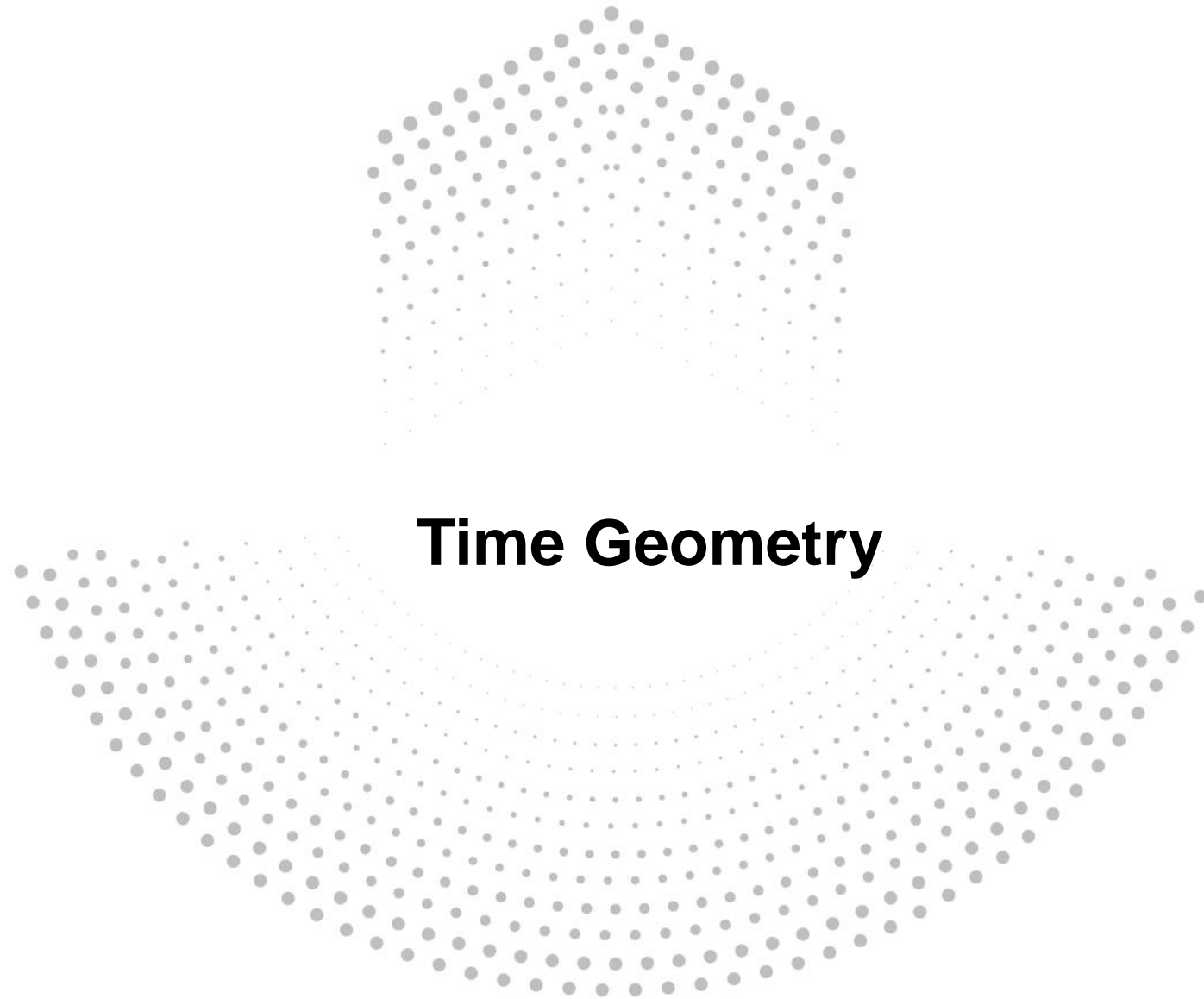
- Replace manual annotation with time properties
- In particular, combine with time-sensitive models



Time Evaluation

- Standardize data
- Evaluate on temporal properties





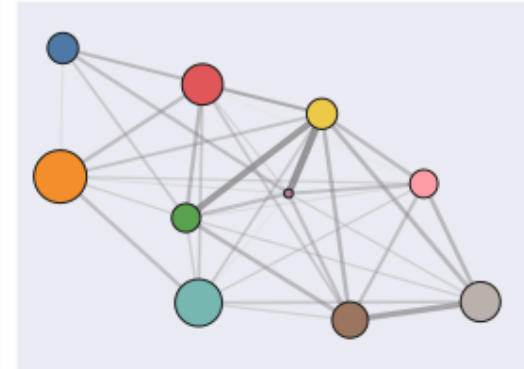
Time Geometry

Types of spatiotemporal geometric learning

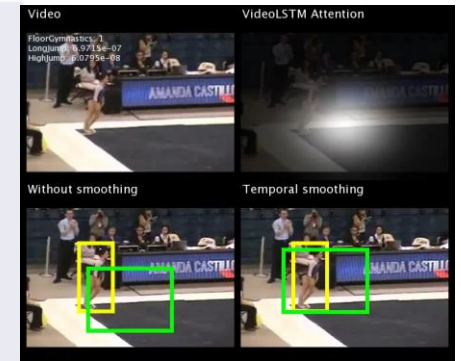
Random walks on space-time graphs



(a) Balance beam



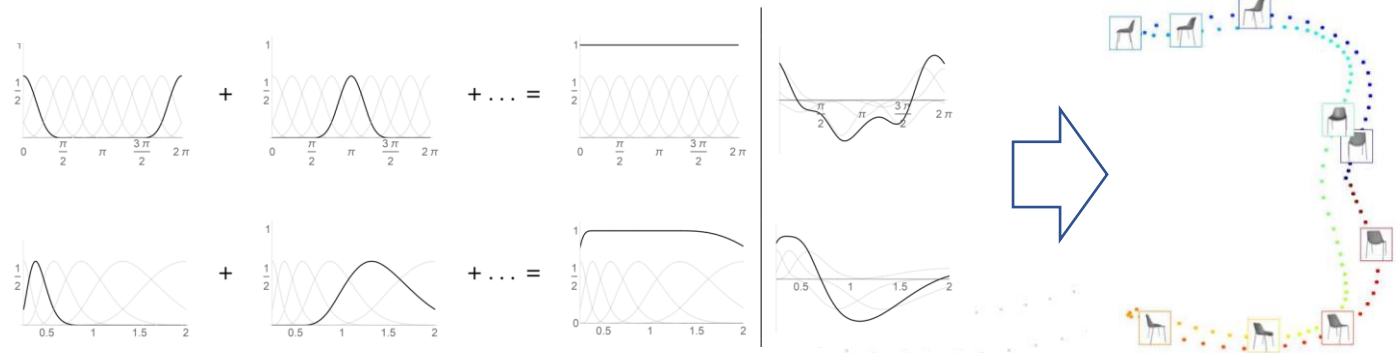
(b) Floor exercise



[1]

● swing ● flip ● flares ● handstands ● twist ● balance

Space-time manifolds



VideoGraph: Recognizing Minutes-Long Human Activities in Videos



N. Hussein

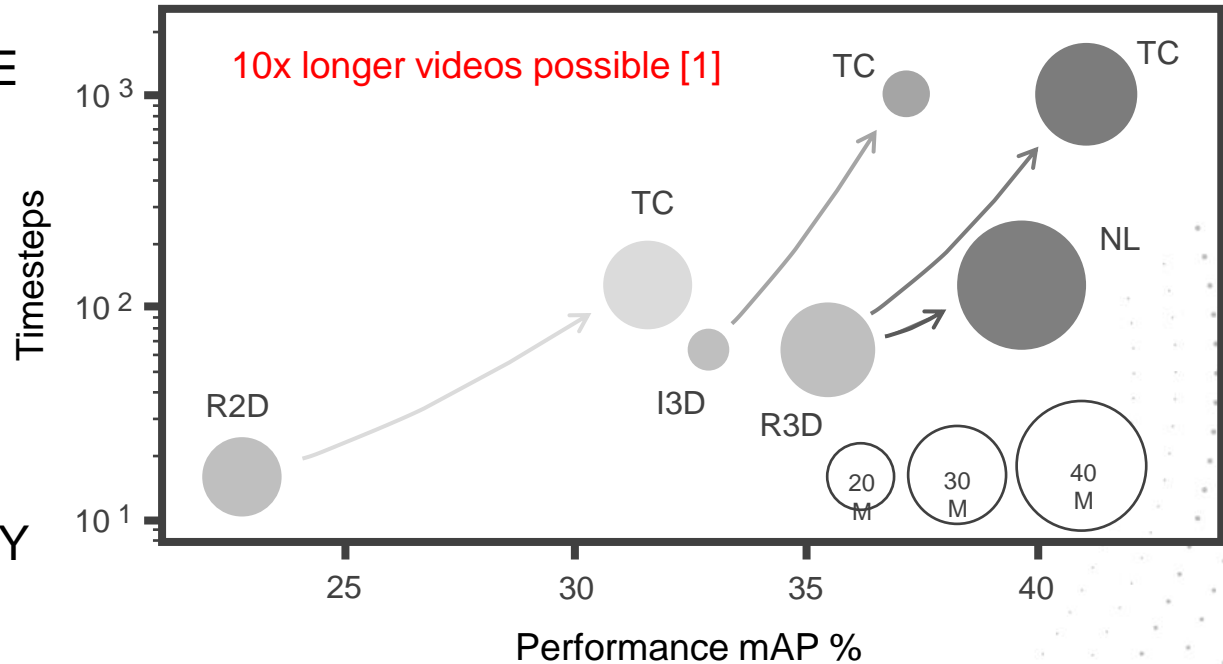


E. Gavves



A. Smeulders

Temporal length is important with complex videos



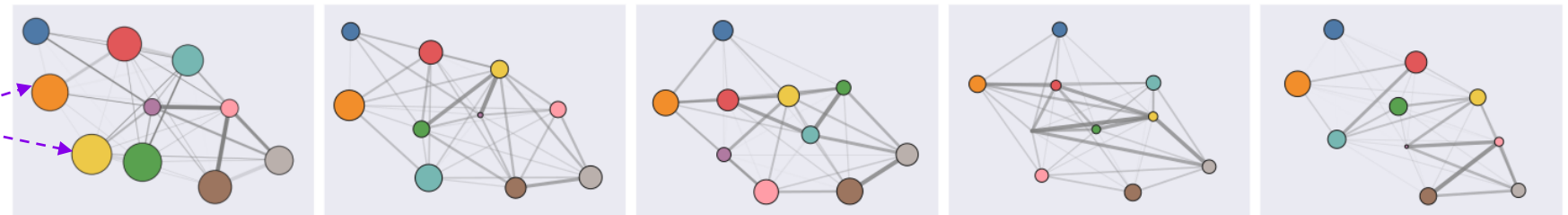
Code available



How to stretch time further → (Time) Graphs!

- Sublinear temporal representation
- Compositionality
- Interpretability

Characteristic one-actions



(a) Making Cereals

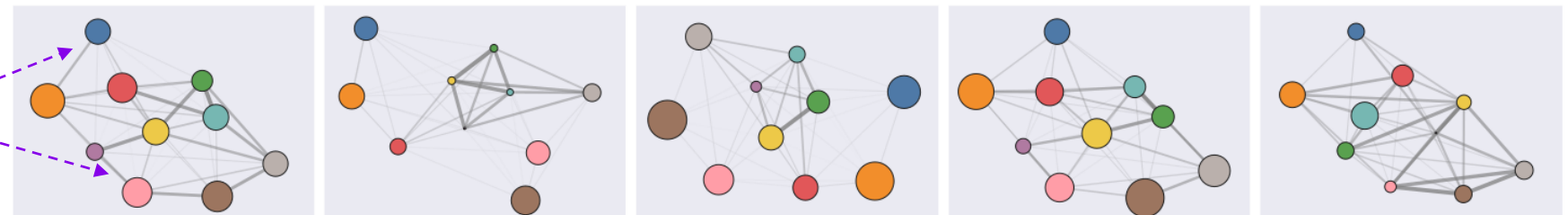
(b) Preparing Coffee

(c) Frying Eggs

(d) Making Juice

(e) Preparing Milk

Temporal relationships



(f) Making Pancake

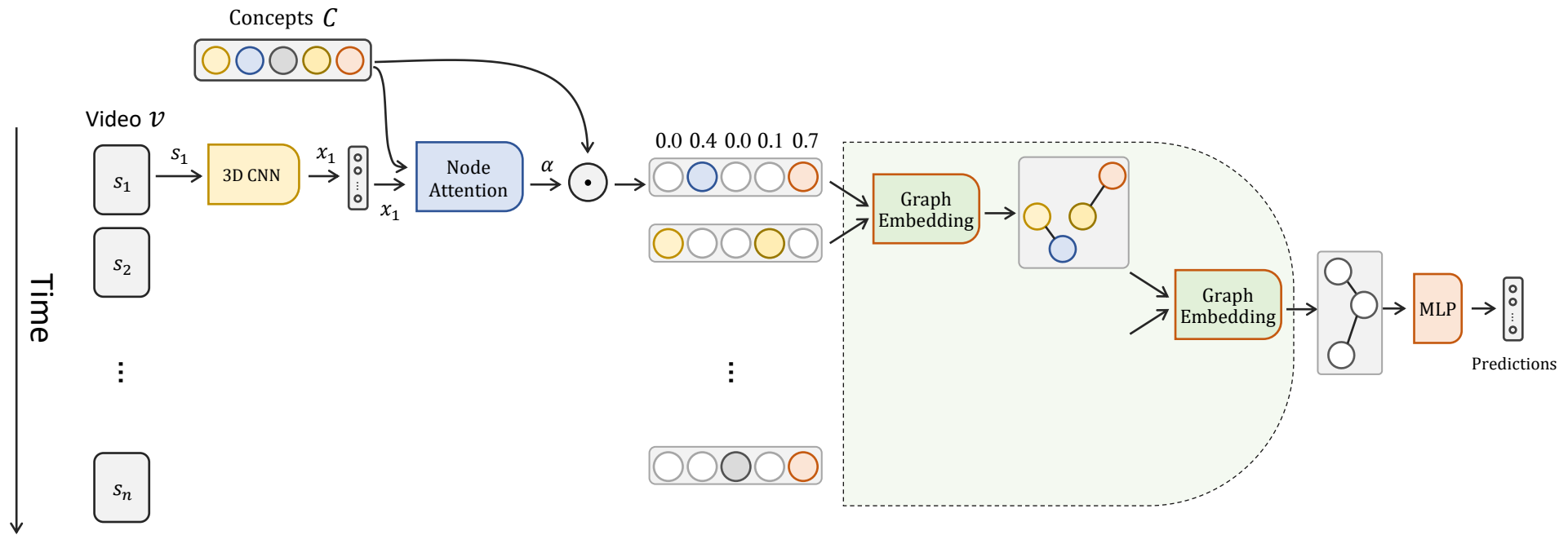
(g) Making Salat

(h) Making Sandwich

(i) Making Scrambled Egg

(j) Preparing Tea

VideoGraph



Some results

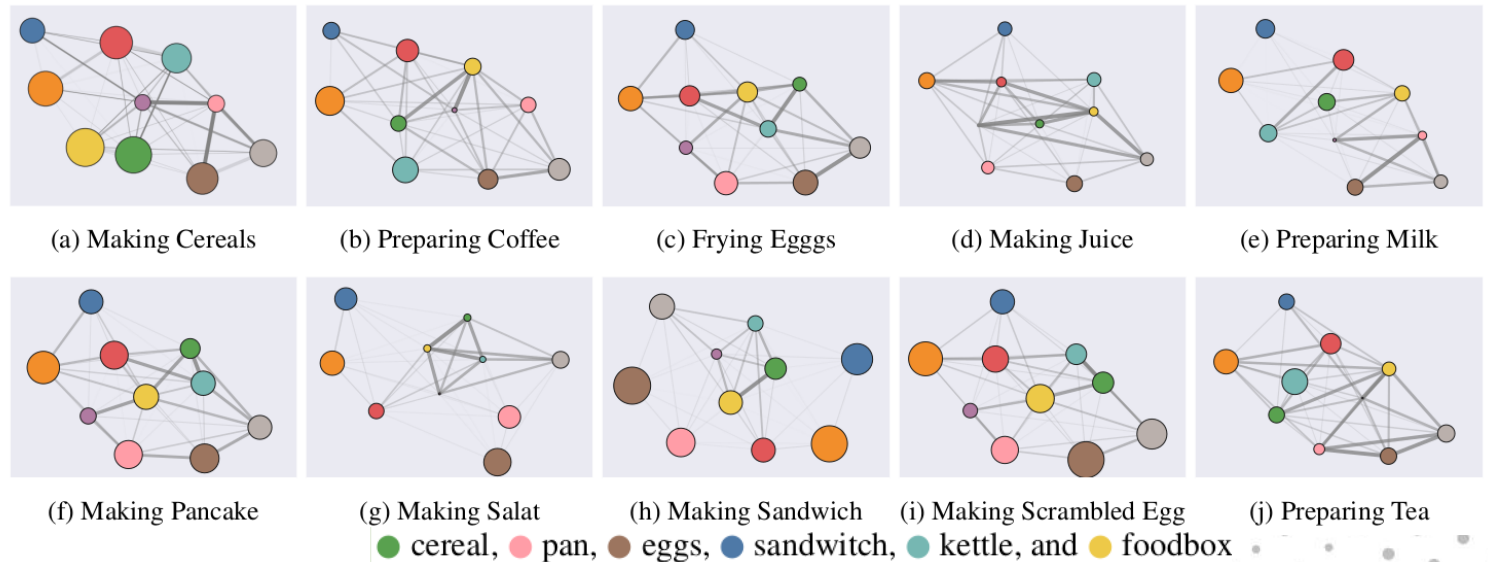
Charades

Method	Modality	mAP (%)
Two-stream	RGB + Flow	18.6
Two-stream + LSTM	RGB + Flow	17.8
ActionVLAD	RGB + iDT	21.0
Temporal Fields	RGB + Flow	22.4
Temporal Relations	RGB	25.2
ResNet-152	RGB	22.8
ResNet-152 + Timeception	RGB	31.6
I3D	RGB	32.9
I3D + ActionVLAD	RGB	35.4
I3D + Timeception	RGB	37.2
I3D + VideoGraph	RGB	37.8

Breakfast

Method	Breakfast Acc. (%)	Breakfast mAP (%)
ResNet-152	41.13	32.65
ResNet-152 + ActionVLAD	55.49	47.12
ResNet-152 + Timeception	57.75	48.47
ResNet-152 + VideoGraph	59.12	49.38
I3D	47.05	58.61
I3D + ActionVLAD	60.20	65.48
I3D + Timeception	61.82	67.07
I3D + VideoGraph	63.14	69.45

Code available



Categorical Normalizing Flows via Continuous Transformations



P. Lippe



E. Gavves

From classifying to generating graphs

- What if we would like to create new, plausible graphs?
- Normalizing Flows is State-of-the-art in generative modelling

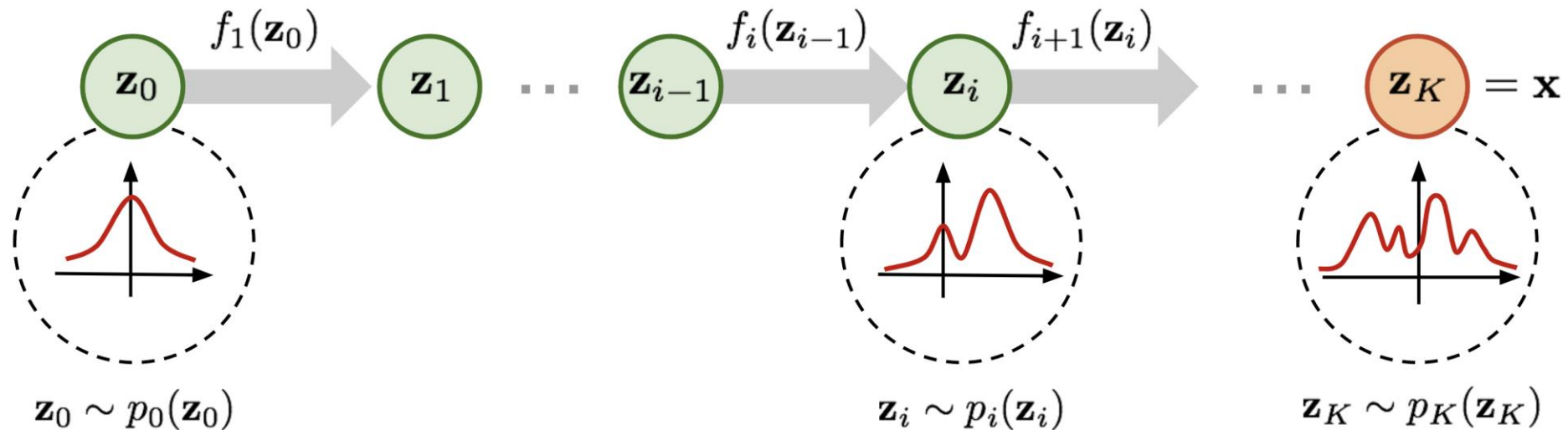
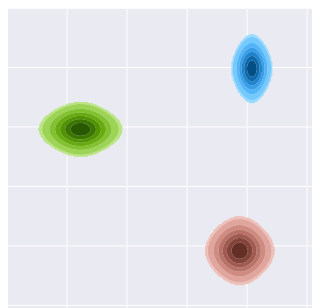


Figure credit: Weng, Lilian. "Flow-based Deep Generative Models", 2018.

- + Universality
- + Exact likelihood estimate
- + Efficient density evaluation and (parallel) sampling
- Does not work on categorical data seamlessly

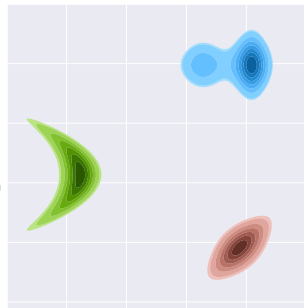
Step 1: Categorical Normalizing Flows

- Learn encoder to represent categorical data in continuous space
 - Must not lose information in the representation
 - Must be smooth and have support for higher dimensions



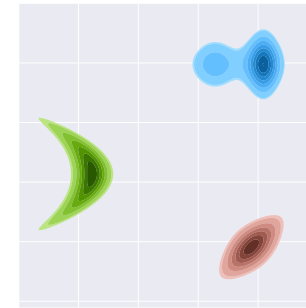
Mixture model

$$q(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^N g(\mathbf{z}_i|\mu(x_i), \sigma(x_i))$$



Linear flows

$$q(\mathbf{z}|\mathbf{x}) = \prod_{i=1}^N q(\mathbf{z}_i|x_i)$$



Variational encoding

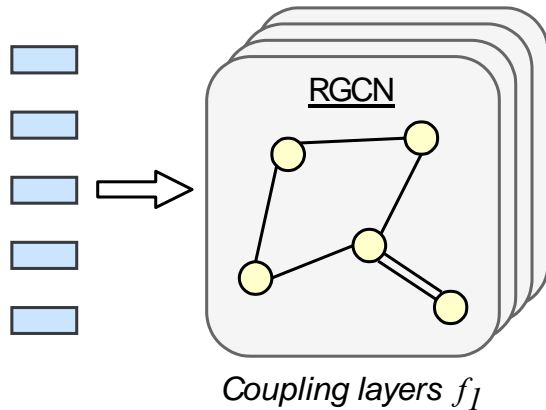
$$q(\mathbf{z}|\mathbf{x})$$

- Learn decoder s.t. continuous \mathbf{z} contains what is in \mathbf{x} exactly
 - Variational inference with factorized decoder

$$p(\mathbf{x}) \geq \mathbb{E}_{\mathbf{z} \sim q(\cdot|\mathbf{x})} \left[\frac{\prod_i p(x_i|\mathbf{z}_i)}{q(\mathbf{z}|\mathbf{x})} p(\mathbf{z}) \right]$$

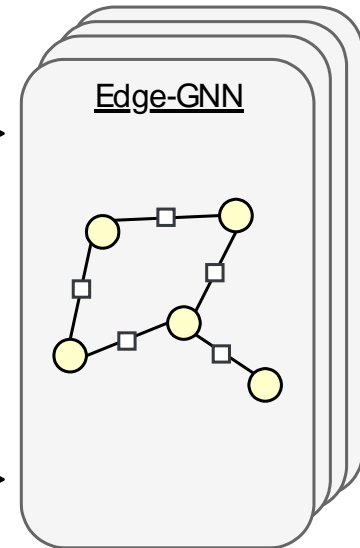
Step 2: Graph generation with CNF

CNF - Node type representation
(discrete \rightarrow continuous)



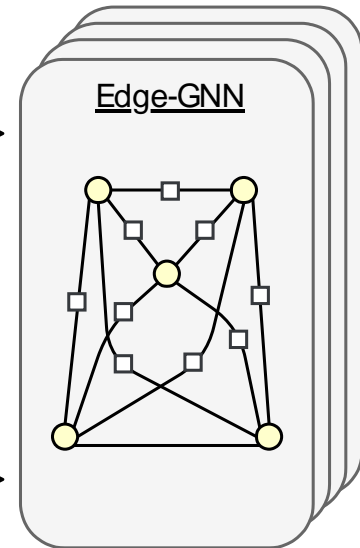
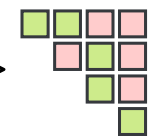
Coupling layers f_1

CNF - Edge attribute representation
(discrete \rightarrow continuous)



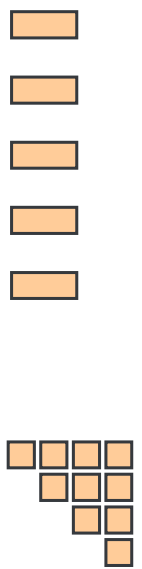
Coupling layers f_2

Adding virtual edge
representation
(CNF)



Coupling layers f_3

Prior distribution



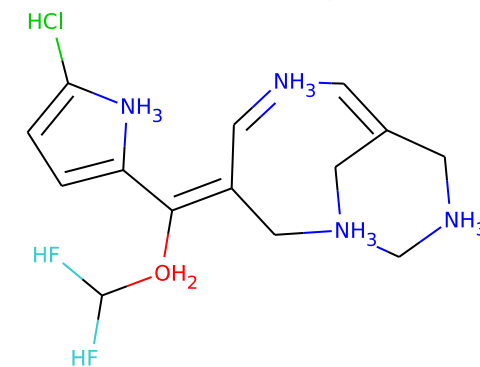
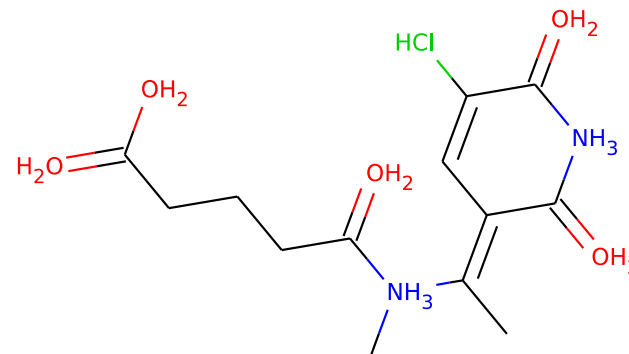
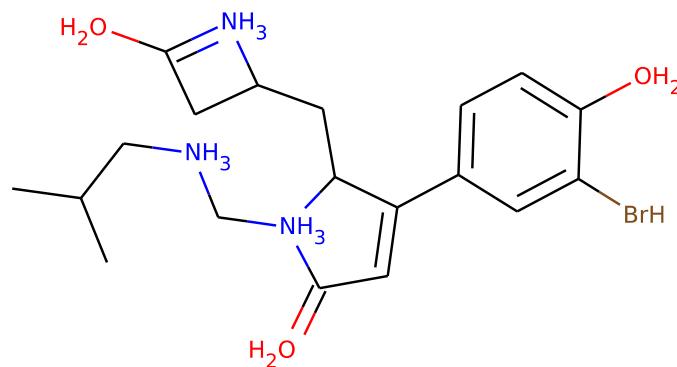
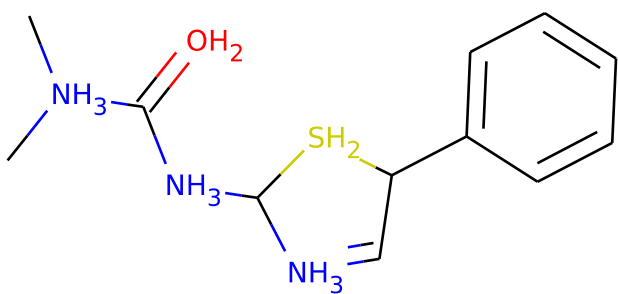
- + Permutation-invariant
- + Efficient three-step approach

Some results

Code available

Molecule generation with Zinc250k dataset (224k examples)

Method	Validity	Uniqueness	Novelty	Reconstruction	Parallel	General
JT-VAE	100%	100%	100%	71%	✗	✗
GraphAF	68%	99.10%	100%	100%	✗	✓
R-VAE	34.9%	100%	—	54.7%	✓	✓
GraphNVP	42.60%	94.80%	100%	100%	✓	✓
GraphCNF	83.41% (±2.88)	99.99% (±0.01)	100% (±0.00)	100% (±0.00)	✓	✓
+ Sub-graphs	96.35% (±2.21)	99.98% (±0.01)	99.98% (±0.02)	100% (±0.00)	✓	✓



Rotation Equivariant Siamese Networks for Tracking



D. Gupta



D. Arya



E. Gavves

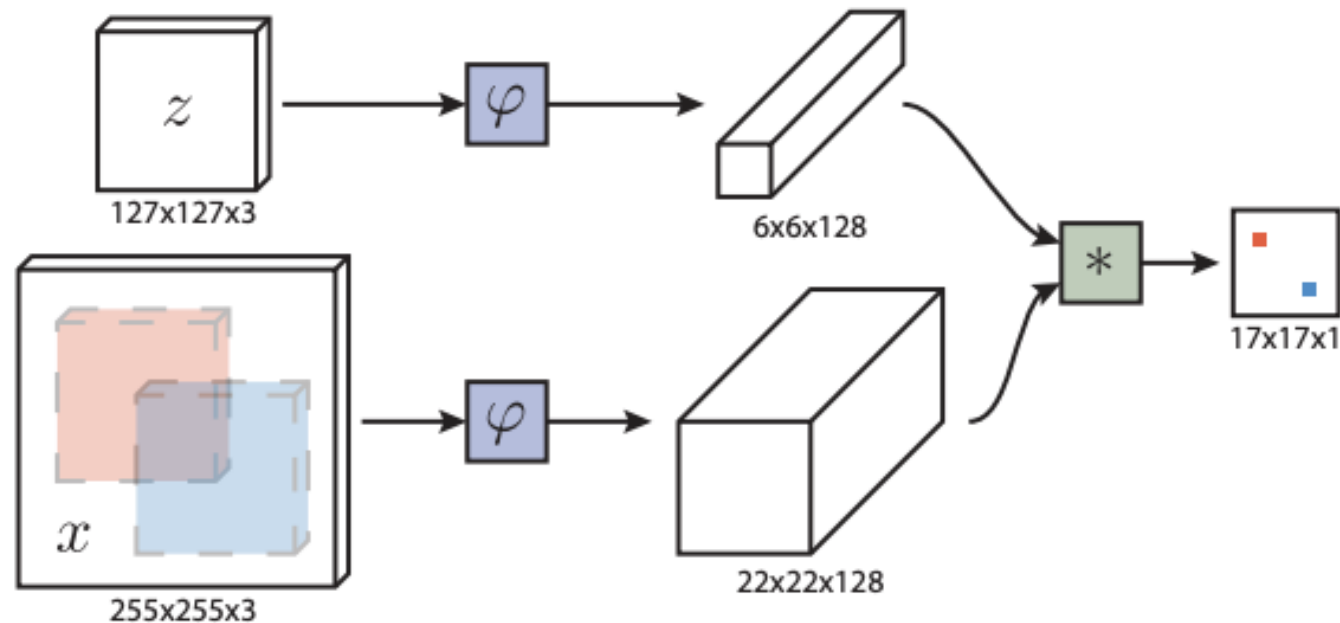
In-plane rotations in tracking

- Drone, surveillance, ego-motion recordings, etc.



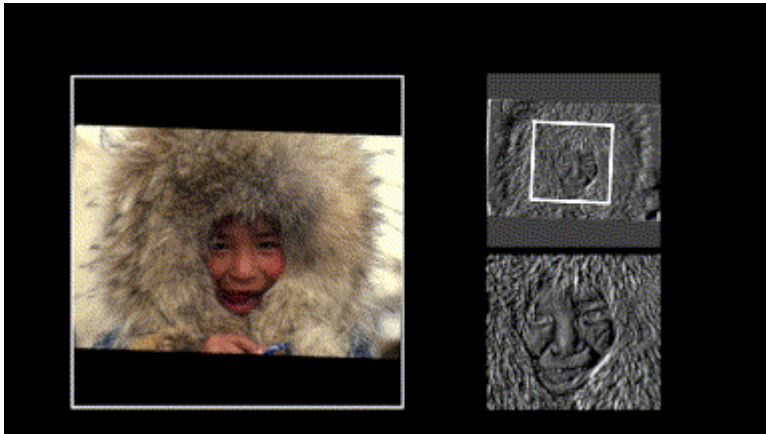
SoTA: Siamese Trackers [1, 2]

- Tracking as matching the target query to per frame instances
- Sensitive to rotations ← Convolutions are sensitive to rotations

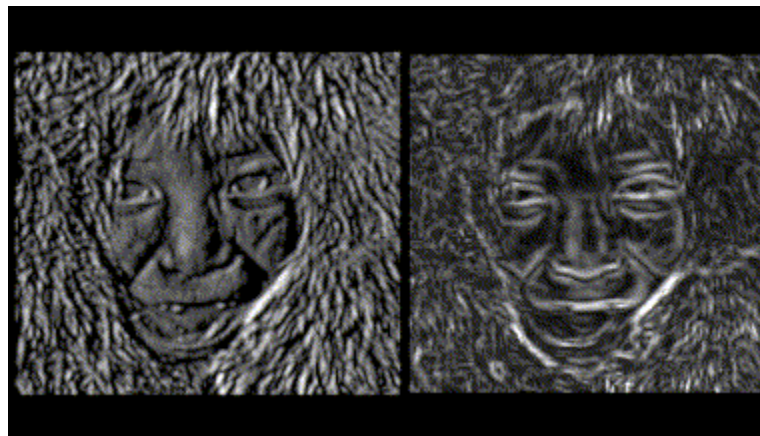
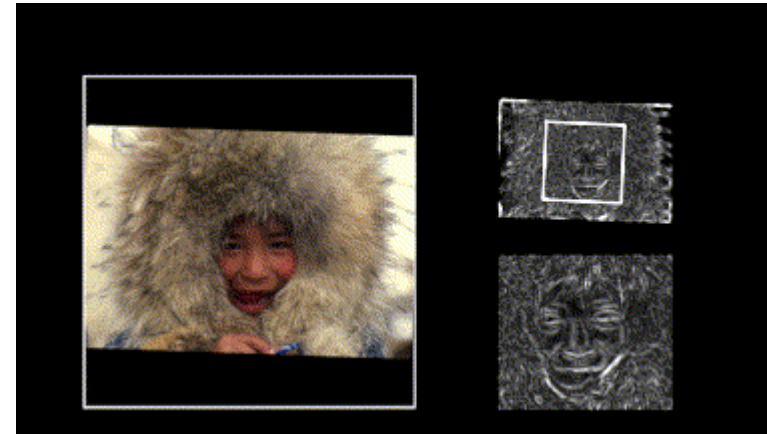


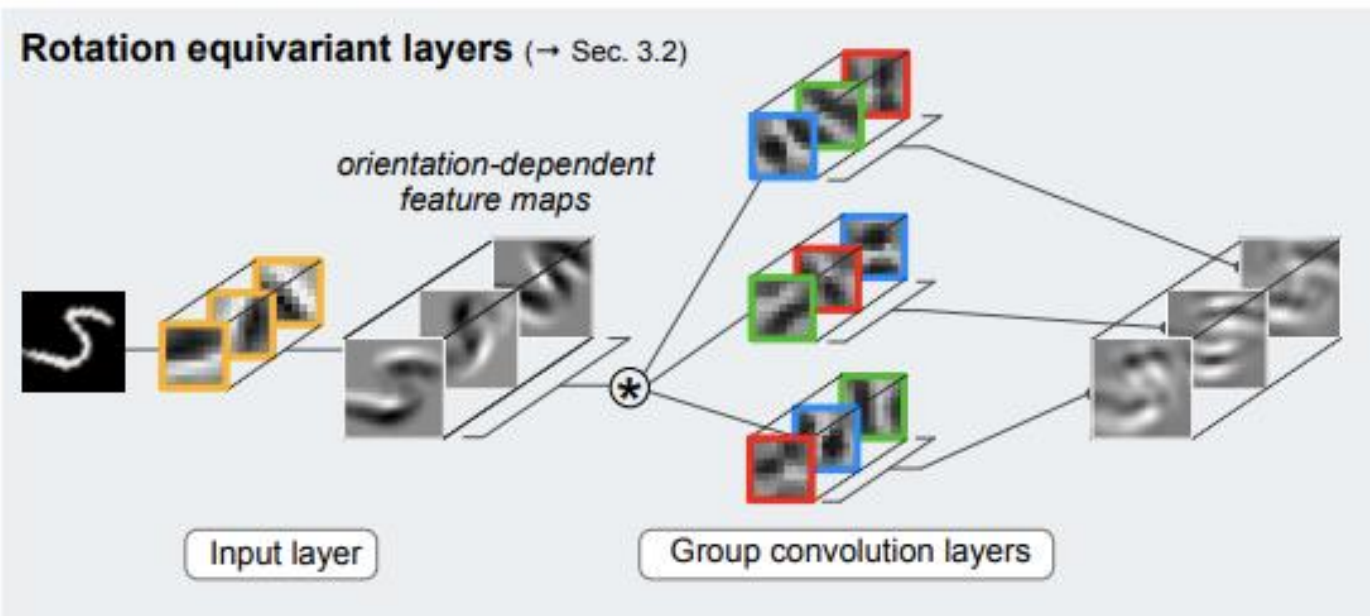
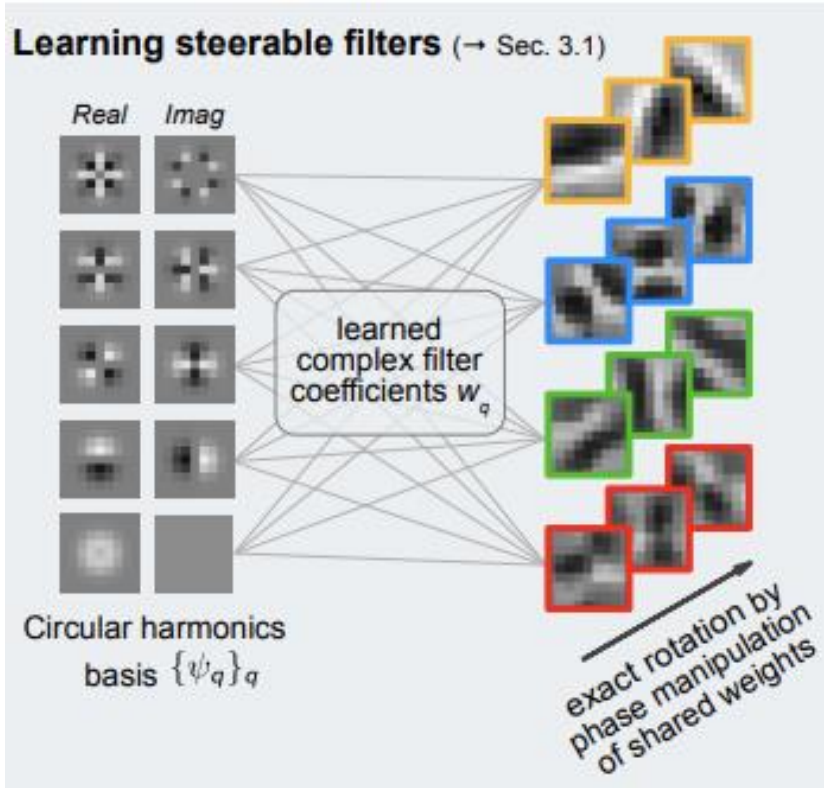
Rotation equivariance in CNNs

Regular CNN



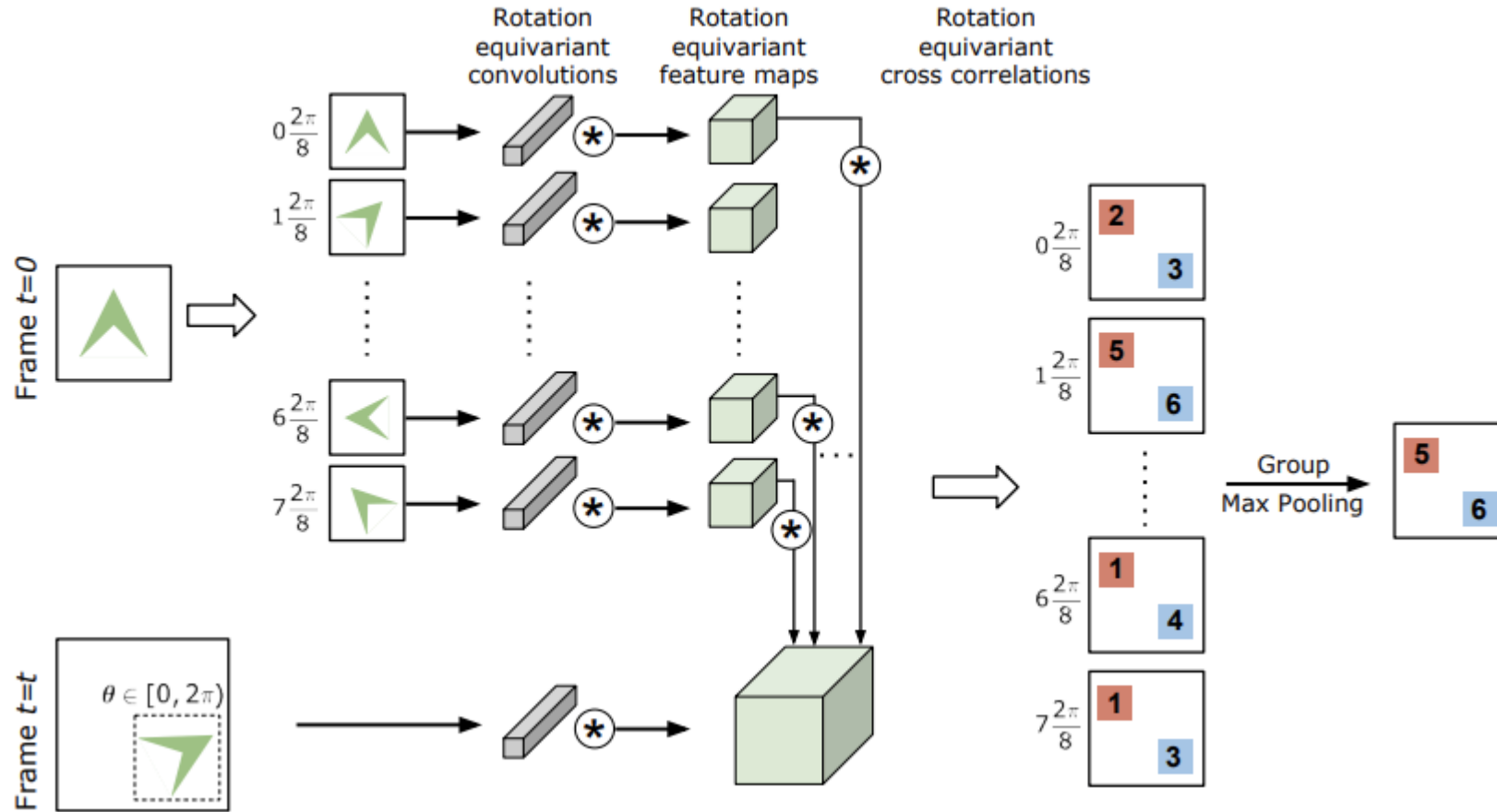
Rotation Equivariant Net
(H-Net; Worrall et al, CVPR 2017)





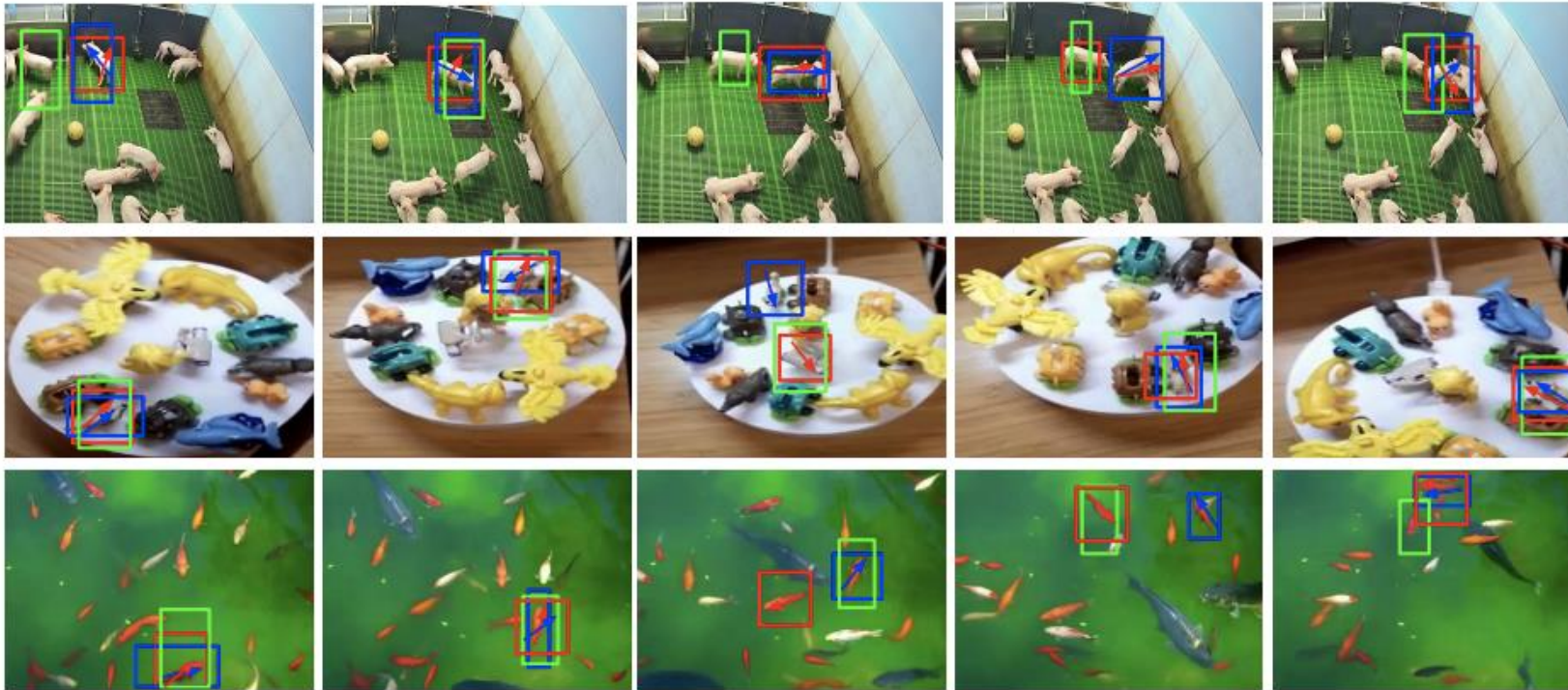
$$\left. \begin{aligned} \psi_{jk}(r, \phi) &= \tau_j(r) e^{ik\phi} \\ \rho_\theta \psi_{jk}(x) &= e^{-ik\theta} \psi_{jk}(x) \end{aligned} \right\} \rho_\theta \Psi(x) = \sum_{j=1}^J \sum_{k=0}^K w_{jk} e^{-ik\theta} \psi_{jk}(x)$$

Rotation Equivariant Siamese Trackers



Some results

Code soon available



Ground-truth, SiamFC, RE-SiamFC

Model	Type	Succ	Pr
SiamFC [1]	-	0.315	0.523
	R4	0.360	0.629
	R8	0.423	0.676
SiamFCv2	-	0.288	0.473
	R4	0.348	0.622
	R8	0.425	0.678
	R16	0.423	0.688
SiamFCv2	aug	0.317	0.541
SiamRPN++ [18]	-	0.461	0.634
SiamRPN++	R4	0.485	0.679
DiMP18 [2]	-	0.429	0.643
DiMP50 [2]	-	0.447	0.668



Time Supervision

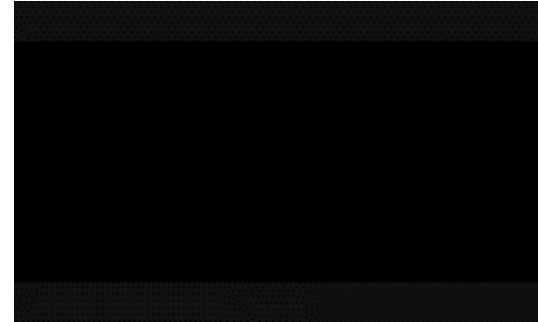


(Some) properties of time

Temporal
Asymmetry



Temporal
Causality

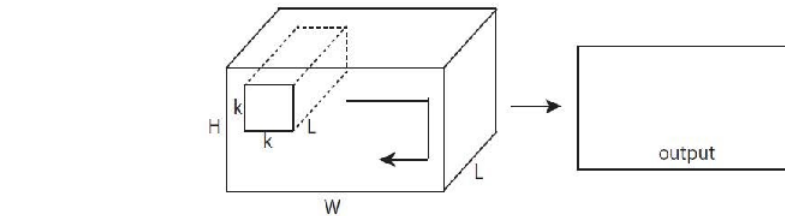
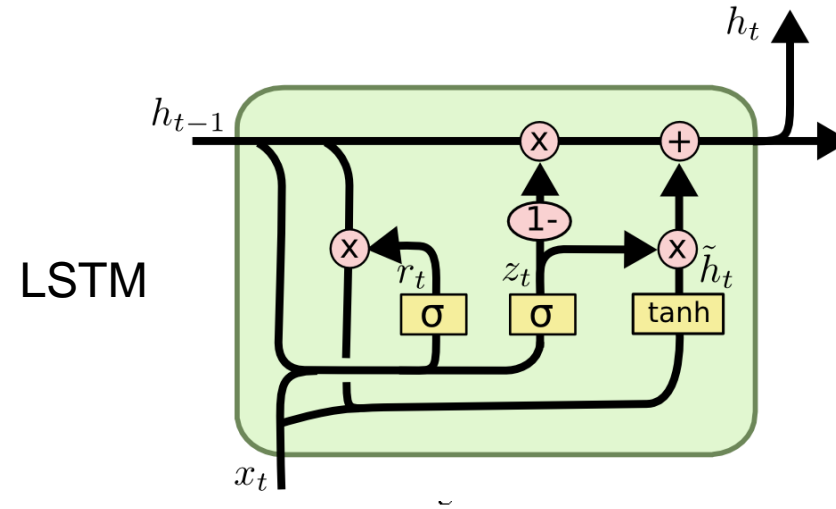
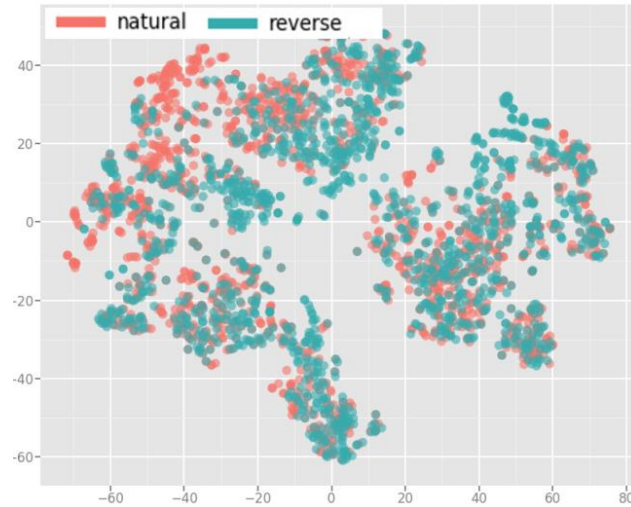


Temporal
Continuity



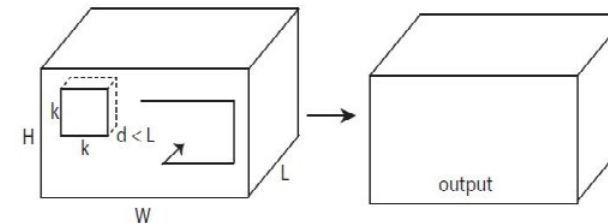
Temporal
Redundancy

Arrow of time: LSTM vs C3D



(a) 2D convolution on multiple frames

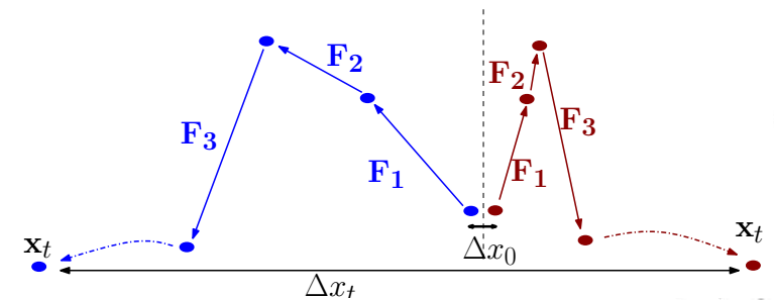
C3D



(b) 3D convolution on multiple frames

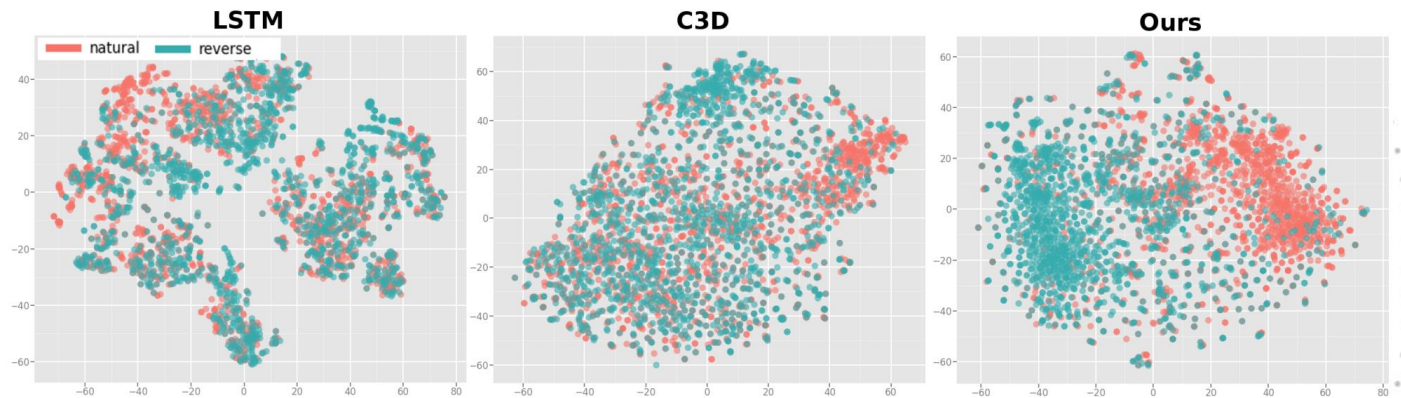
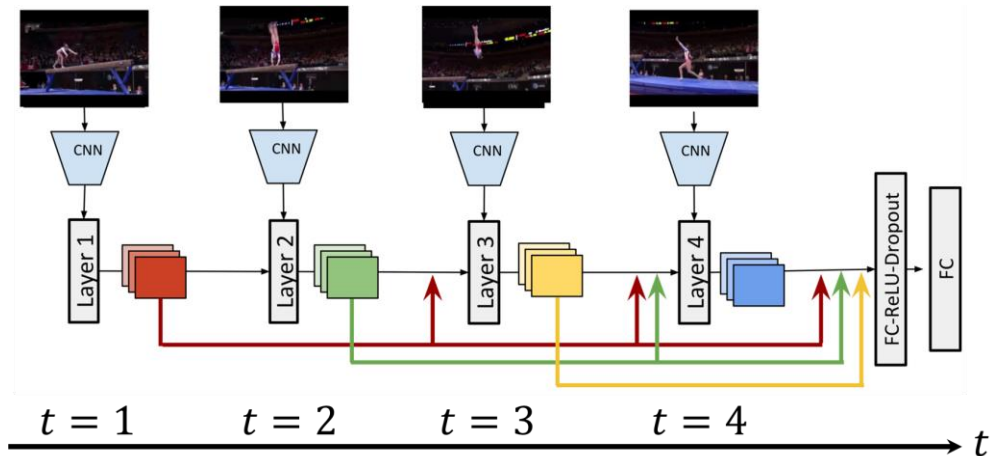
Revisiting recurrent neural networks

- Recurrent Nets are highly sensitive dynamical systems [1]
 - Even with highly discriminative symbolic (one-hot vector) inputs
 - Gradients very sensitive to initialization → Poor learning! → No generalization
- Visual features are
 - much noisier, less discriminative, much more redundant
- Learning LSTM on videos is orders of magnitude harder
 - Chaotic regime → no useful gradients → no learning
 - Forward/Backward/Shuffling of frames → LSTM performs the same on arrow of time



Time-aligned neural networks [1]

- Idea: Why not flip the ConvNet to align the layers with time steps?
- No vanishing/exploding gradients, no problems with noise/redundancy



Conclusion: Poor temporal modelling could be due to hard –and thus unsuccessful- optimization



Time Evaluation



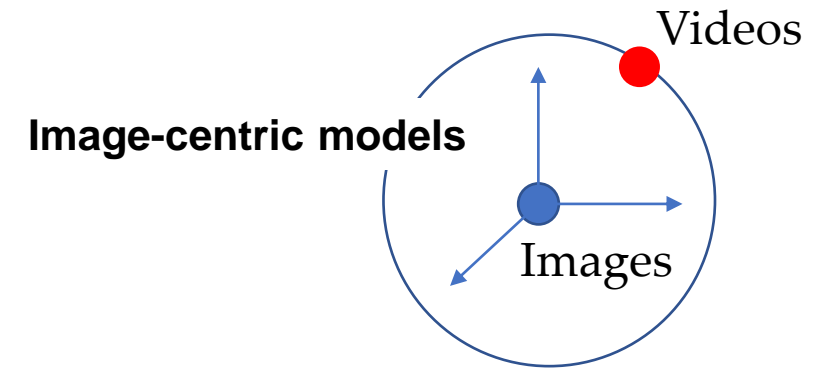
Engage-Generalize-Scale

- Current validation paradigm in video is hard to sustain
- Models for (n -dimensional + time) signals, e.g. scientific recordings
 - Particles through time, climate, astronomy, ecosystems, biology
- Possible great advantage: scientific knowledge as groundtruth



Conclusion

- Time largely ignored in model building and validation



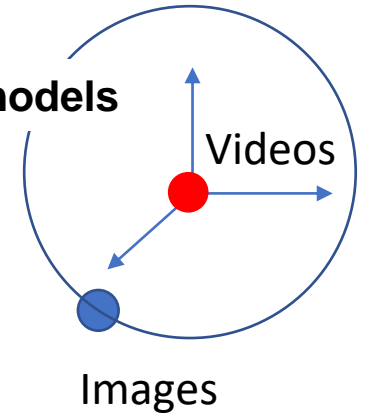
Conclusion

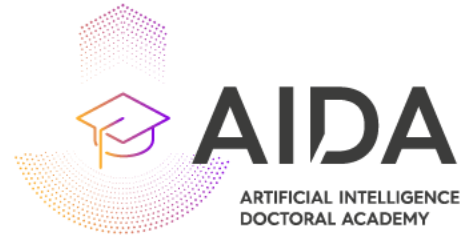
- Time largely ignored in model building and validation

Static → Temporal

- Geometry, generative & time supervision will be the key
- Hopefully, impact on any field with spatiotemporal complex data

Spatiotemporal-centric models





The International AI Doctoral Academy (**AIDA**) is a joint initiative of the European R&D Projects: **AI4media**, **VISION** and **HumanE-AI Net**



www.ai4media.eu



www.vision4eu.eu



www.humane-ai.eu

These projects have received funding from the European Union's Horizon 2020 research and innovation programme under the following Grant agreements: No 951911 (AI4media), No. 952070 (VISION) and No. 952026 (HumanE-AI Net)

