

# Structured Deep Learning for Video Analysis

Fabien Baradel  
PhD Candidate

Advisors: Christian Wolf & Julien Mille

June, 29th, 2020

# What is video understanding?



Human actions

→ Sitting on the floor

Entity-level interactions

→ Grabing a silencer

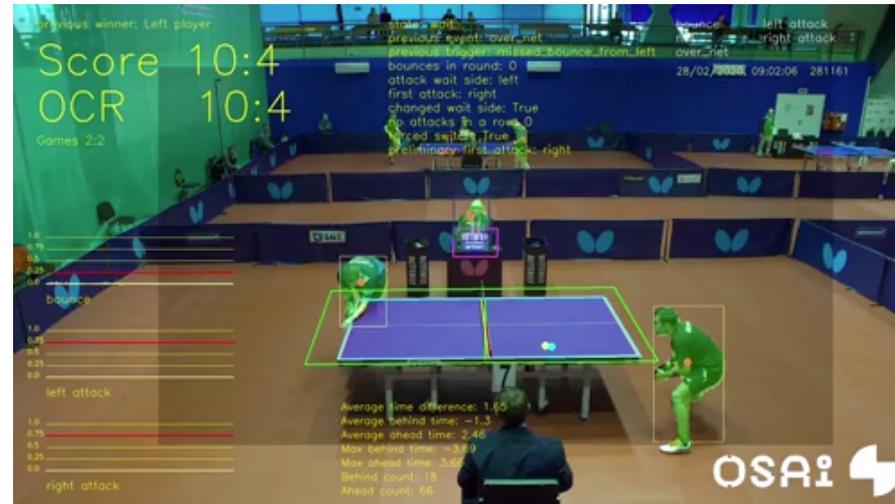
Temporal reasoning / Causality

→ The baby starts crying because of the silencer

# Why video understanding?



Indexing  
Retrieval



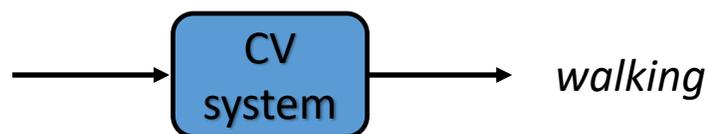
Recommendation  
Analysis



Human-robot interactions

# A video understanding task

## *Action Recognition*



Classification task  
Pre-defined labels  
Similar to object recognition

# Human Pose

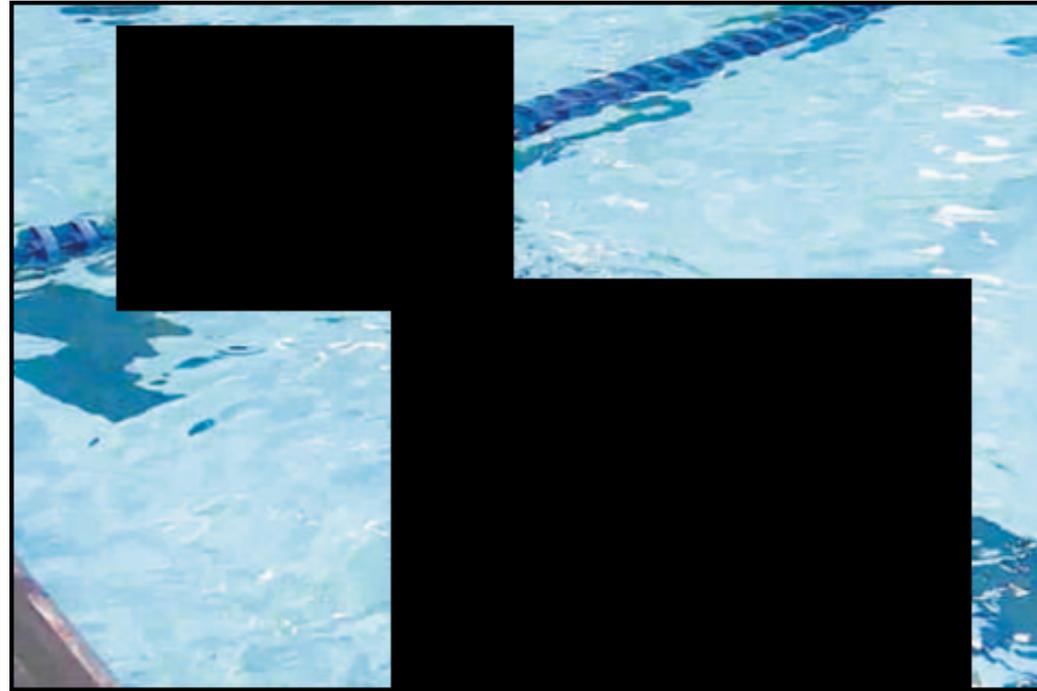


*Walking*  
Pose is enough



*Chopping*  
What is being chopped?

# Context & Appearance



*Swimming?*

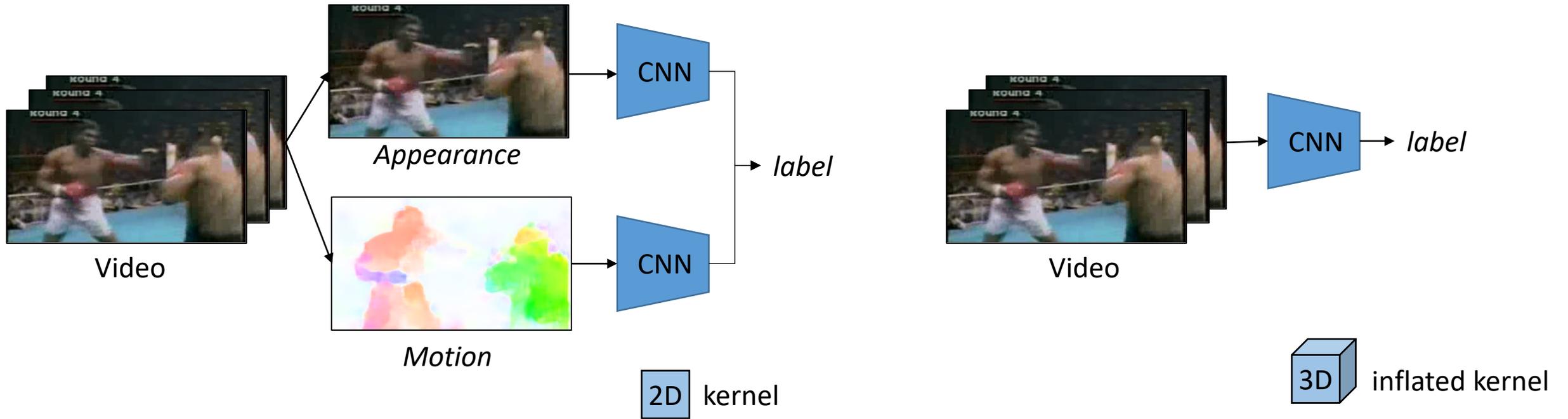
*Handshaking*

# Action Recognition

## *Recent works*

Two-stream

I3D



### Limitations

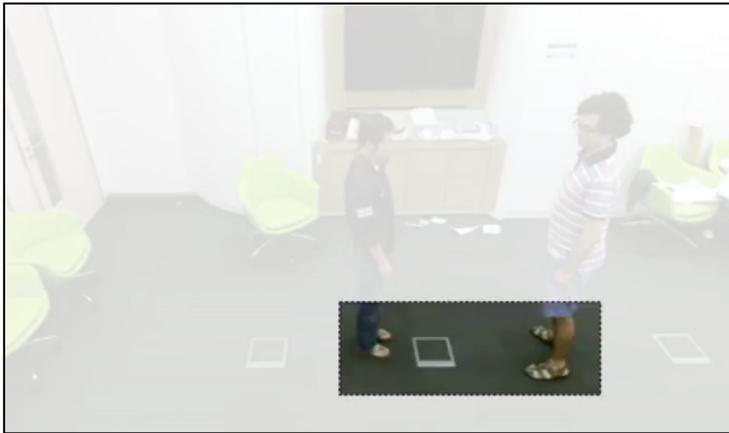
- Biased towards context
- Lack of explainability
- Human pose? Objects? Scene?

[Two-stream, Simonyan, NIPS'16]

[I3D, Carreira et al, CVPR'18]

# Structured Deep Learning

# Outline



## Visual Attention



**Christian Wolf**  
INSA Lyon - LIRIS

Glimpse Clou  
*F. Baradel, C. Wolf,*  
*G. Taylor*  
CVPR'18

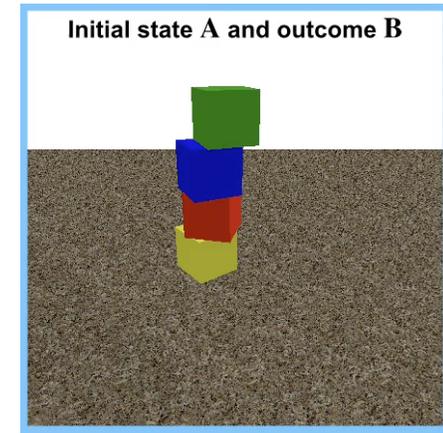


**Julien Mille**  
INSA CVL - LI Tours



## Entity-level interactions

« Object level Reasoning »  
*F. Baradel, N. Neverova, C. Wolf,*  
*J. Mille, G. Mori*  
ECCV'18



## Reasoning



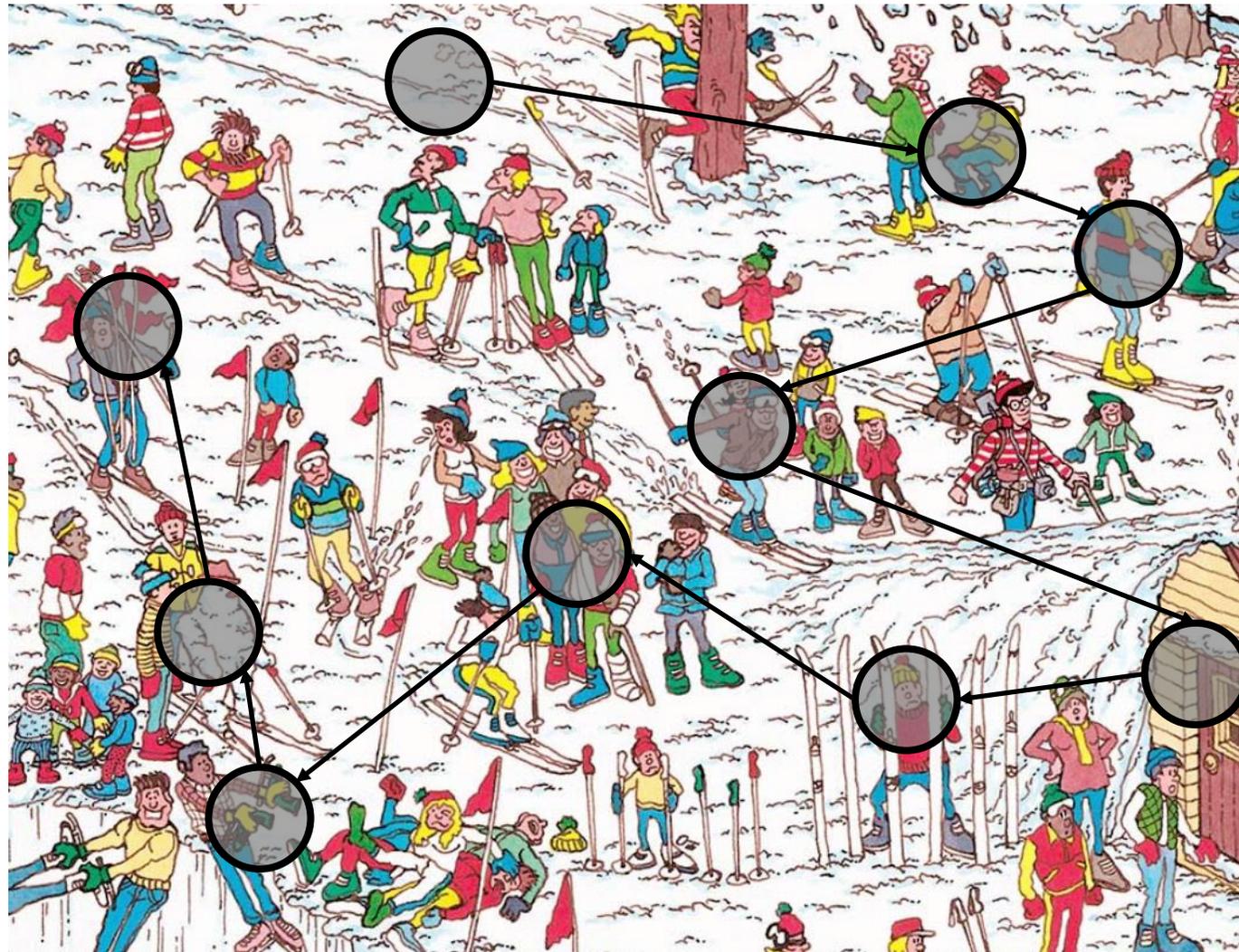
**Graham W. Taylor**  
University of Guelph  
Vector Institute

« Cross-factual learning »  
*F. Baradel, N. Neverova, J. Mille,*  
*J. Mori, C. Wolf*  
ICLR'20 (spotlight)

# Visual Attention

What is happening?

*Winter activities*

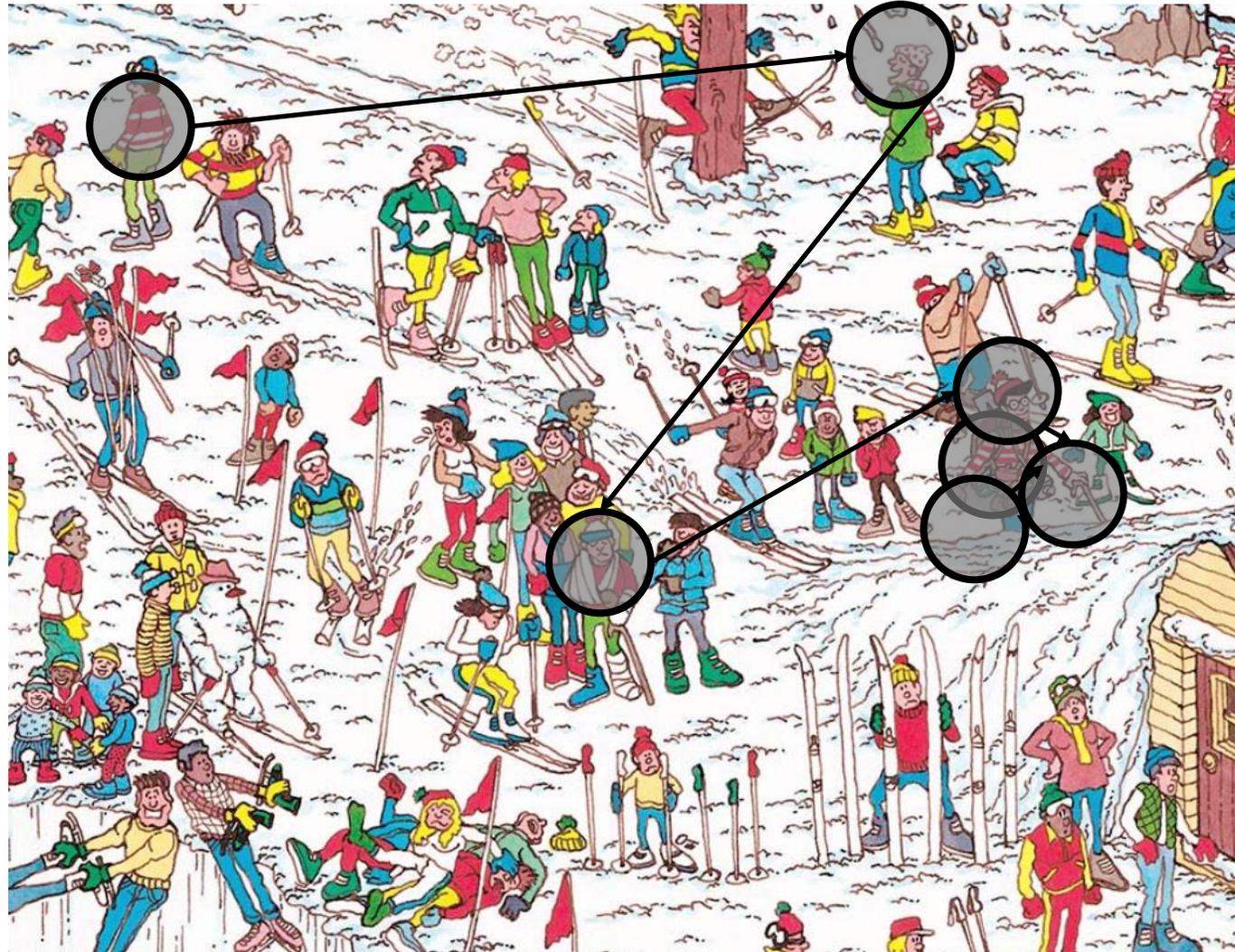


[Yarbus, 1976]  
[Roger et al, 2012]

# Visual Attention

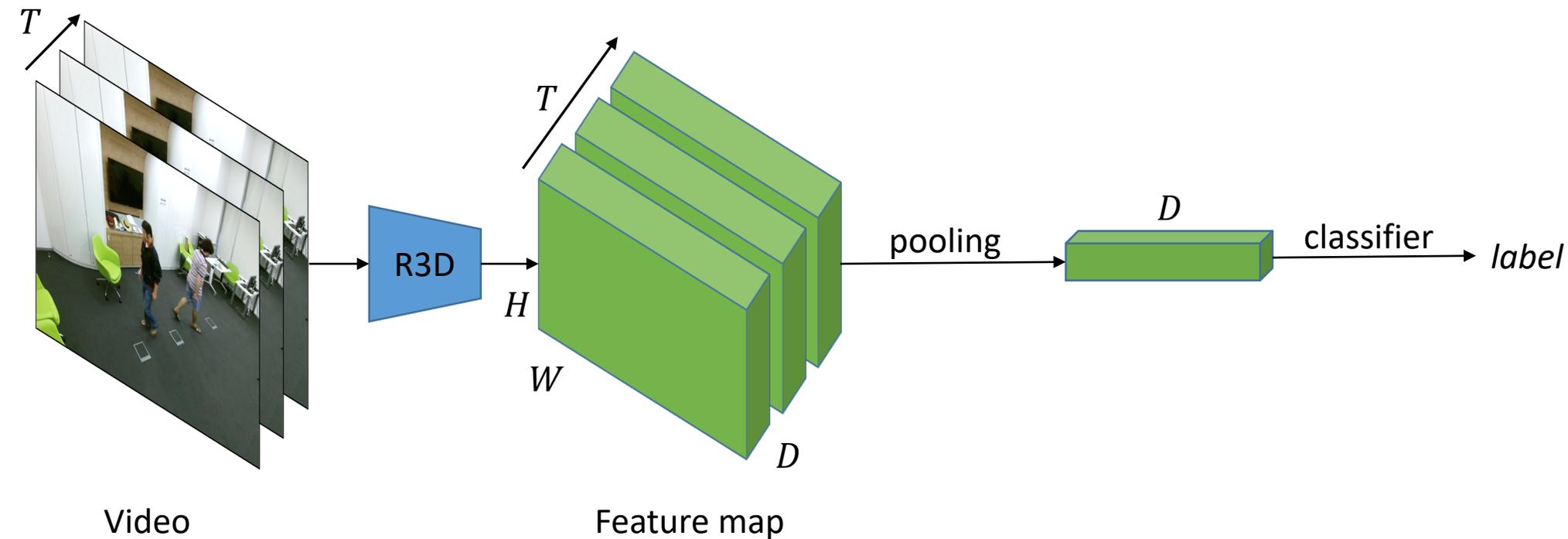
What is Charlie doing?

*Walking*



[Yarbus, 1976]  
[Roger et al, 2012]

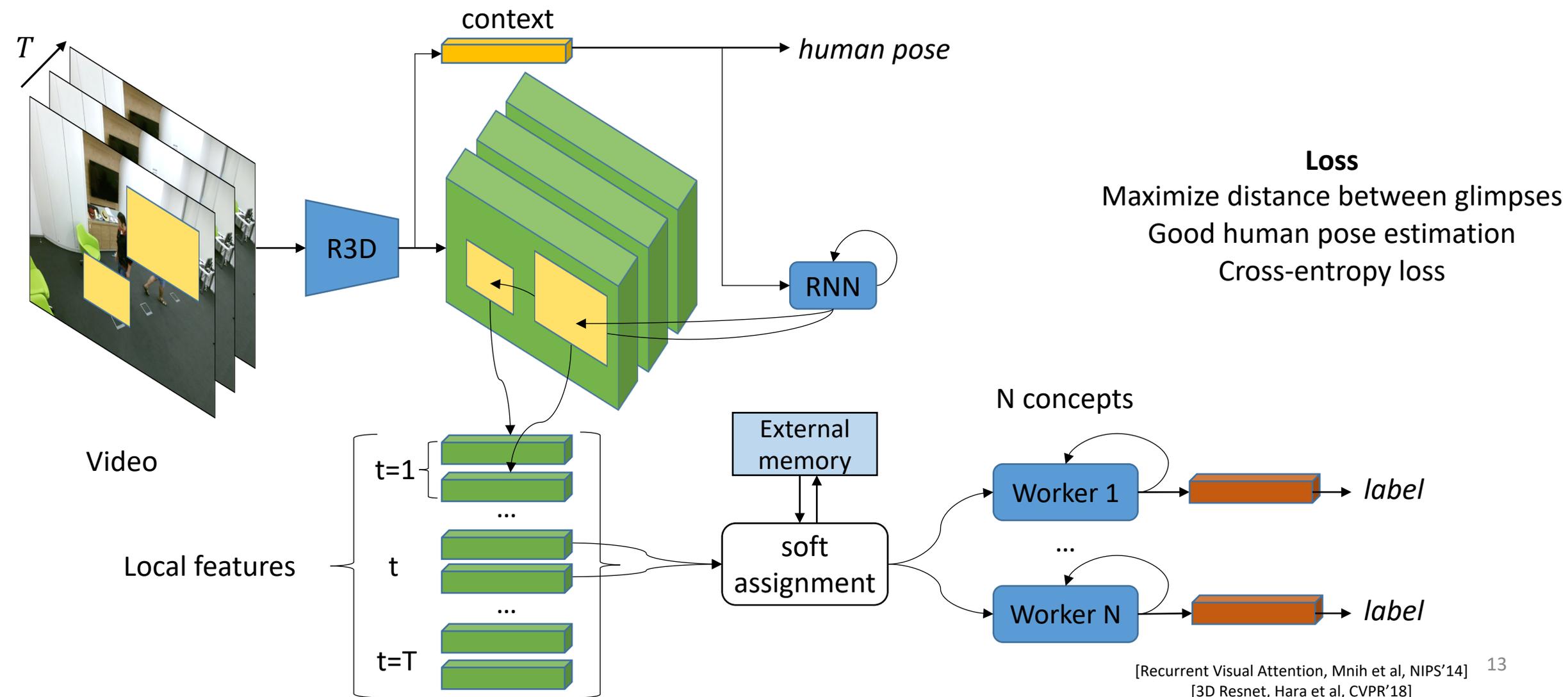
# Action Recognition *Baseline*



## Limitations

What about fine-grained human actions?  
How to focus on relevant parts of the video?

# Glimpse Clouds Method



# Glimpse Clouds

## *State-of-the-art results*

<i>Method</i>	<i>Modality</i>	<i>CS</i>	<i>CV</i>
Ensemble TS-LSTM	skeleton	74.6	81.3
View invariant	skeleton	80.0	87.2
Hands-Attention (ours)	skeleton+ RGB	84.8	90.6
<b>Glimpse Clouds (ours)</b>	<b>RGB</b>	<b>88.4</b>	<b>93.2</b>

*Accuracy on NTU-RGB+D*

<i>Method</i>	<i>Modality</i>	<i>V1</i>
Enhanced viz.	Skeleton	86.1
Ensemble TS-LSTM	Skeleton	89.2
NKTM	RGB	75.8
<b>Glimpse Clouds (ours)</b>	<b>RGB</b>	<b>90.1</b>

*Accuracy on Northwestern-UCLA*

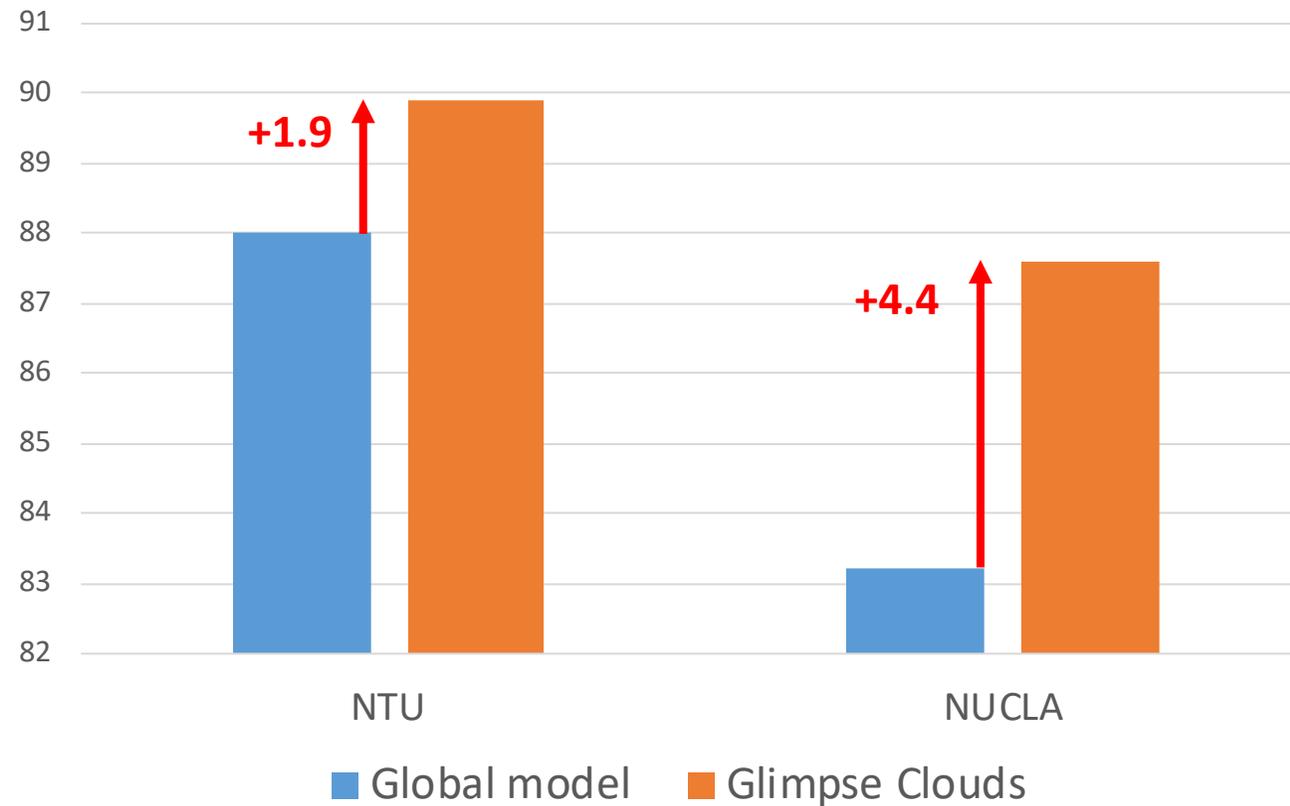


[fabienbaradel/glimpse\\_clouds](https://github.com/fabienbaradel/glimpse_clouds)

# Glimpse Clouds

## *Ablation study*

Impact of the attention mechanism



**Resolution matters**  
**Local fine-grained features**

# Glimpse Clouds *Visualization*



Raw video



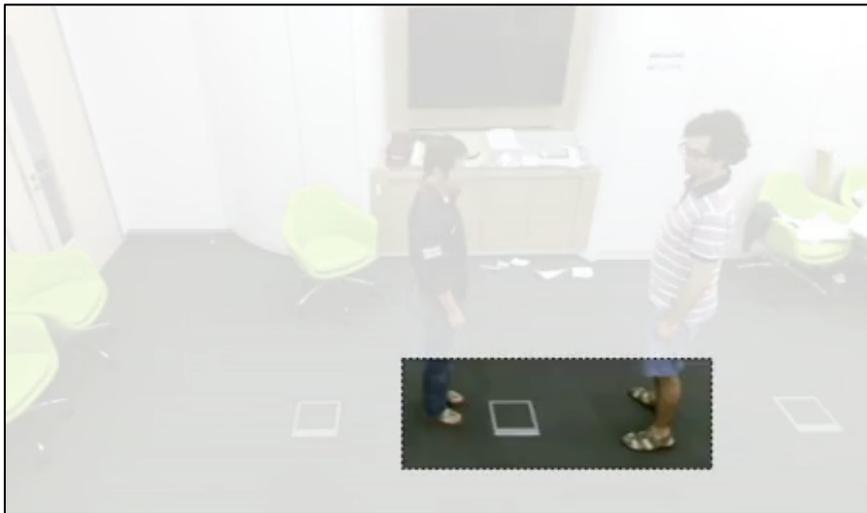
Attended regions

Worker 1 → ~Hands

Worker 2 → ~Heads

Worker 3 → ~Legs

# Outline



## Visual Attention



**Christian Wolf**  
INSA Lyon - LIRIS



**Julien Mille**  
INSA CVL - LI Tours

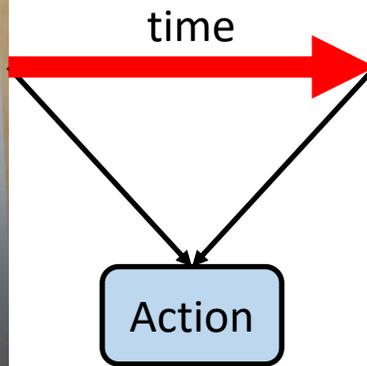
Unstructured local features...  
Incorporate structure from images?  
Leverage visual entities interactions?



## Entity-level interactions

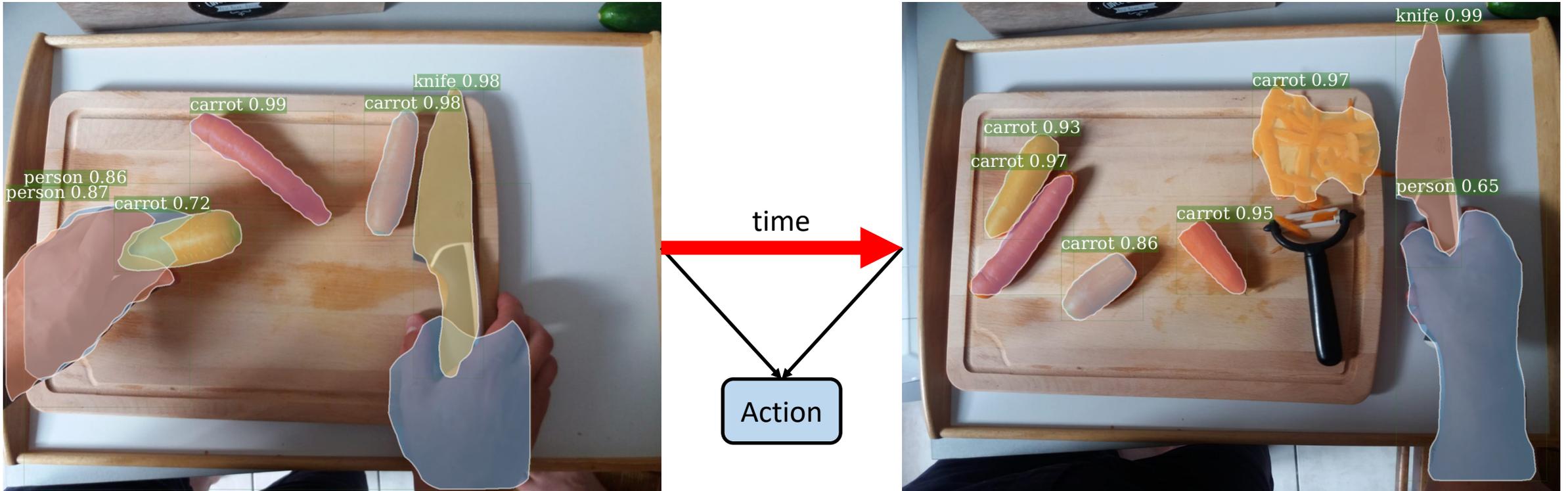


# Object-level Reasoning



Often possible to infer what happened from few frames

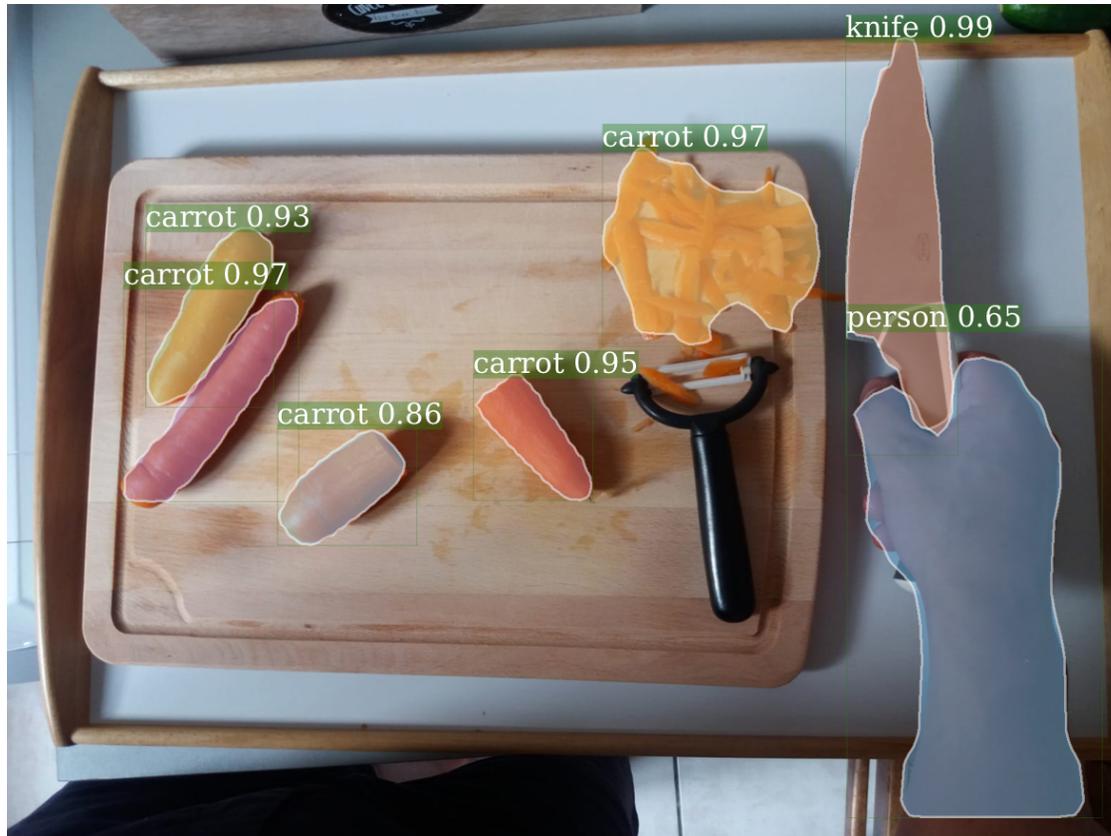
# Object-level Reasoning



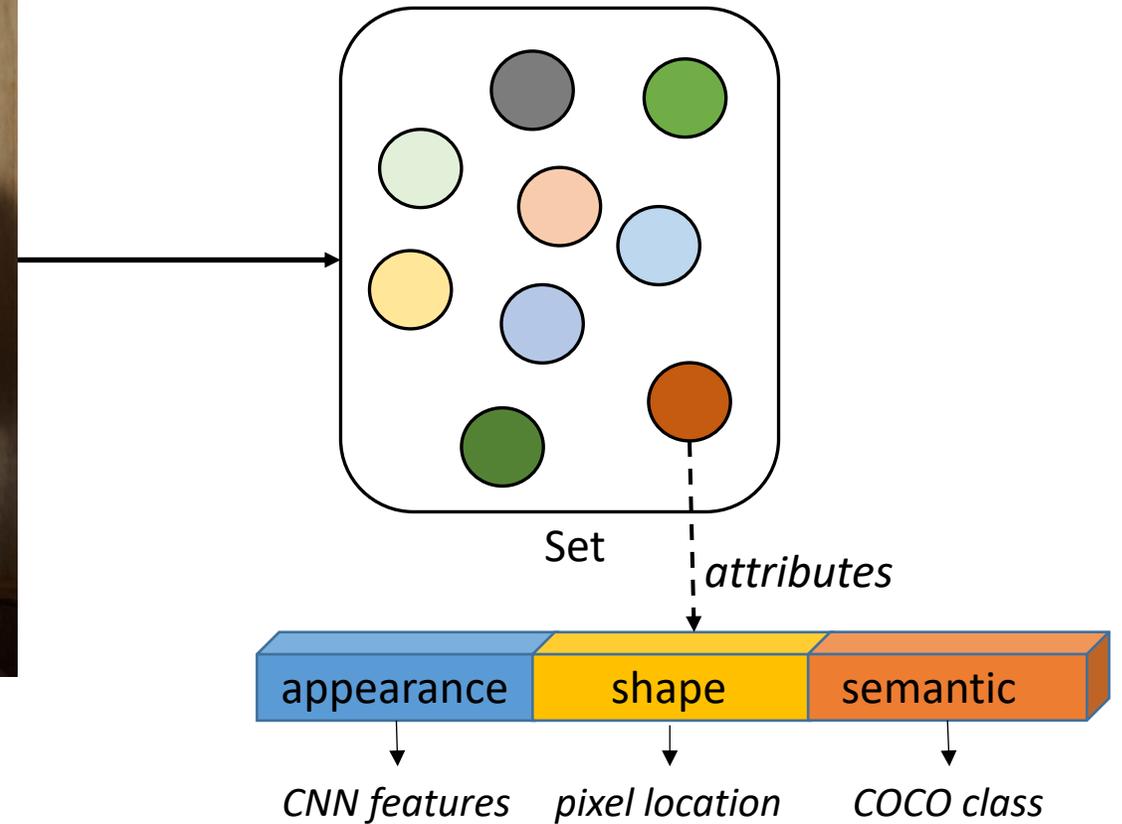
Often possible to infer what happened from few frames  
Visual entities interactions

[Mask-RCNN, He et al, ICCV'17]

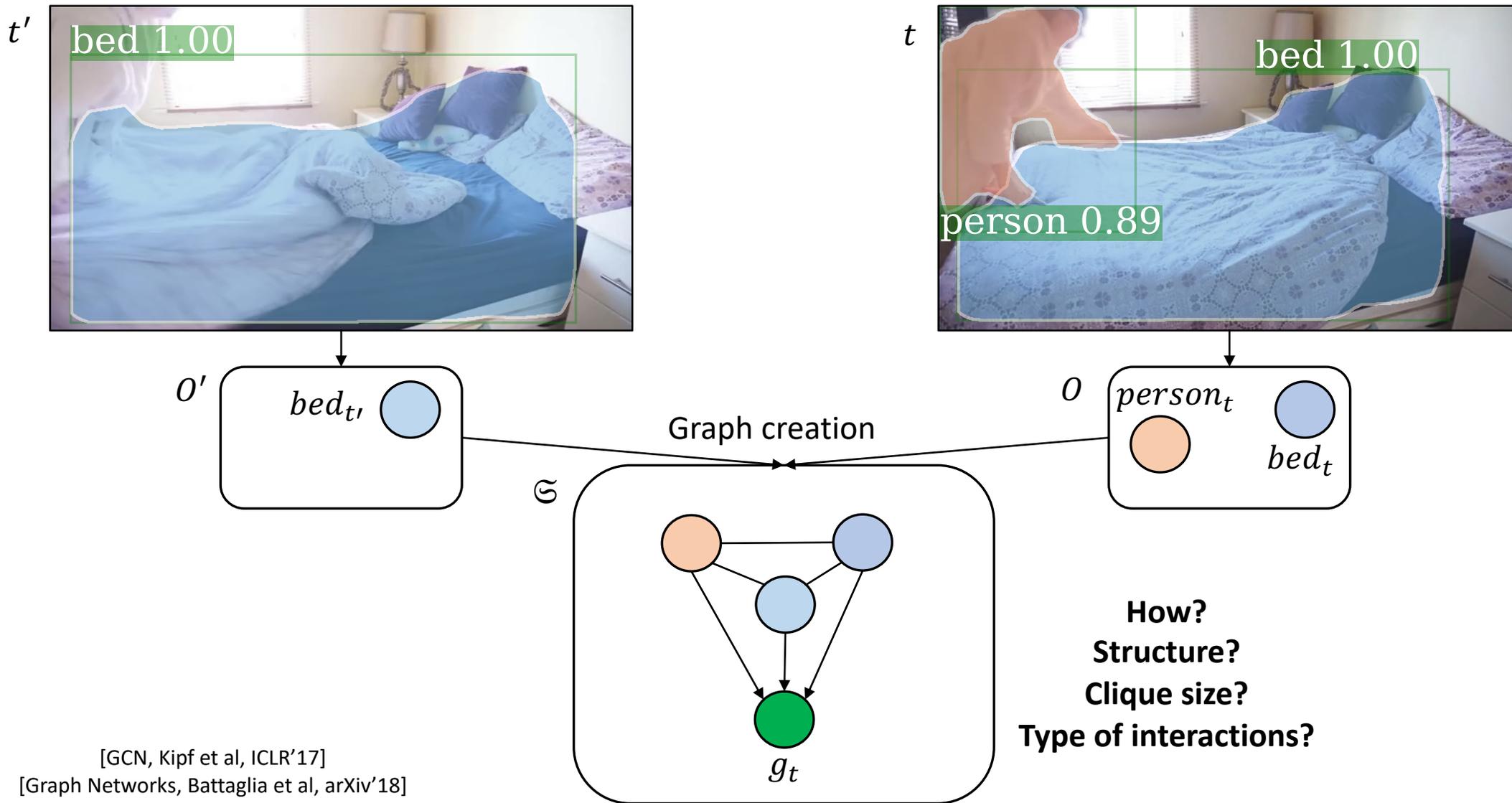
# Image as a set of objects



RGB  
Mask-RCNN



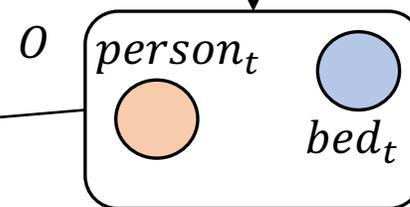
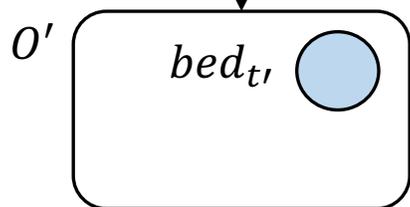
# Object Relation Network



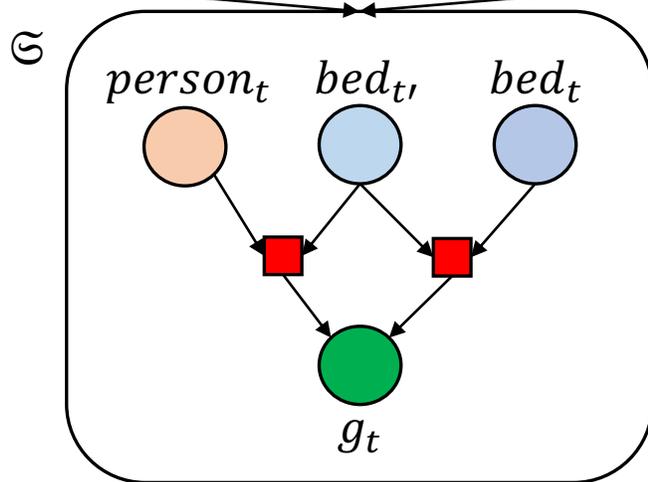
[GCN, Kipf et al, ICLR'17]

[Graph Networks, Battaglia et al, arXiv'18]

# Object Relation Network



Graph creation



Shared MLP

$$g_t = f(\text{bed}_{t'}, \text{person}_t) + f(\text{bed}_{t'}, \text{bed}_t)$$

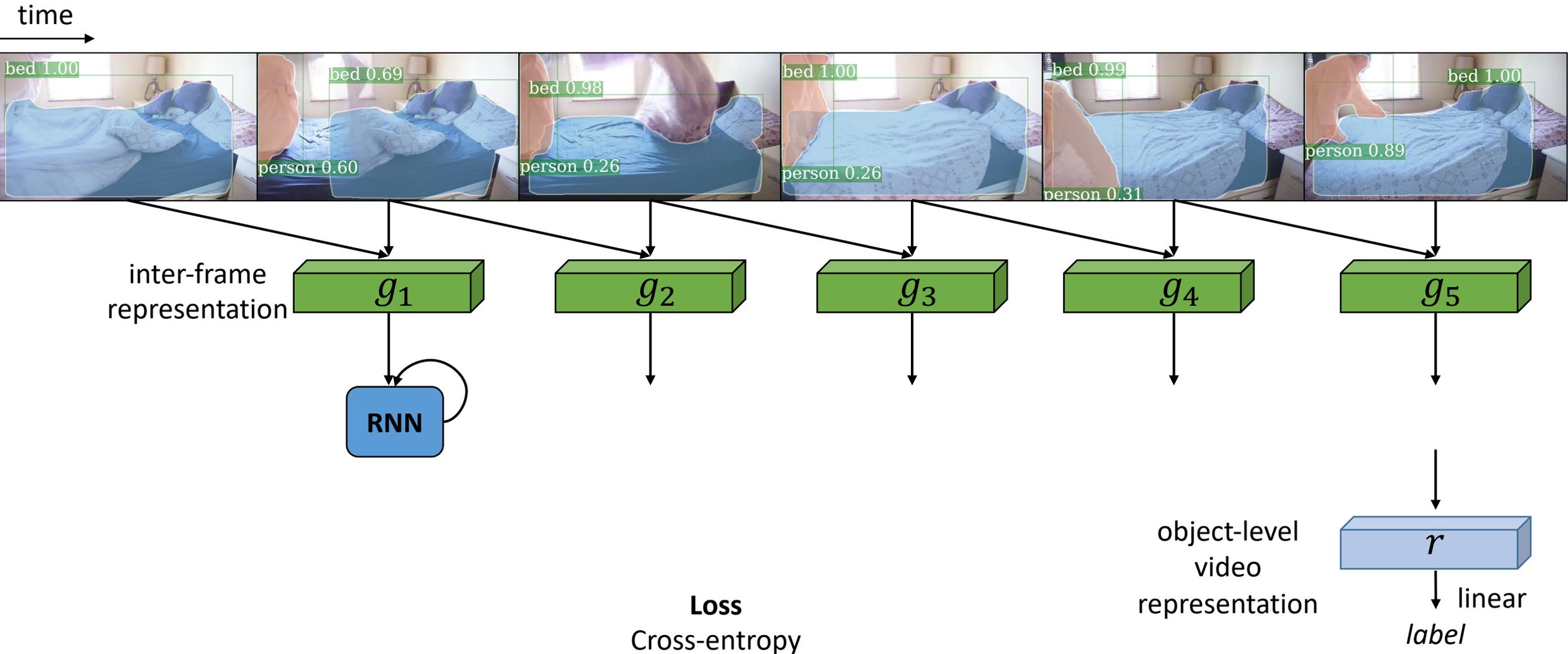
$$g_t = \sum_{o \in \mathcal{O}} \sum_{o' \in \mathcal{O}'} f(o', o)$$

- Data efficient
- Invariant to the number of objects
- Inter-frame object relations
- Semantic meaning
- Clique of size 2

[GCN, Kipf et al, ICLR'17]

[Graph Networks, Battaglia et al, arXiv'18]

# Object Relation Network



# Object Relation Network

## *State-of-the-art*

<i>Method</i>	<i>Acc.</i>
C3D	21.50
I3D	27.63
Multiscale TRN	33.60
<b>Object Relation Network</b>	<b>35.97</b>

*Accuracy on Something-Something*

<i>Method</i>	<i>mAP</i>
Resnet50	40.5
I3D	39.7
<b>Object Relation Network</b>	<b>44.7</b>

*Mean Average Precision on VLOG*

<i>Method</i>	<i>Acc.</i>
Resnet18	32.05
Resnet3D-18	34.20
<b>Object Relation Network</b>	<b>40.89</b>

*Verb accuracy on EPIC Kitchens*

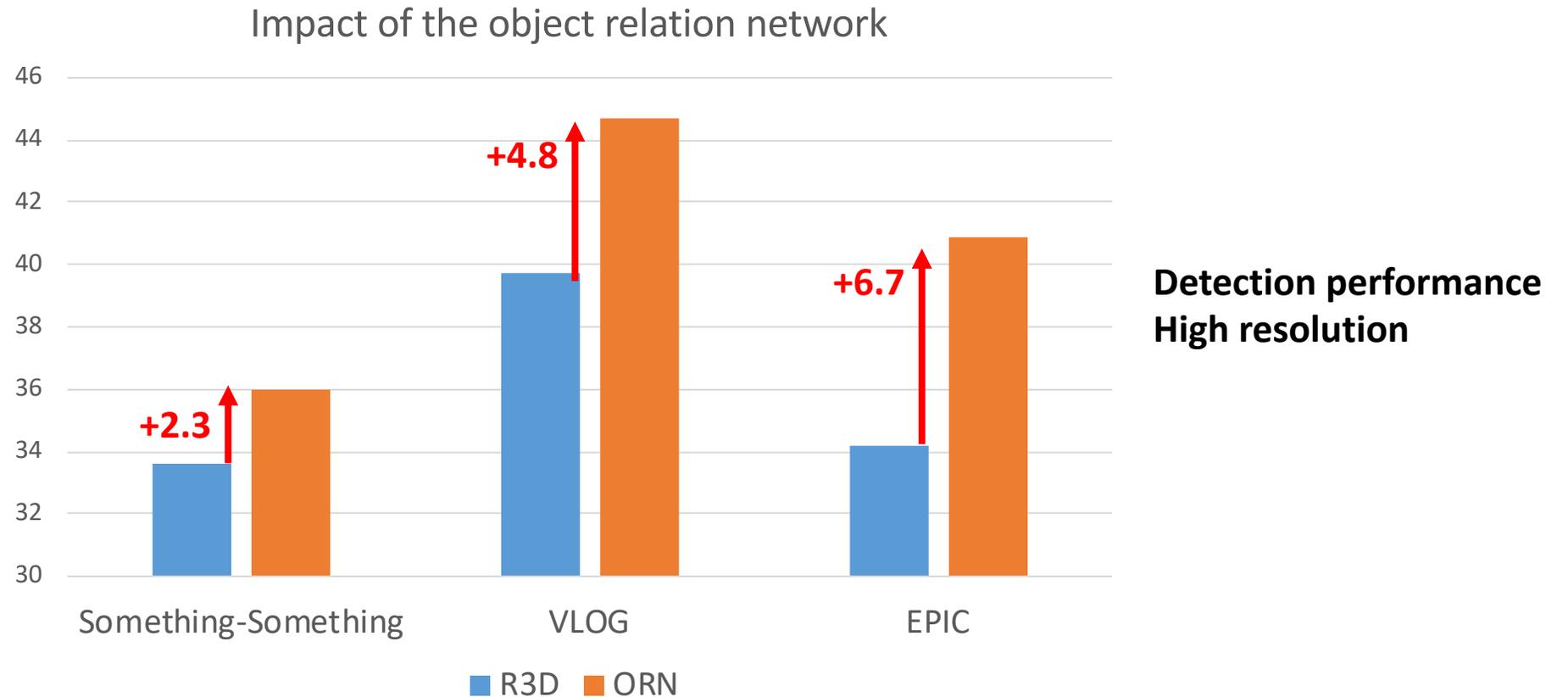


[fabienbaradel/object\\_level\\_visual\\_reasoning](https://github.com/fabienbaradel/object_level_visual_reasoning)



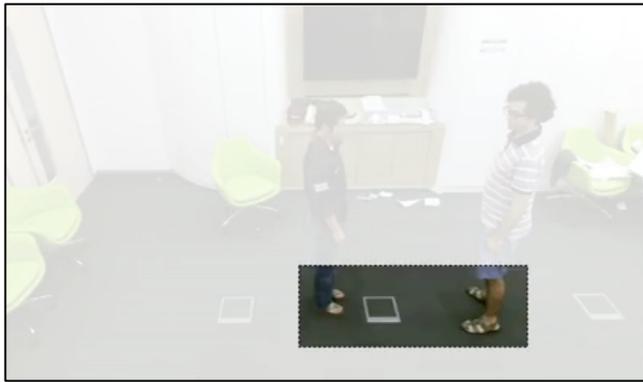
*Object masks detected by Mask-RCNN*

# Object Relation Network *Ablation study*





# Outline

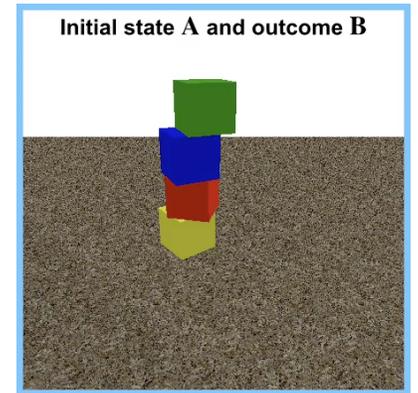


**Visual Attention**



**Entity-level interactions**

Structure matters...  
Can we go one step further?  
Beyond supervised learning?  
Learning underlying latent concepts?



**Reasoning**

« Counterfactual learning »

*F. Baradel, N. Neverova, J.*

*Mille, C. Wolf*  
ICLI (light)



**Christian Wolf**  
INSA Lyon - LIRIS



**Julien Mille**  
INSA CVL - LI Tours

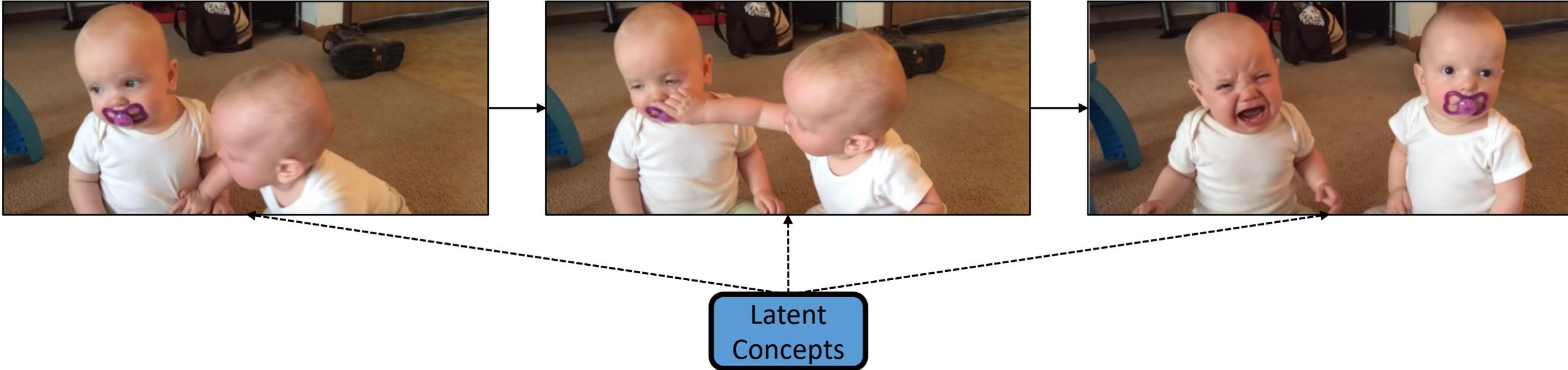


**Natalia Neverova**  
Facebook



**Greg Mori**  
SFU

# Reasoning & Causation

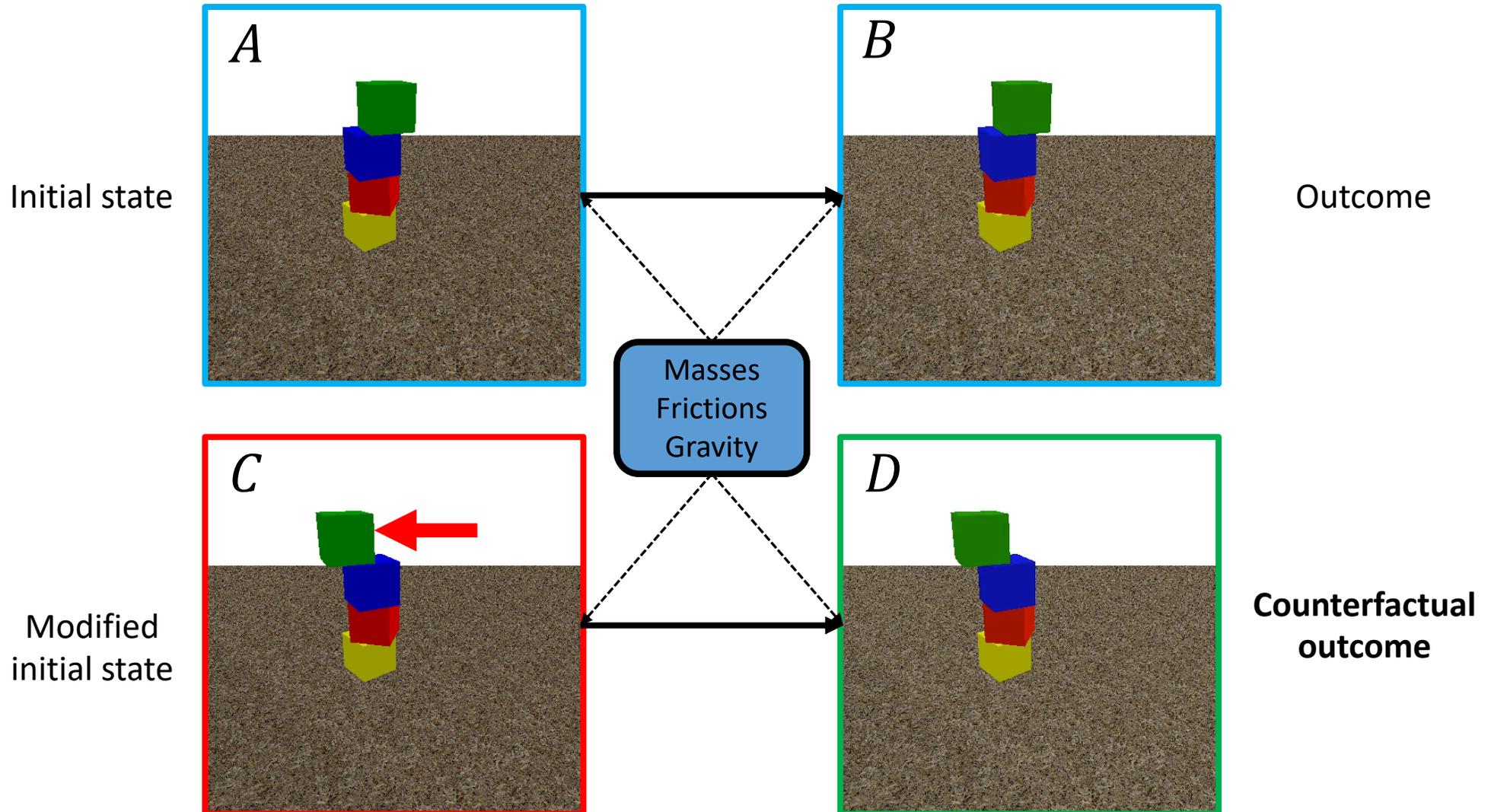


Understanding of complex relationships  
Cause-effect

What would happened if?  
Counterfactual statement

# Counterfactual

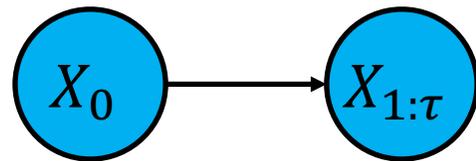
*Future forecasting*



# Counterfactual

*Future forecasting*

Feedforward



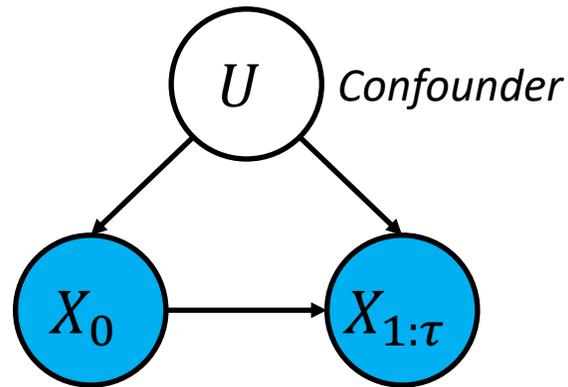
*A*

*B*

Initial state

Outcome

Counterfactual

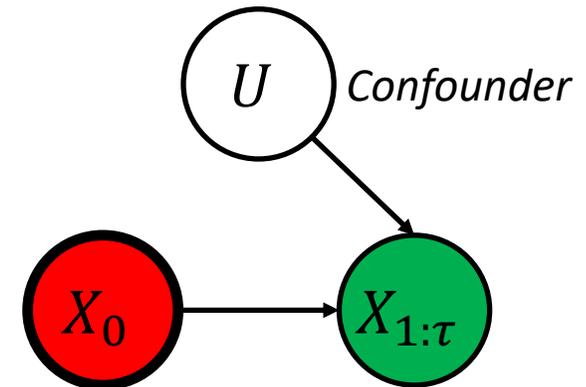


*A*

*B*

Initial state

Outcome



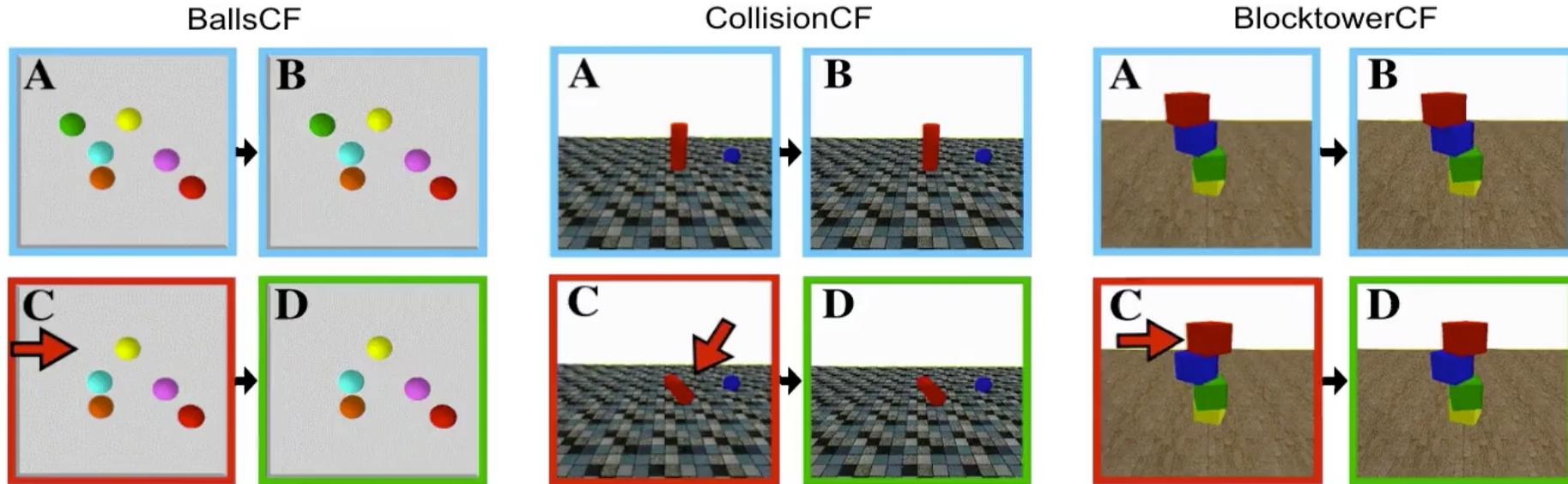
*do*( $X_0 = C$ )

*D*

Modified initial state

Counterfactual outcome

# CoPhy benchmark



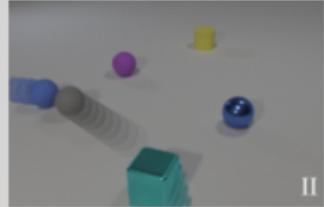
Large-scale datasets  
250k examples ((A,B), (C,D))  
7 millions of frames  
Supervision of the do-operator (  $\Rightarrow$  )  
Confounders are necessary for future prediction

**IntPhys**



- physical plausibility
- out of distributions events

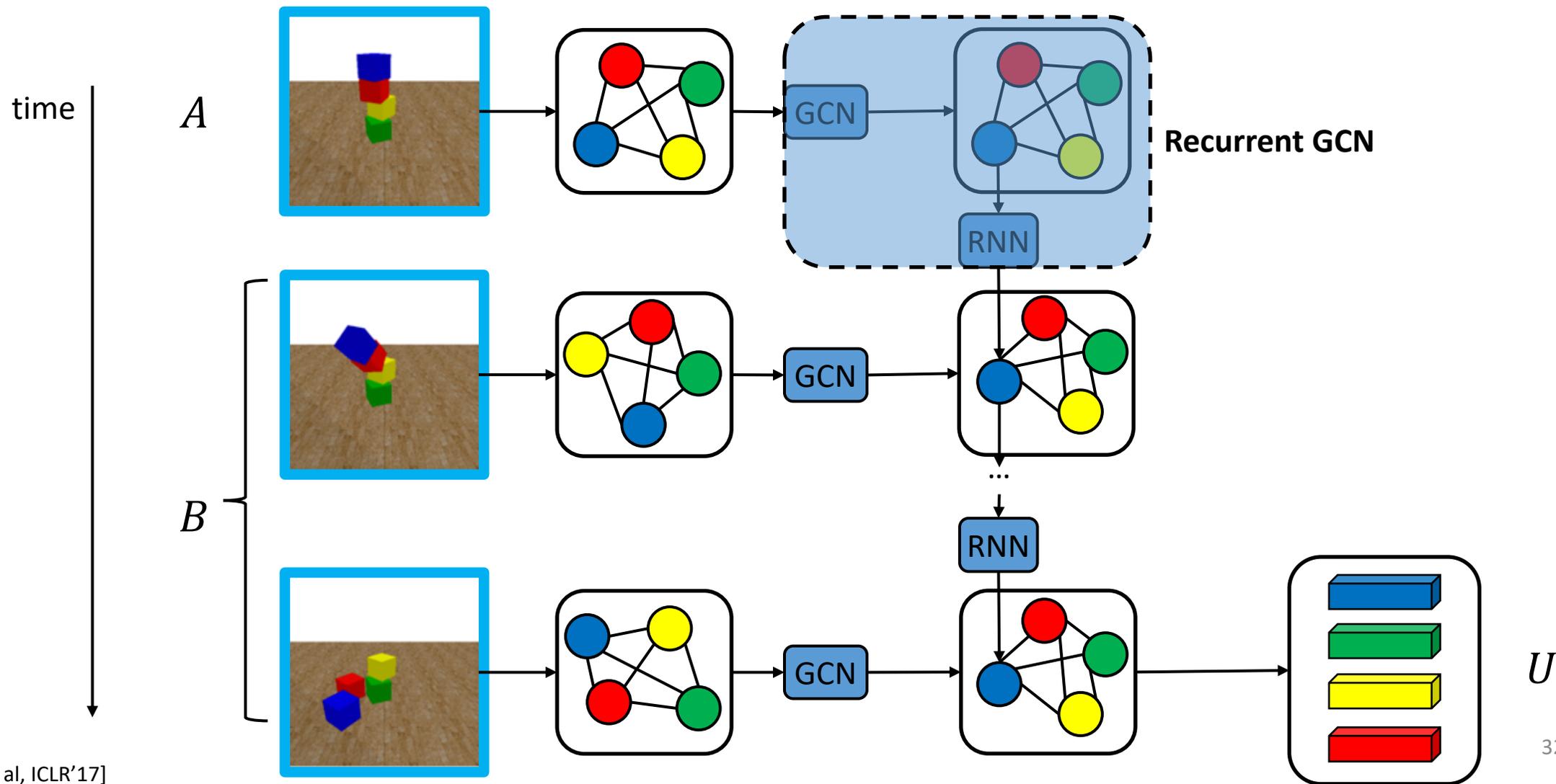
**CLEVRER**



- choose between counterfactual answers
- wide range of tasks

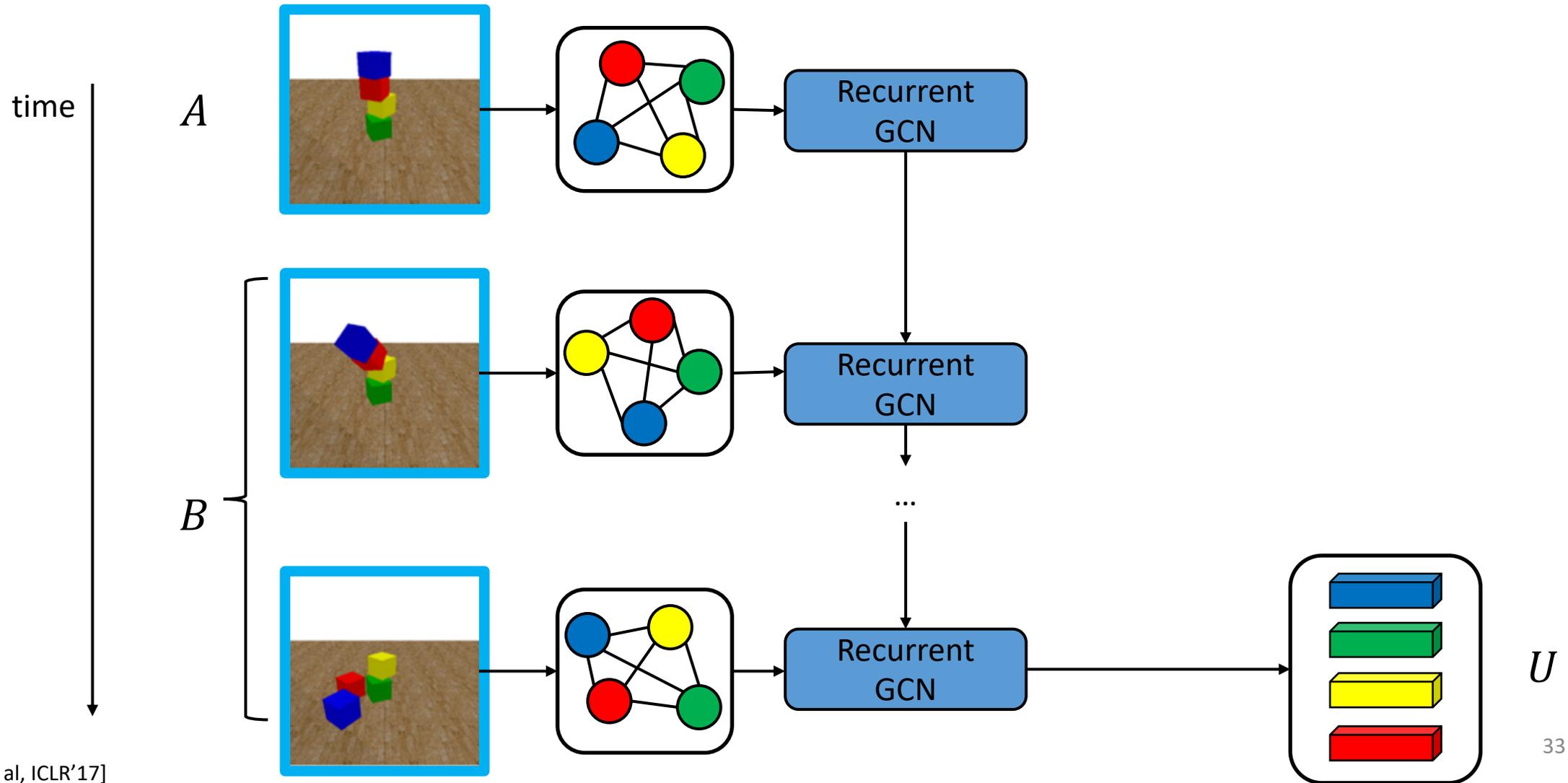
# CoPhyNet

*Unsupervised confounders estimations*



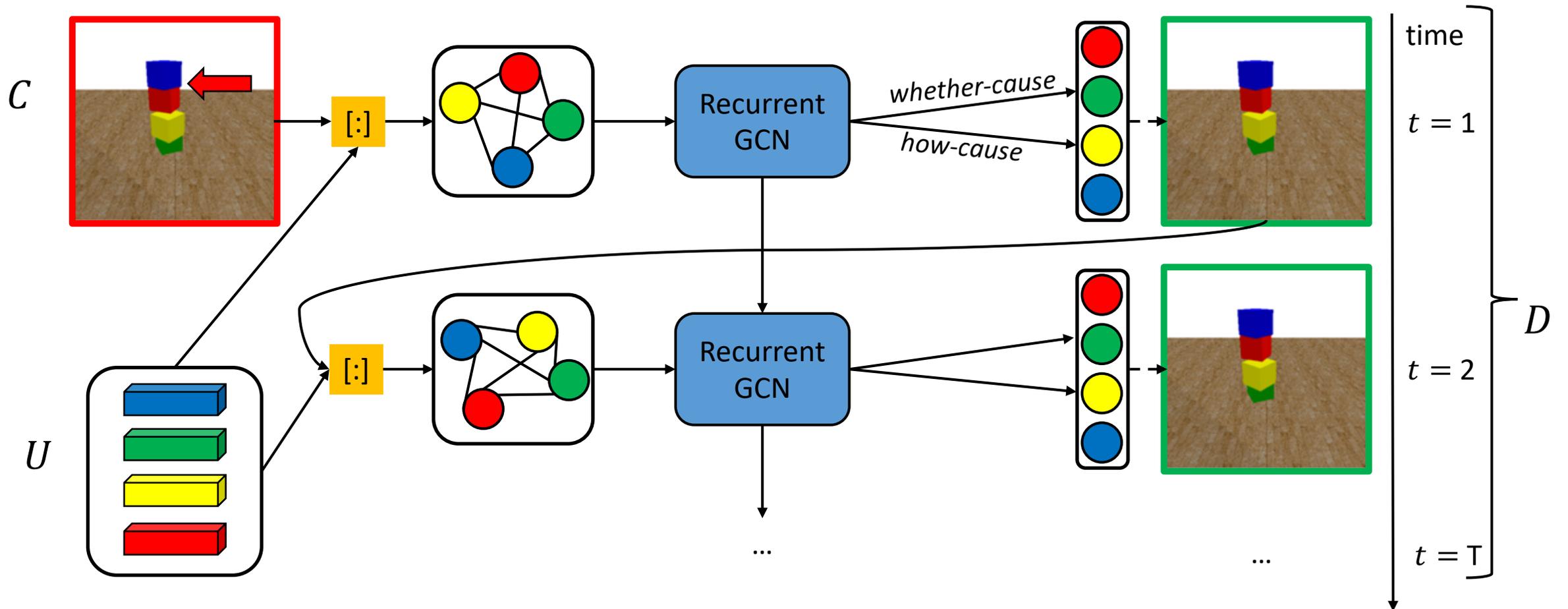
# CoPhyNet

*Unsupervised confounders estimations*

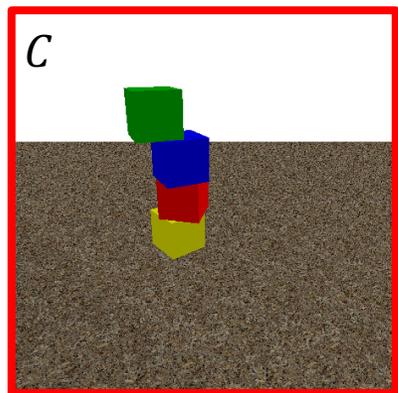


# CoPhyNet

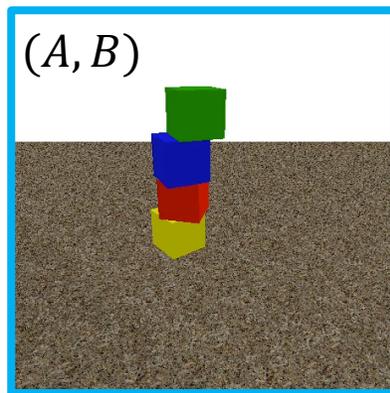
## Trajectory prediction



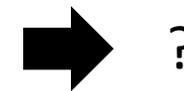
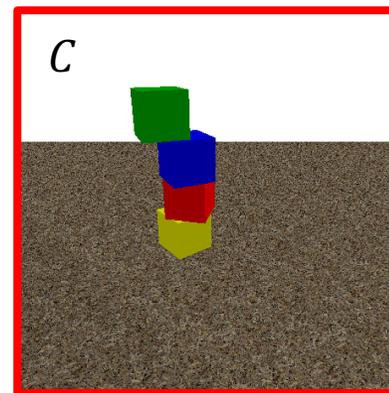
# Human study



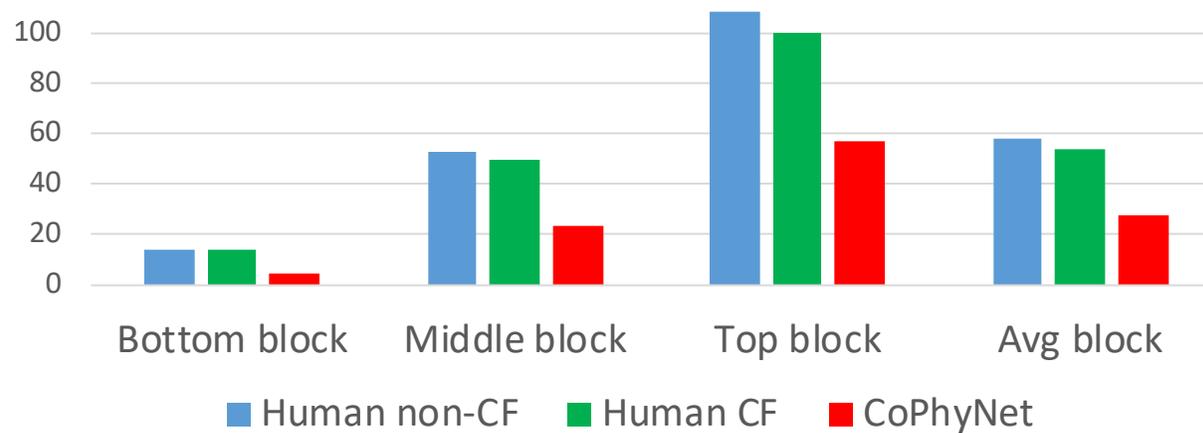
Human non-CF



Human CF



2D pixel error for each block



# Cophynet Results

NOT COMPARABLE!

Copying baselines

Feedforward models

Soft-upper bound

Train → Test	Copy C	Copy B	IN	NPE	CoPhyNet	IN Sup.
3 → 3	0.470	0.601	0.318	0.331	<b>0.294</b>	0.296
3 → 3*	0.365	0.592	0.298	0.319	<b>0.289</b>	0.282
3 → 4	0.754	0.846	0.524	0.523	<b>0.482</b>	0.467

Unseen confounders

3 → 3\*

Unseen number of blocks

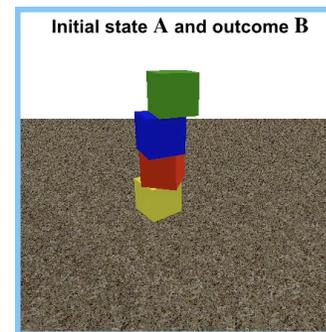
3 → 4

MSE on 3D positions (average over time)



[fabienbaradel/cophy](https://github.com/fabienbaradel/cophy)

+

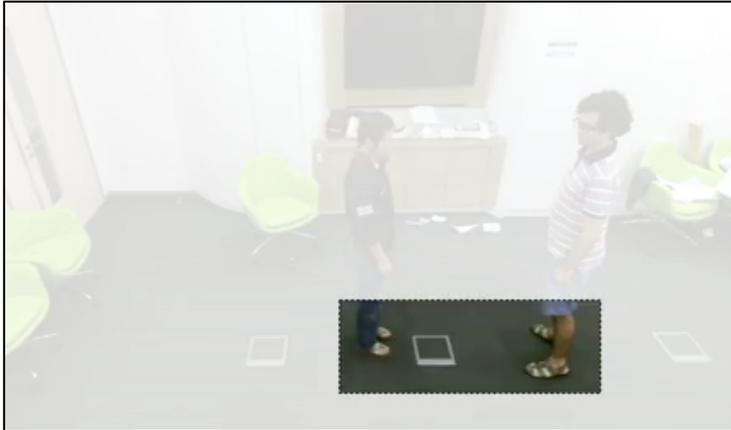


CoPhy benchmark

[Interaction Network, Battaglia et al, NIPS'17]

[Neural Physic Engine, Chang et al, ICLR'17]

# Conclusion



## Visual Attention

« Glimpse Clouds »  
*F. Baradel, C. Wolf, J. Mille,  
G. Taylor, CVPR'18*

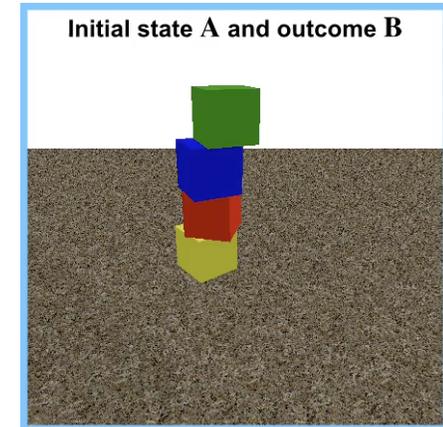
Focus on important parts  
Automatic selection  
Distributed recognition



## Entity-level interactions

« Object level Reasoning »  
*F. Baradel, N. Neverova, C. Wolf,  
J. Mille, G. Mori, ECCV'18*

Object-centric modeling  
Intra-time interactions  
Learned relations



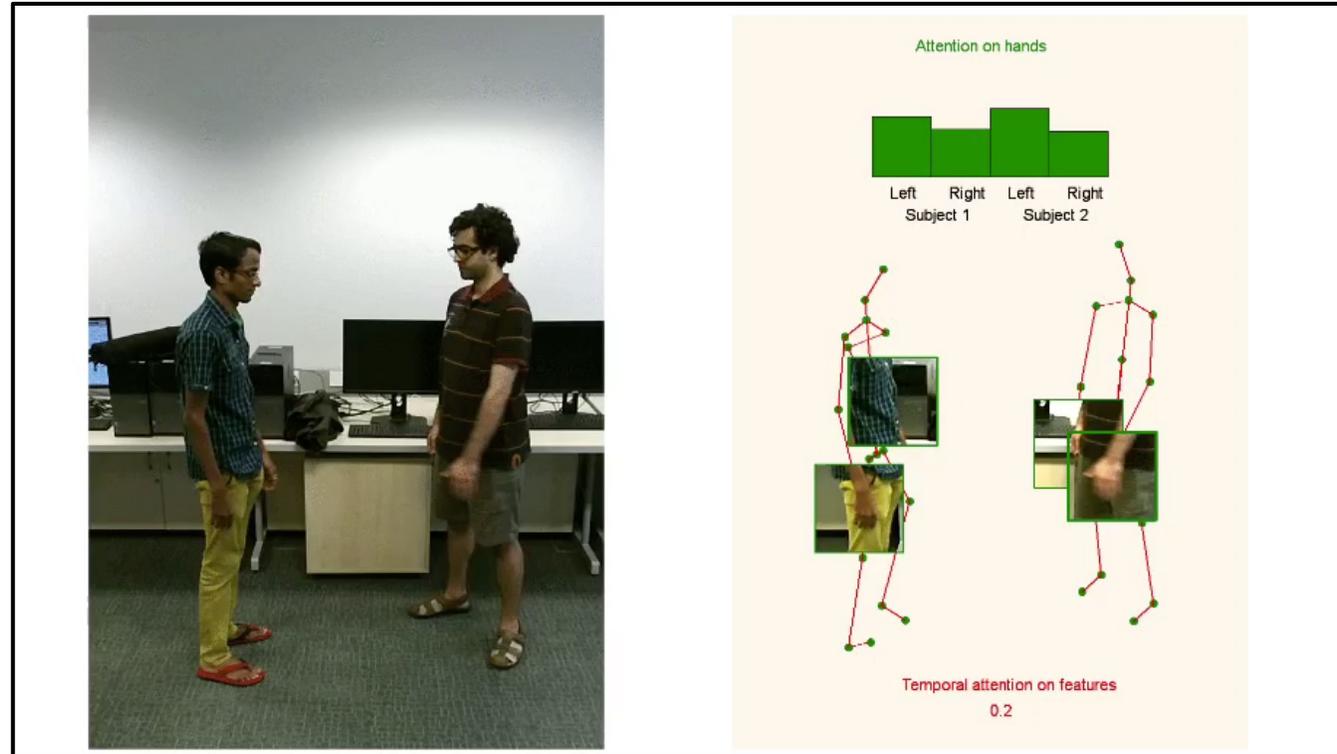
## Reasoning

« Counterfactual learning »  
*F. Baradel, N. Neverova, J. Mille,  
G. Mori, C. Wolf, ICLR'20 (spotlight)*

Unsupervised latent discovery  
Future trajectory  
New task in visual space <sup>37</sup>

# Other works

## *Pose-driven Attention*



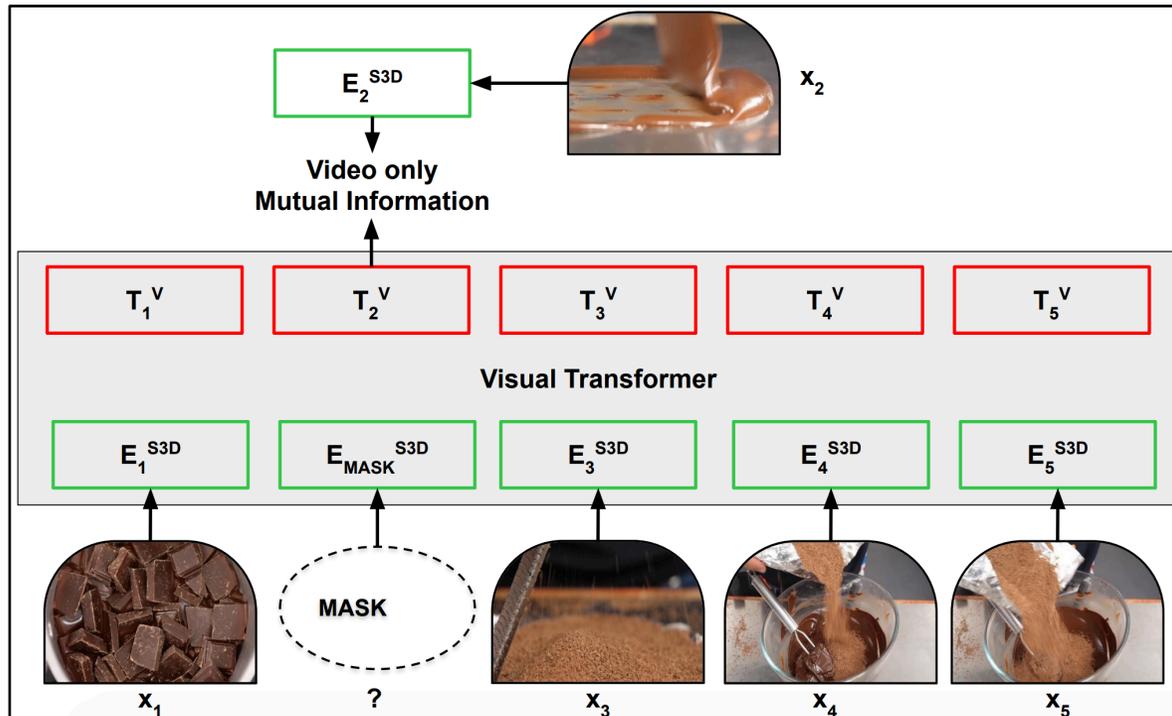
« Focus to Hands »  
F. Baradel, C. Wolf, J. Mille,  
ICCVW'17

Skeleton & RGB  
Handcrafted/Learned context features  
Focus around hands

« Pose-driven Attention to RGB »  
F. Baradel, C. Wolf, J. Mille,  
BMVC'18

# Other works

## *Self-supervised learning*



« Contrastive Bidirectional Transformer »

*C. Sun, F. Baradel, K. Murphy, C. Schmid*

Under review ECCV'20

Human learning  
Beyond large-scale annotated datasets  
Efficient  
Discover regularities  
Prediction of missing parts  
Instructional videos  
Vision-Text alignment



**Cordelia Schmid**  
INRIA – Google



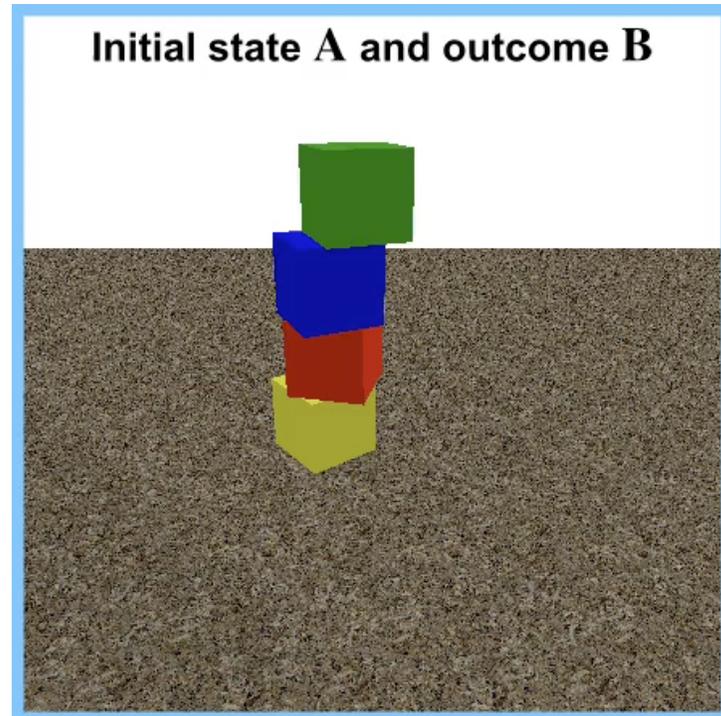
**Chen Sun**  
Google



**Kevin P. Murphy**  
Google

# What next?

## *Real word counterfactual predictions*

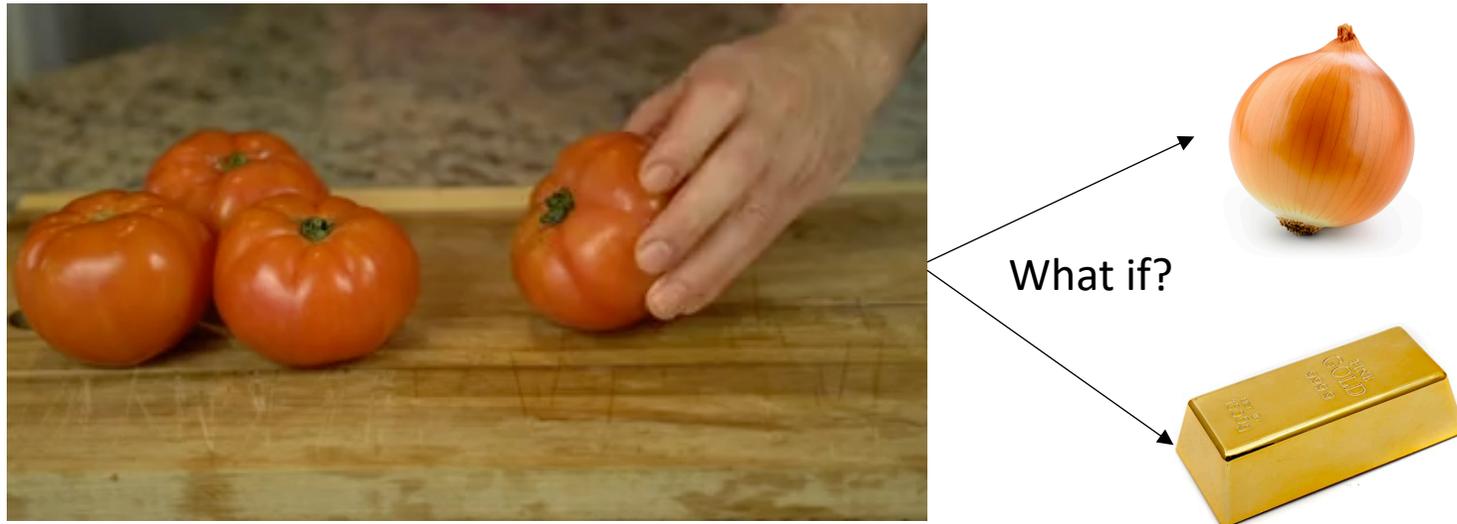


« CoPhy++ »  
*F. Baradel, N. Neverova, J. Mille,*  
*G. Mori, C. Wolf*  
**To be submitted to TPAMI**

No object supervision  
Unsupervised keypoints  
Predictions in image space

# What next?

## *Real word counterfactual predictions*



Beyond correlation and dataset biases  
Latent concept  
Generalization  
Semantical structure  
Ontology

# What next?

## *Disentangled representation*

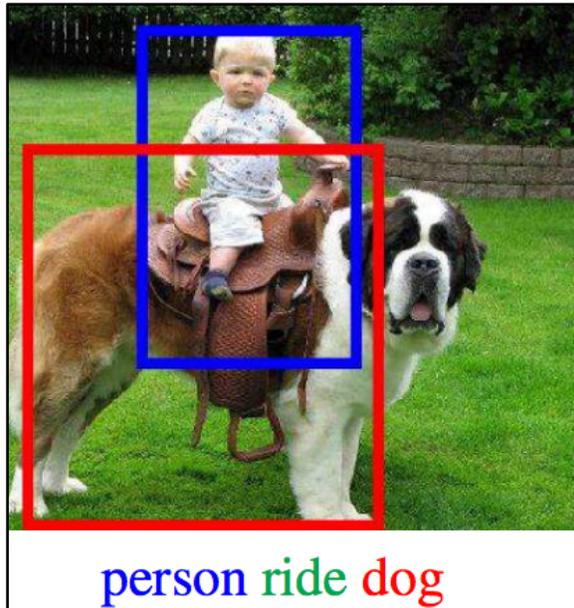


Image vs Video  
Appearance vs Motion  
Human Object Interaction vs Long-range activities  
Efficient representation

# Thank you!



**Christian Wolf**  
INSA Lyon - LIRIS



**Julien Mille**  
INSA CVL - LI Tours



**Natalia Neverova**  
Facebook AI Research



**Graham W. Taylor**  
University of Guelph  
Vector Institute



**Greg Mori**  
Simon Fraser University  
Borealis AI



**Cordelia Schmid**  
INRIA – Google



**Chen Sun**  
Google



**Kevin P. Murphy**  
Google