

DEEP LEARNING FOR ARTIFICIAL INTELLIGENCE

3rd Master Course UPC ETSETB TelecomBCN Barcelona. Autumn 2019.



Instructors



Xavier
Giró-i-Nieto



Marta R.
Costa-jussà



Noé
Casas



Verónica
Vilaplana



Ramon
Morros



Javier
Ruiz



Albert
Pumarola



Jordi
Torres

Organizers



UNIVERSITAT POLITÈCNICA
DE CATALUNYA
BARCELONATECH



Supporters

Google Cloud
GitHub Education

+ info: <http://bit.ly/dlai2019>

[\[course site\]](#)



#DLBCN

Day 10 Lecture 1

Self-supervised Learning



Xavier Giro-i-Nieto

xavier.giro@upc.edu

@DocXavi

Associate Professor

Universitat Politècnica de Catalunya
Technical University of Catalonia



Video lecture



UAB UOC upf.

Master in Computer Vision Barcelona

[\[http://pàgines.uab.cat/mcv/\]](http://pàgines.uab.cat/mcv/)



Xavier Giro-i-Nieto
xavier.giro@upc.edu

Associate Professor
Universitat Politècnica de
Catalunya



Module 6 - Day 8 - Lecture 1 Self-supervised Learning from Video Sequences 28th March 2019



UAB UOC upf.

Master in Computer Vision Barcelona

[\[http://pàgines.uab.cat/mcv/\]](http://pàgines.uab.cat/mcv/)



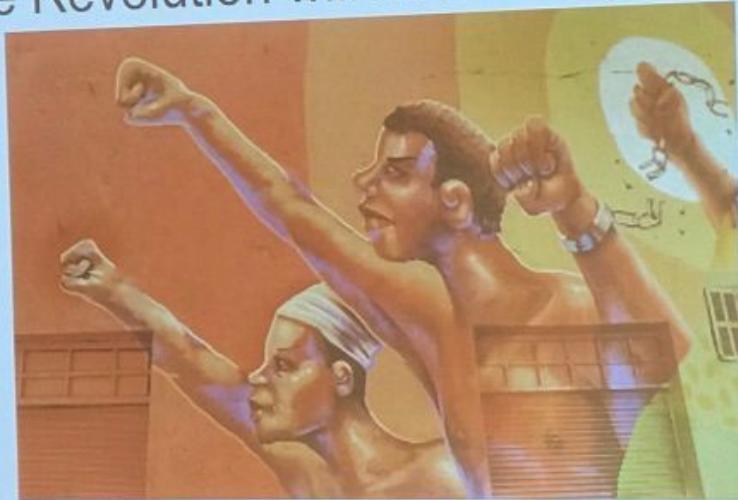
Xavier Giro-i-Nieto
xavier.giro@upc.edu
Associate Professor
Universitat Politècnica de
Catalunya



Module 6 - Day 9 - Lecture 2 Self-supervised Audiovisual Learning 4th April 2019

Lectures on applied deep learning for video in [Master in Computer Vision Barcelona 2019](#)

The Revolution will not be Supervised



Alexei A. Efros
UC Berkeley



Outline

- 1. Unsupervised Learning**
- 2. Self-supervised Learning**
 - a. Autoencoder
 - b. Temporal regularisations
 - c. Temporal verifications
 - d. Predictive Learning
 - e. Miscellaneous: optical flow, color & multiview

Types of machine learning

Yann Lecun's Black Forest cake

■ "Pure" Reinforcement Learning (cherry)

- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**



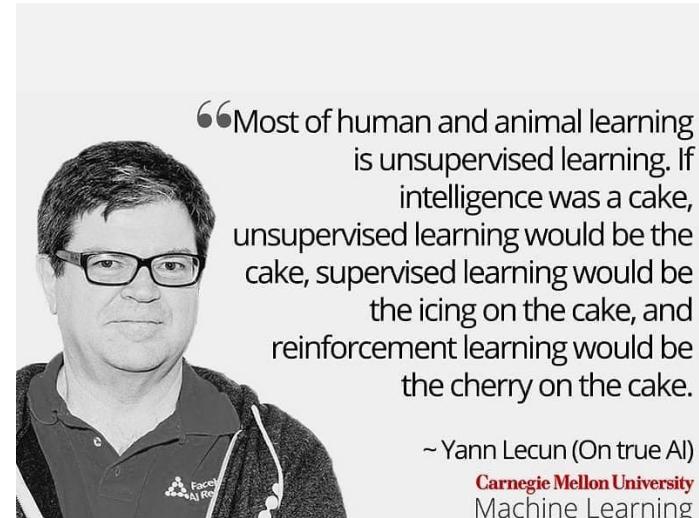
■ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ Predicts future frames in videos
- ▶ **Millions of bits per sample**

■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)



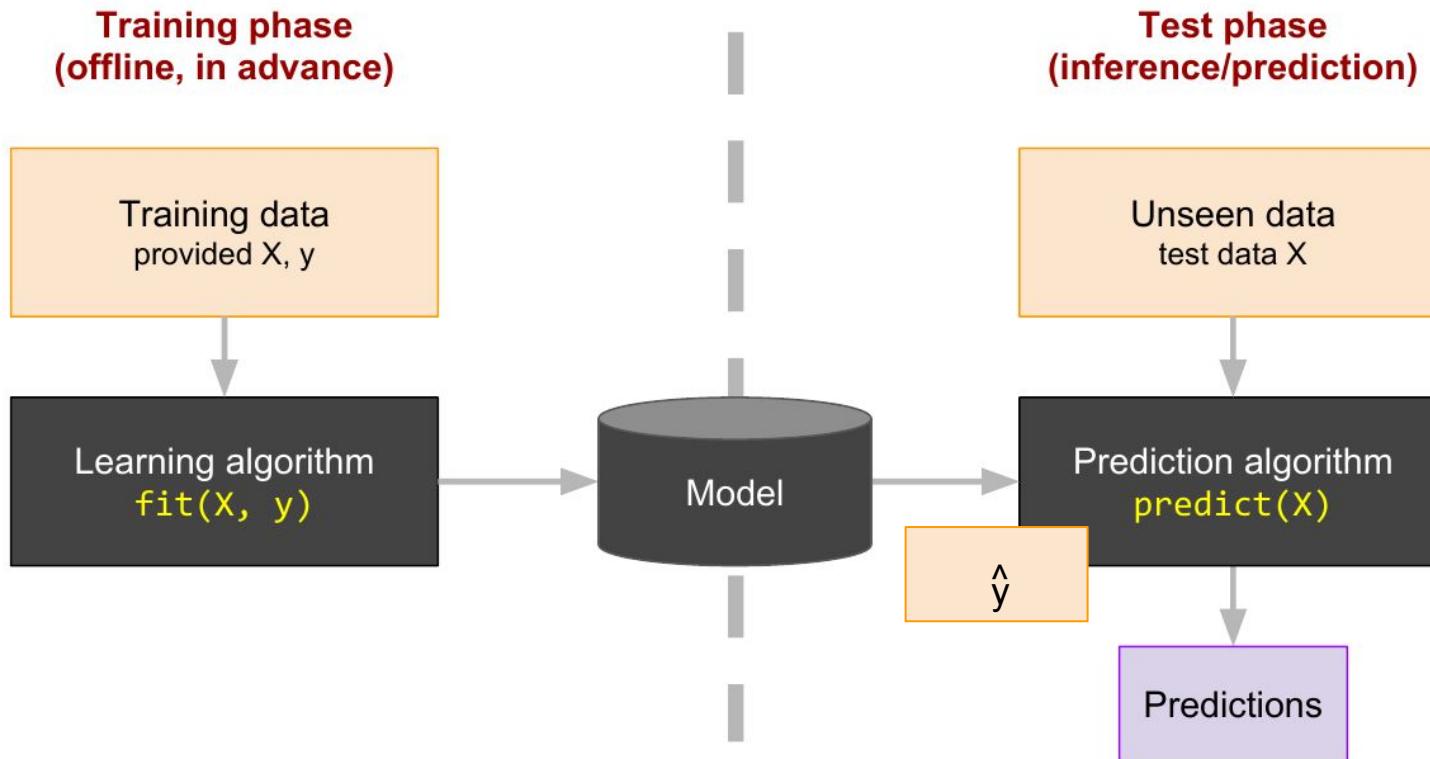
A broader picture of types of learning...

	...with a teacher	...without a teacher
Active agent...	Reinforcement learning (with extrinsic reward)	Intrinsic motivation / Exploration.
Passive agent...	Supervised learning	Unsupervised learning



Slide inspired by Alex Graves (Deepmind) at
["Unsupervised Learning Tutorial"](#) @ NeurIPS 2018.

Supervised learning



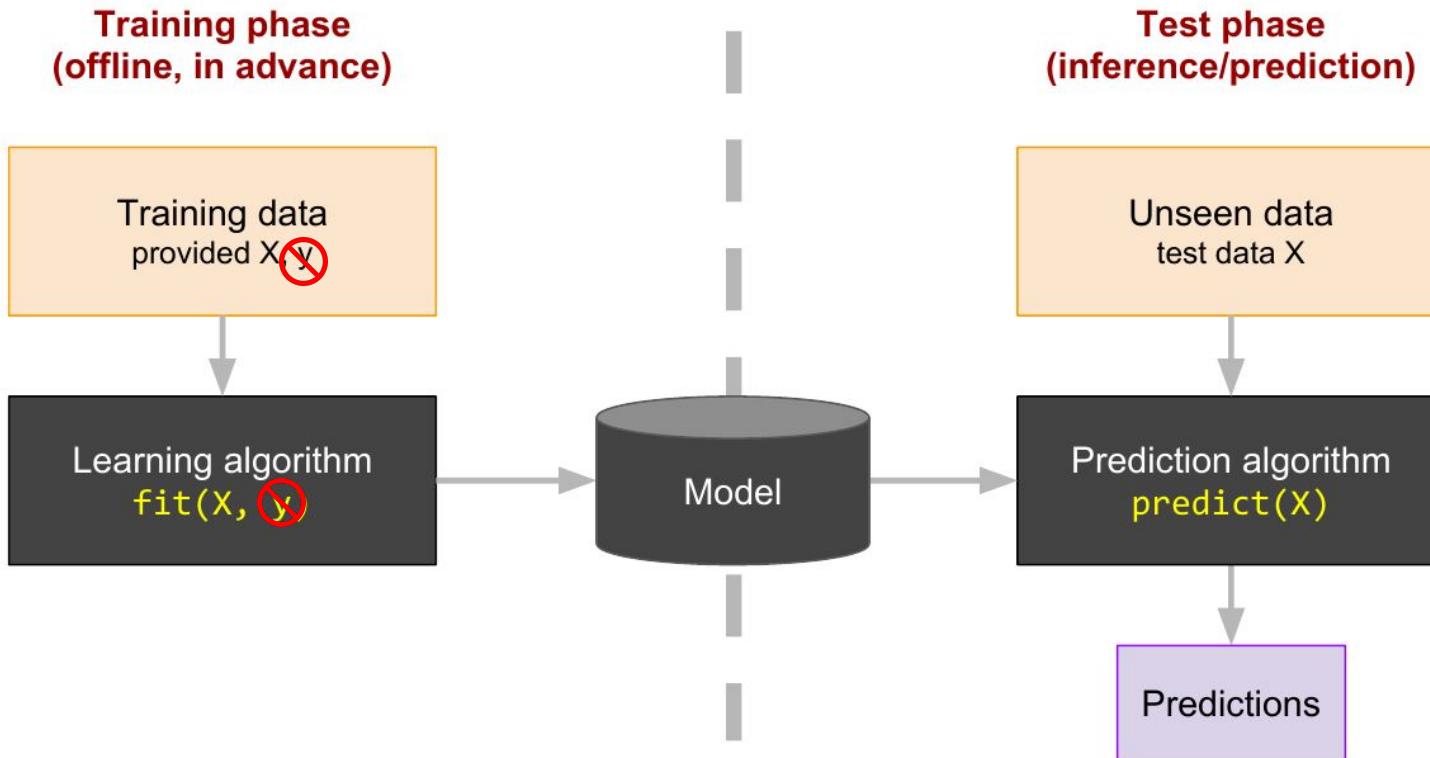
A broader picture of types of learning...

	...with a teacher	...without a teacher
Active agent...	Reinforcement learning (with extrinsic reward)	Intrinsic motivation / Exploration.
Passive agent...	Supervised learning	Unsupervised learning



Slide inspired by Alex Graves (Deepmind) at
["Unsupervised Learning Tutorial"](#) @ NeurIPS 2018.

Unsupervised learning



Unsupervised learning

- It is the nature of how intelligent beings percept the world.
- It can save us tons of efforts to build an AI agent compared to a totally supervised fashion.
- Vast amounts of unlabelled data.

WHY?



Video lectures on Unsupervised Learning



The manifold hypothesis

The data distribution lie close to a low-dimensional manifold

Example: consider image data

- Very high dimensional (1,000,000)
- A randomly generated image will almost certainly not look like any real-world scene
 - The space of images that occur in nature is almost completely empty
- Hypothesis: real world images lie on a smooth, low-dimensional manifold
 - Mahalanobis distance is a good measure of similarity

Similar for audio and text



UPC
UNIVERSITAT POLITÈCNICA DE CATALUNYA
BarcelonaTech
Departament de Teoria del Senyal i Comunicacions



DEEP LEARNING FOR ARTIFICIAL INTELLIGENCE

Master Class UPC ETSAT-TELMOIN-CN Barcelona, Autumn 2017

Instructors



Organizers



Supporters



+ info: <http://daii.deeplearning.barcelona>

[course site]



Xavier Giro-i-Nieto
savent.giro@upc.edu

Associate Professor
Universitat Politècnica de Catalunya
Technical University of Catalonia



Kevin McGuinness, [UPC DLCV 2016](#)

Xavier Giró, [UPC DLAI 2017](#)

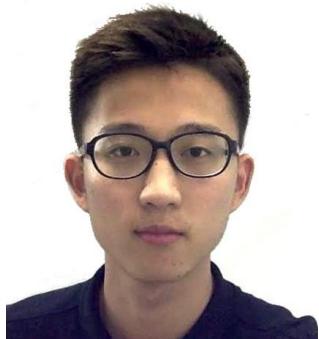
Outline

1. Unsupervised Learning
2. **Self-supervised Learning**
 - a. Autoencoder
 - b. Temporal regularisations
 - c. Temporal verifications
 - d. Predictive Learning
 - e. Miscellaneous: optical flow, color & multiview

Acknowledgements



Víctor Campos



Junting Pan



Xunyu Lin



Carlos
Arenas



Sebastian
Palacio



Barcelona
Supercomputing
Center
Centro Nacional de Supercomputación

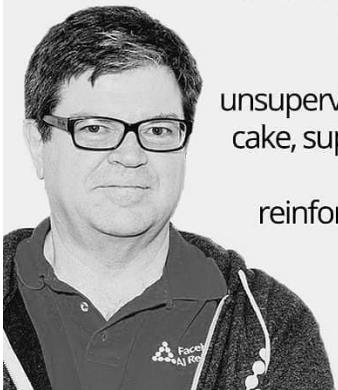


UNIVERSITAT POLITÈCNICA
DE CATALUNYA



German
Research Center
for Artificial
Intelligence

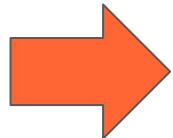
Self-supervised learning



“Most of human and animal learning is unsupervised learning. If intelligence was a cake, unsupervised learning would be the cake, supervised learning would be the icing on the cake, and reinforcement learning would be the cherry on the cake.”

~ Yann Lecun (On true AI)

Carnegie Mellon University
Machine Learning



Yann LeCun
@ylecun

I Now call it "self-supervised learning", because "unsupervised" is both a loaded and confusing term.

In self-supervised learning, the system learns to predict part of its input from other parts of its input. In...
facebook.com/722677142/post...

4:40 PM · Apr 30, 2019 · Facebook

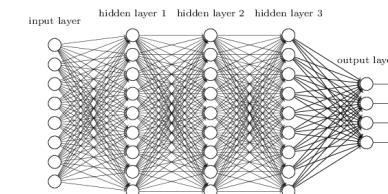
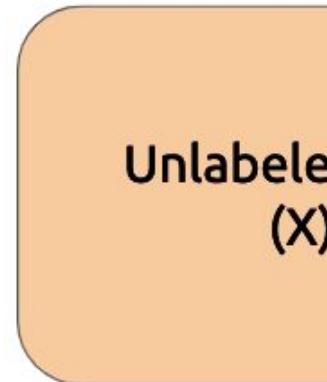
(...)

“Doing this properly and reliably is **the greatest challenge in ML and AI** of the next few years in my opinion.”

Self-supervised learning

Self-supervised learning is a form of unsupervised learning where the data provides the supervision.

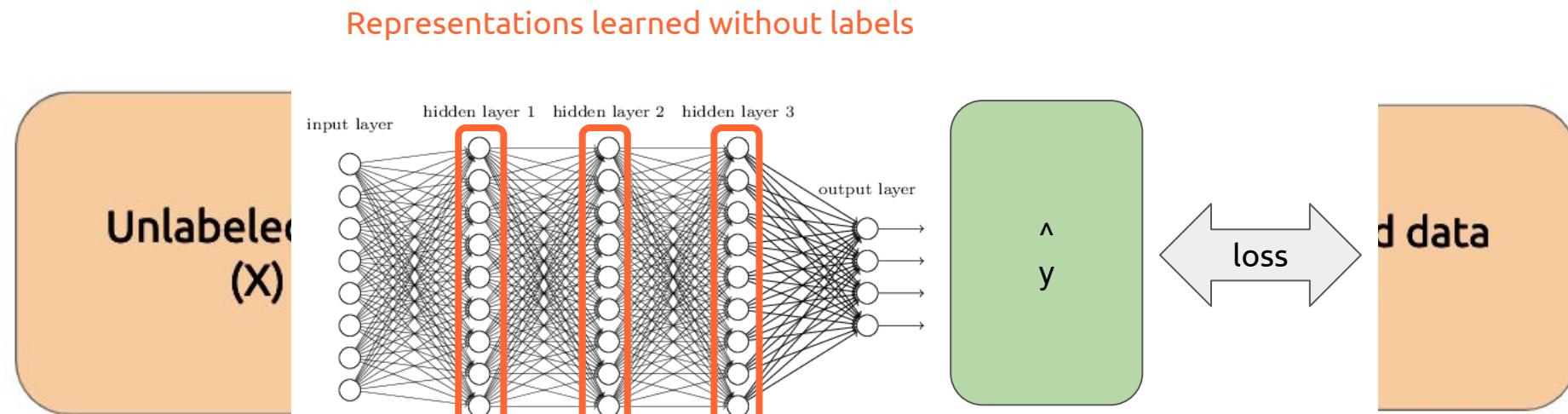
- A **surrogate task** must be invented by withholding a part of the unlabeled data and training the NN to predict it.



Self-supervised learning

Self-supervised learning is a form of unsupervised learning where the data provides the supervision.

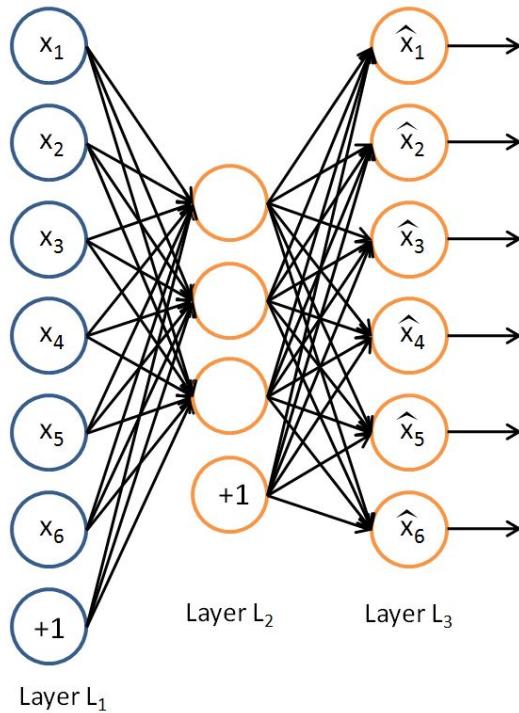
- By defining a proxy loss, the NN learns **representations**, which should be valuable for the actually target task.



Outline

1. Unsupervised Learning
2. Self-supervised Learning
 - a. **Autoencoder**
 - b. Temporal regularisations
 - c. Verifications
 - d. Predictive Learning
 - e. Miscellaneous: optical flow, color & multiview

Autoencoder (AE)



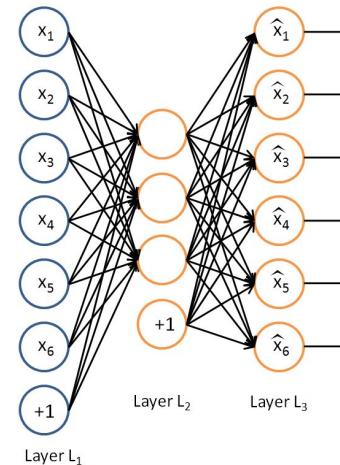
Autoencoders:

- Predict at the output the same input data.
- Do not need labels.

Autoencoder (AE)

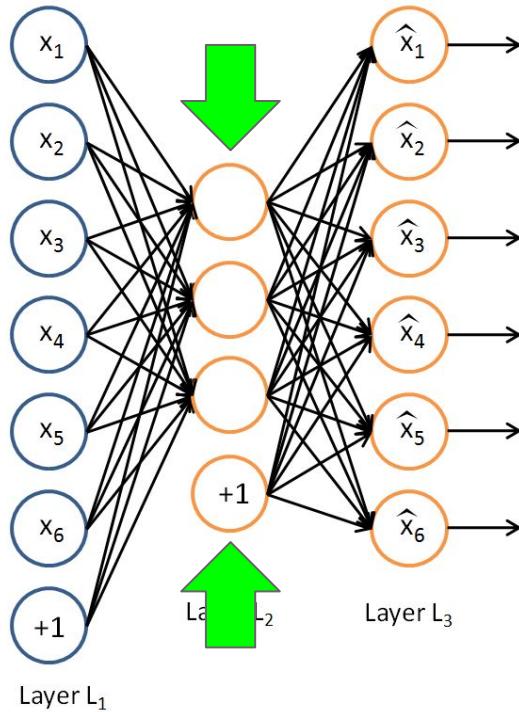
What is the use of an autoencoder ?

WHY?



Autoencoder (AE)

WHY?



Dimensionality reduction:

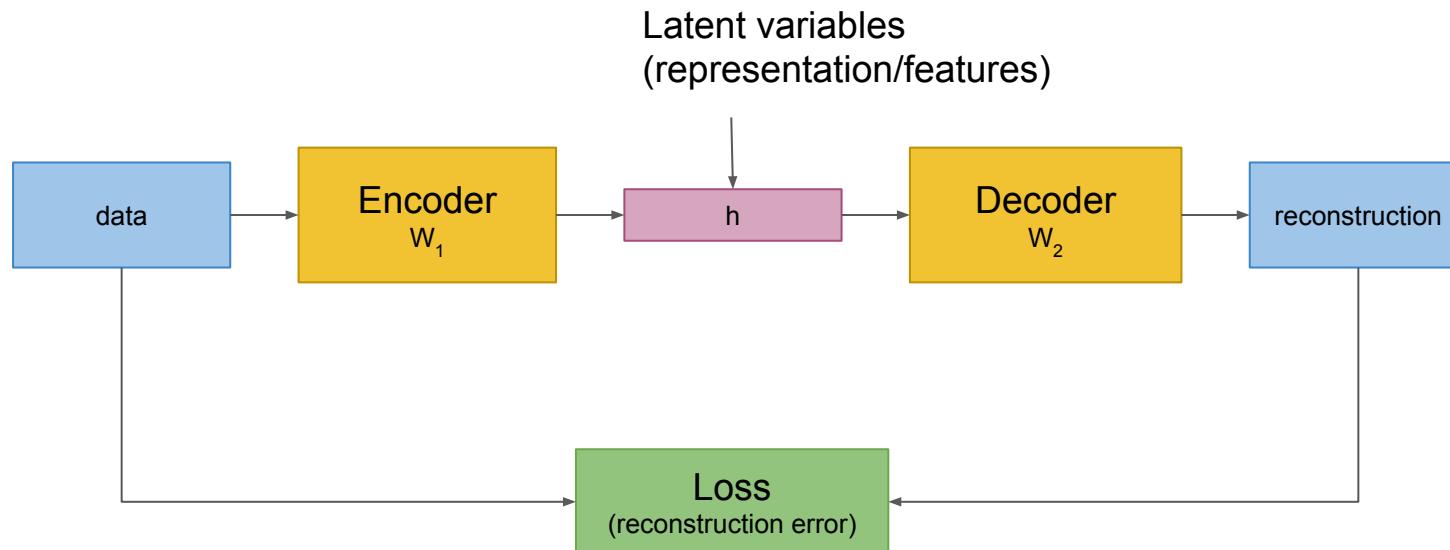
Use the hidden layer as a feature extractor of any desired size.

Autoencoder (AE)

WHY?



1. Initialize a NN by solving an autoencoding problem.

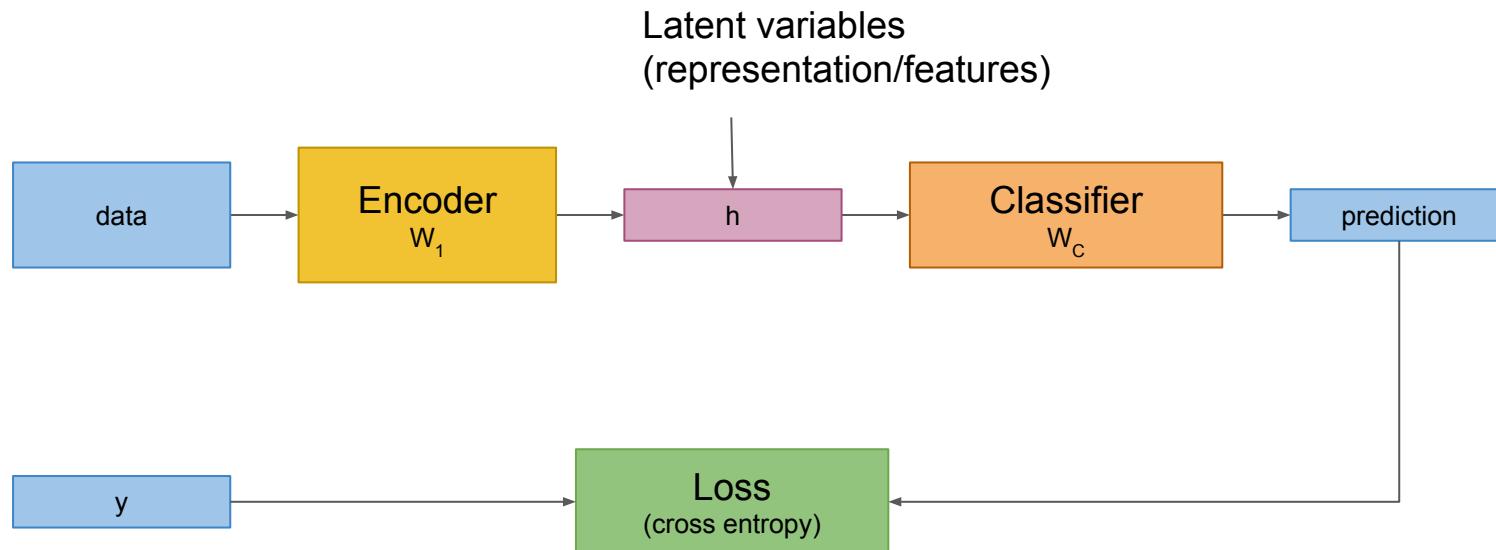


Autoencoder (AE)

WHY?

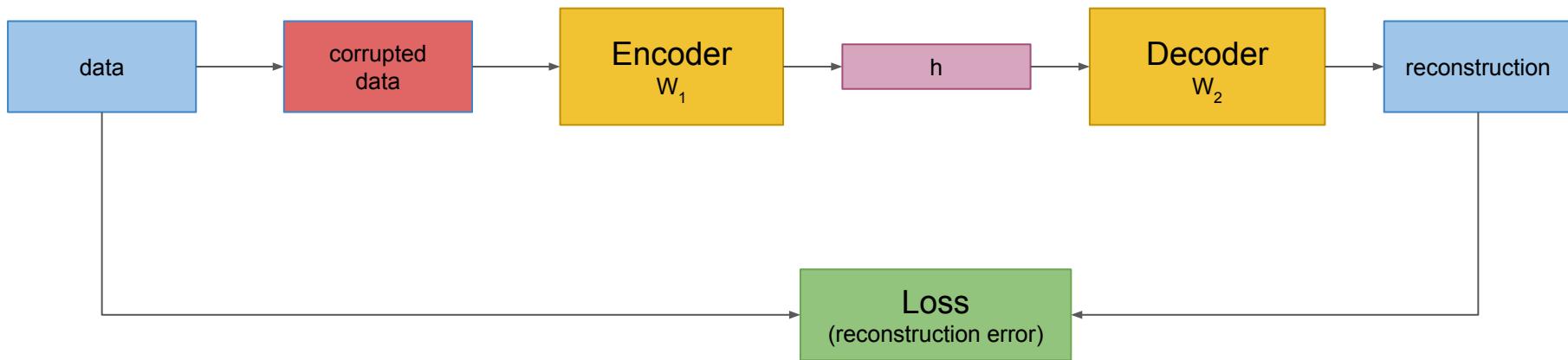


1. Initialize a NN solving an autoencoding problem.
2. Train for final task with “few” labels.



Denoising Autoencoder (DAE)

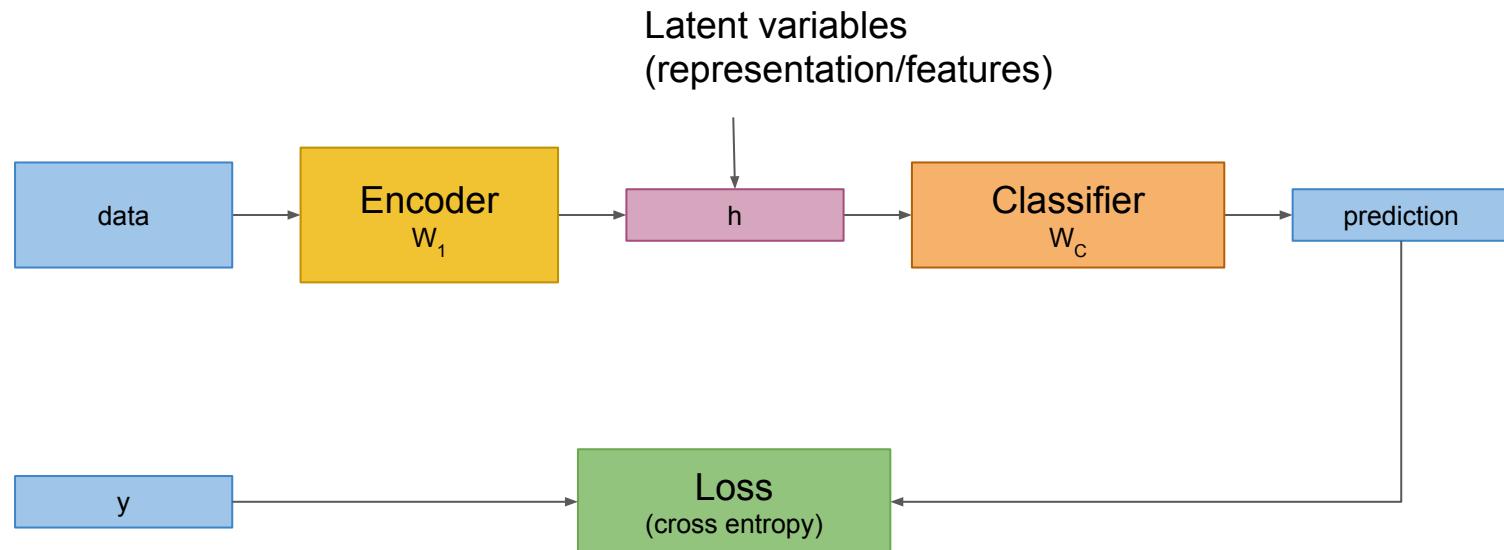
1. Corrupt the input signal, but predict the clean version at its output.
2. Train for final task with “few” labels.



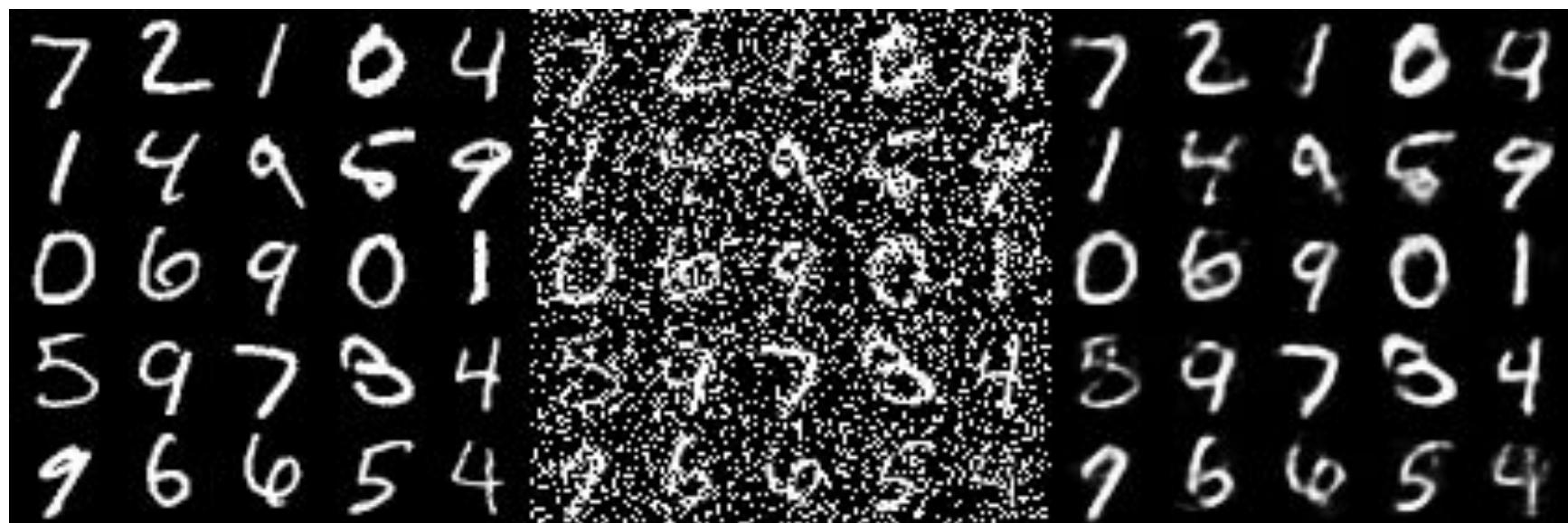
Vincent, Pascal, Hugo Larochelle, Yoshua Bengio, and Pierre-Antoine Manzagol. ["Extracting and composing robust features with denoising autoencoders."](#) ICML 2008.

Denoising Autoencoder (DAE)

1. Corrupt the input signal, but predict the clean version at its output.
2. Train for final task with “few” labels.



Denoising Autoencoder (DAE)





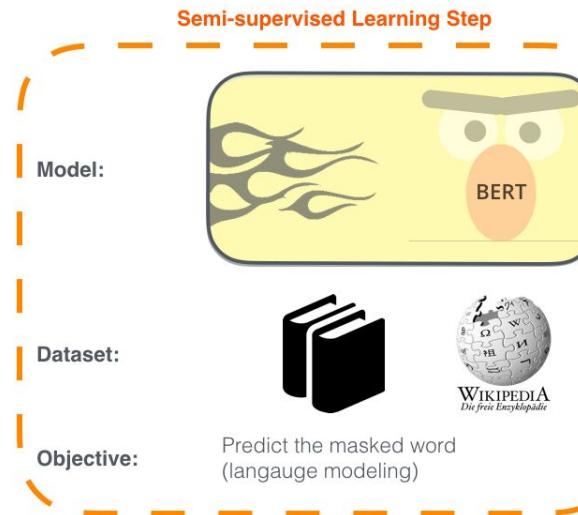
Eigen, David, Dilip Krishnan, and Rob Fergus. ["Restoring an image taken through a window covered with dirt or rain."](#) ICCV 2013.

Masked Autoencoder (MAE)

Show it a sequence of words on input, mask out 15% of the words, and ask the system to predict the missing words (or a distribution of words)

1 - **Semi-supervised** training on large amounts of text (books, wikipedia..etc).

The model is trained on a certain task that enables it to grasp patterns in language. By the end of the training process, BERT has language-processing abilities capable of empowering many models we later need to build and train in a supervised way.



2 - **Supervised** training on a specific task with a labeled dataset.

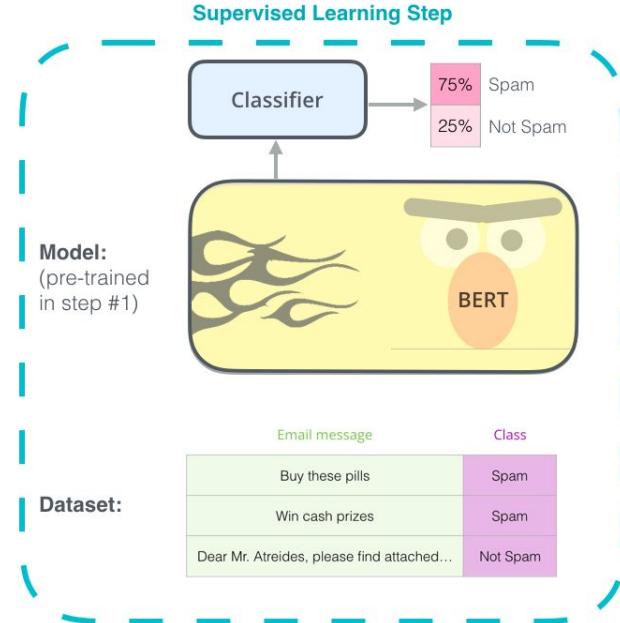


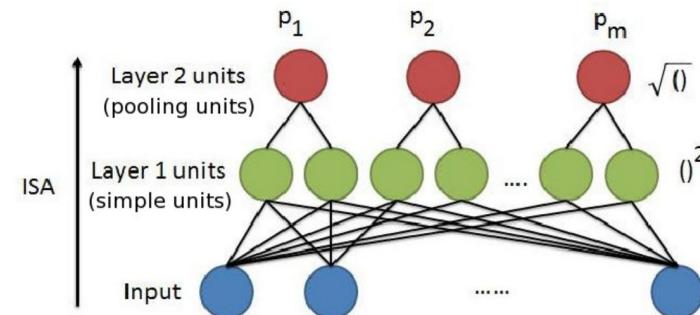
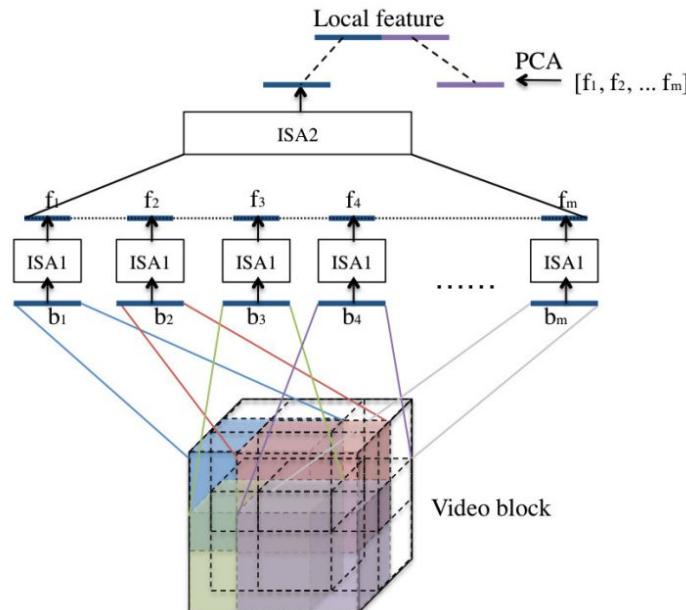
Figure: [Jay Alammar](#)

Outline

1. Unsupervised Learning
2. Self-supervised Learning
 - a. Autoencoder
 - b. **Temporal regularisations**
 - c. Verifications
 - d. Predictive Learning
 - e. Miscellaneous: optical flow, color & multiview

Temporal regularization: ISA

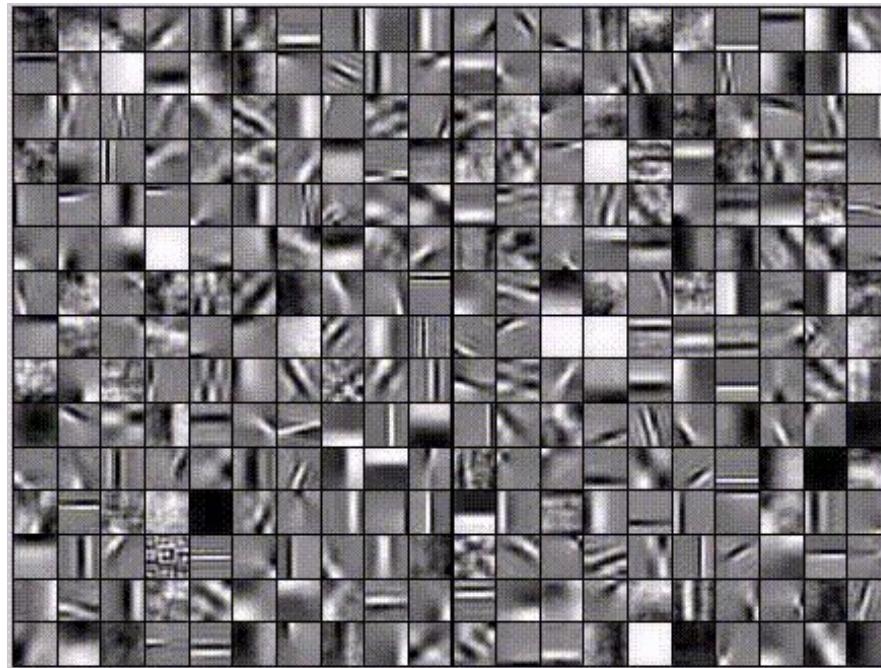
Video features are learned with convolutions and poolings combined with a Independent Subspace Analysis (ISA).



	KTH	Hollywood2	UCF	YouTube
Best published results	92.1%	50.9%	85.6%	71.2%
Our results	93.9%	53.3%	86.5%	75.8%

Temporal regularization: ISA

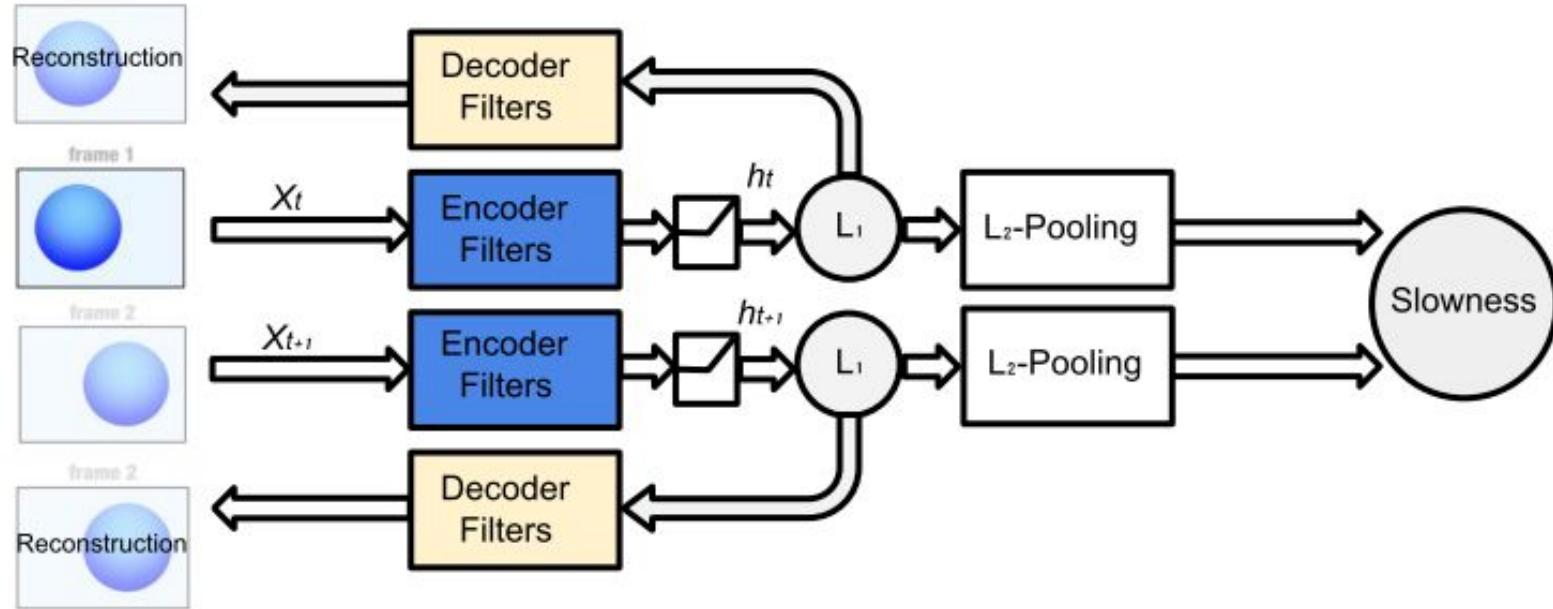
Feature visualizations.



Le, Quoc V., Will Y. Zou, Serena Y. Yeung, and Andrew Y. Ng. "[Learning hierarchical invariant spatio-temporal features for action recognition with independent subspace analysis.](#)" CVPR 2011

Temporal regularization: Slowness

Assumption: adjacent video frames contain semantically similar information.
Autoencoder trained with regularizations by slowness and sparsity.



Goroshin, Ross, Joan Bruna, Jonathan Tompson, David Eigen, and Yann LeCun. "[Unsupervised learning of spatiotemporally coherent metrics.](#)" ICCV 2015.

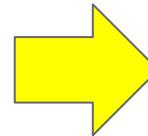
Temporal regularization: Slowliness

Slow feature analysis

- Temporal coherence assumption: features should change slowly over time in video

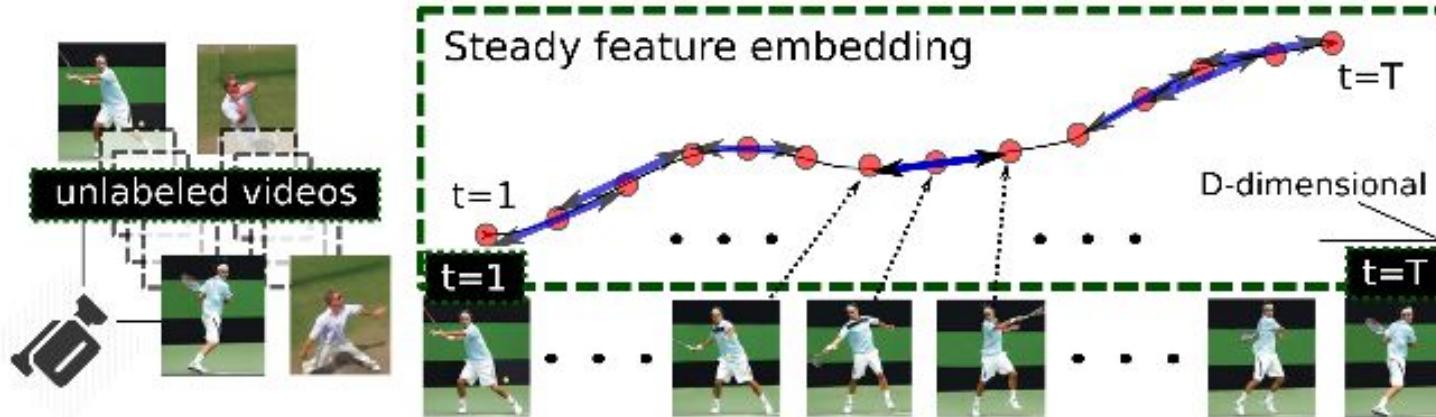
Steady feature analysis

- Second order changes also small: changes in the past should resemble changes in the future



Train on triplets of frames from video

Loss encourages nearby frames to have slow and steady features, and far frames to have different features

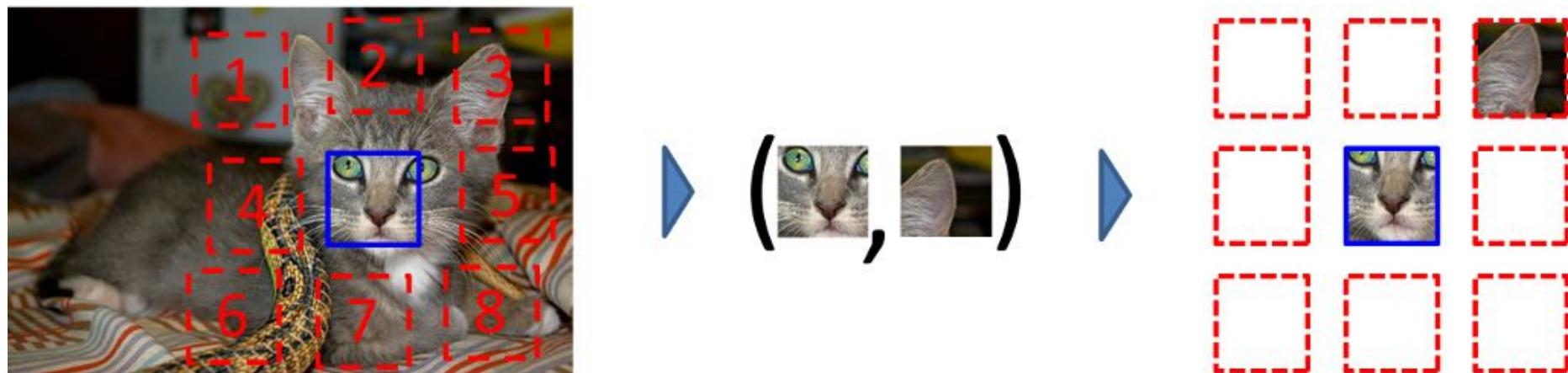


Outline

1. Unsupervised Learning
2. Self-supervised Learning
 - a. Autoencoder
 - b. Temporal regularisations
 - c. Verification
 - d. Predictive Learning
 - e. Miscellaneous

Spatial verification

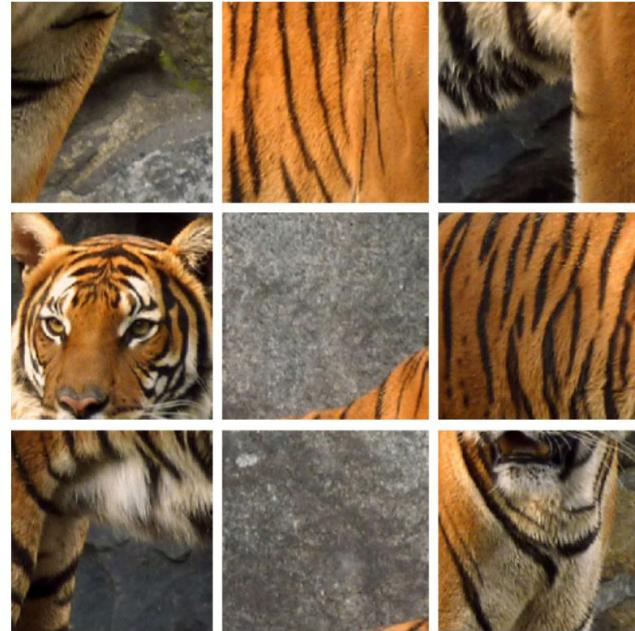
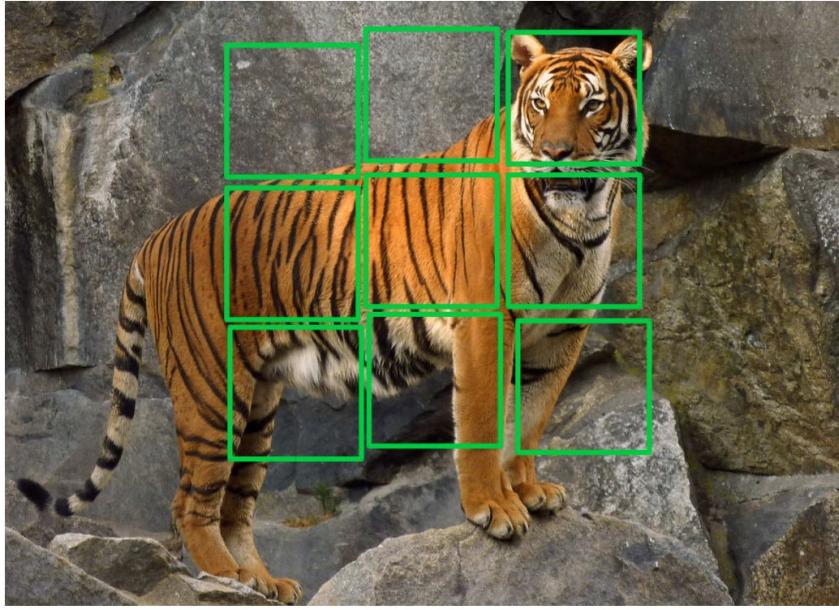
Predict the relative position between two image patches.



#**RelativePosition** Doersch, Carl, Abhinav Gupta, and Alexei A. Efros. "[Unsupervised visual representation learning by context prediction.](#)" ICCV 2015.

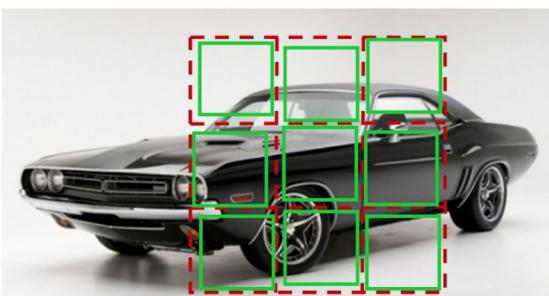
Spatial verification

Train a neural network to solve a jigsaw puzzle.



Spatial verification

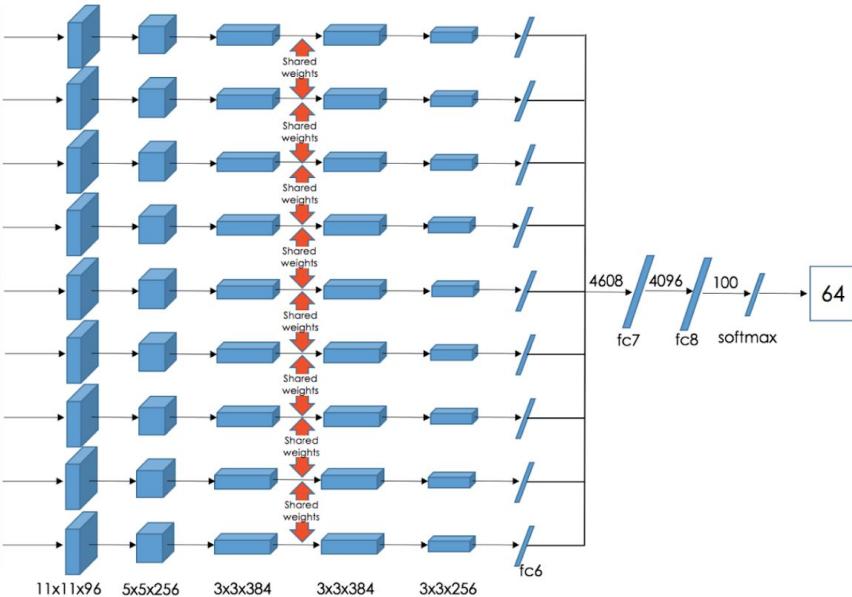
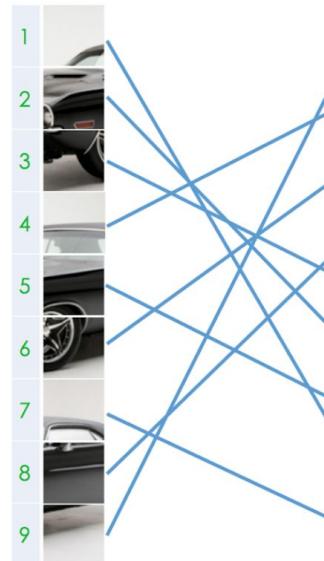
Train a neural network to solve a jigsaw puzzle.



Permutation Set

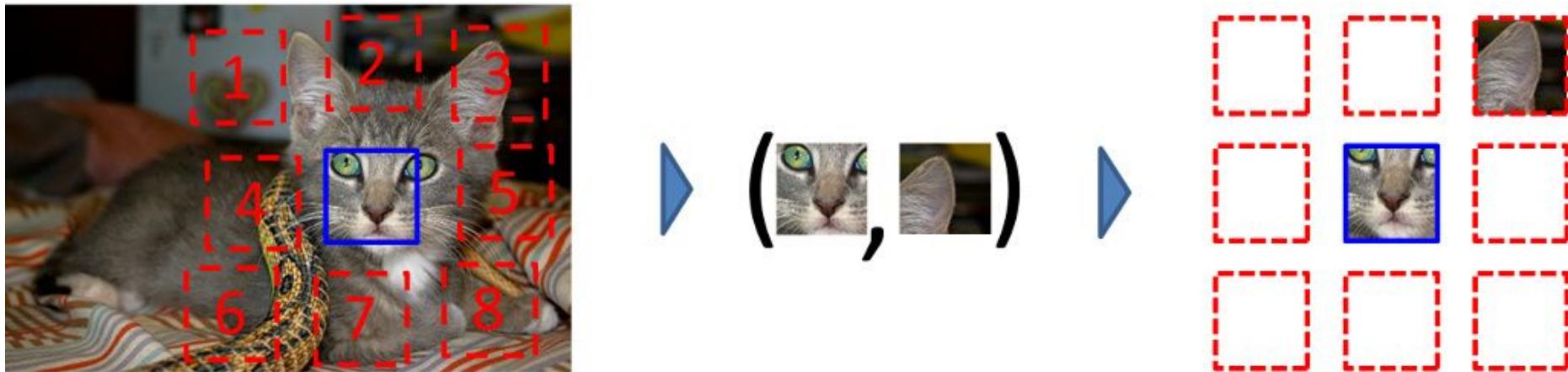
index	permutation
64	9,4,6,8,3,2,5,1,7

Reorder patches according to the selected permutation



Spatial verification

What video-specific surrogate tasks could you think about?

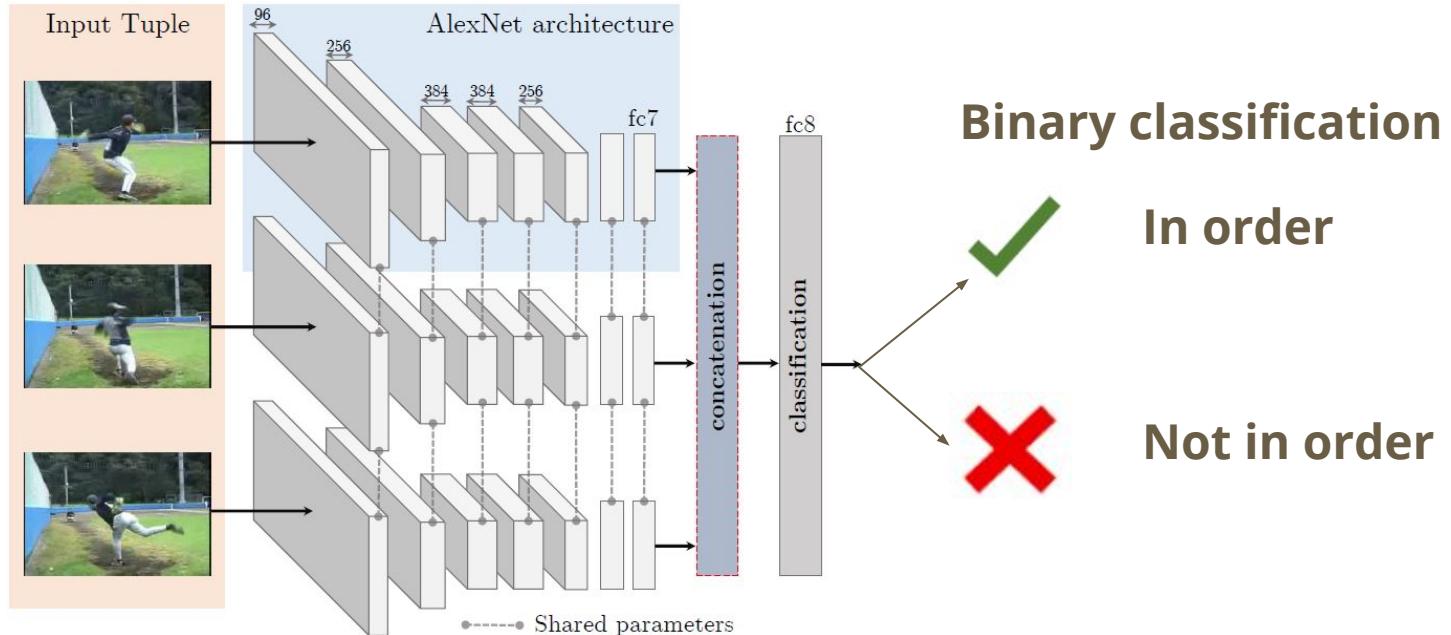


#**RelativePosition** Doersch, Carl, Abhinav Gupta, and Alexei A. Efros. "[Unsupervised visual representation learning by context prediction.](#)" ICCV 2015.

Temporal coherence

Take temporal order as the supervisory signals for learning

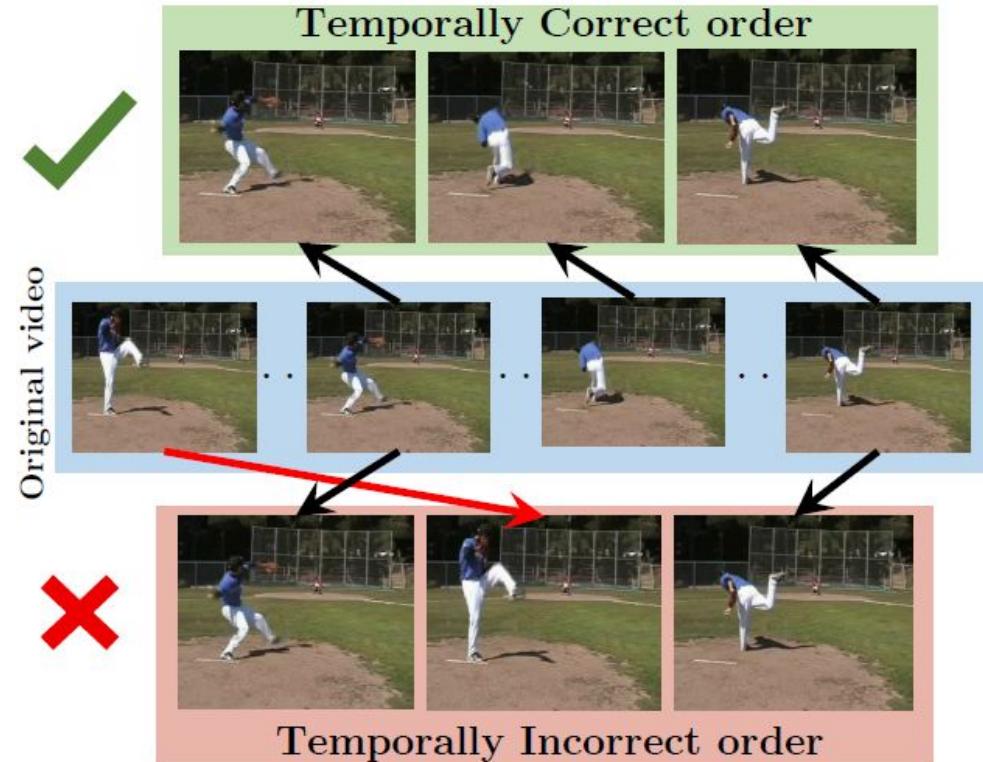
Shuffled
sequences



(Slides by Xunyu Lin): Misra, Ishan, C. Lawrence Zitnick, and Martial Hebert. ["Shuffle and learn: unsupervised learning using temporal order verification."](#) ECCV 2016. [\[code\]](#)

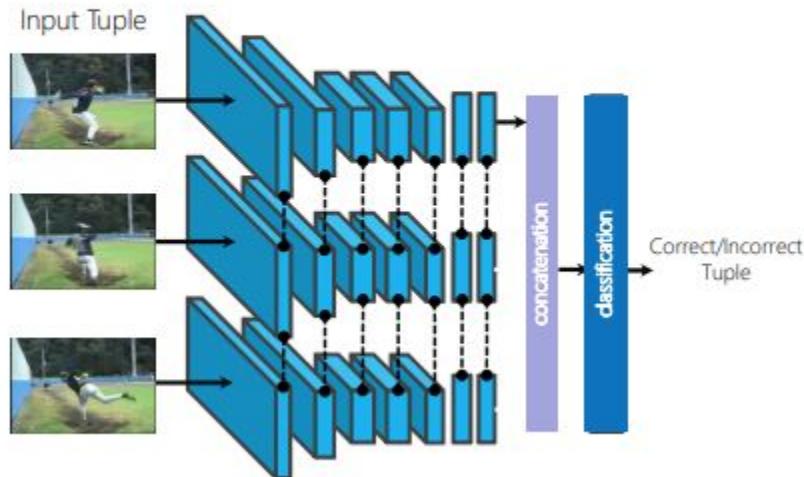
Temporal coherence

Temporal order of frames is exploited as the supervisory signal for learning.

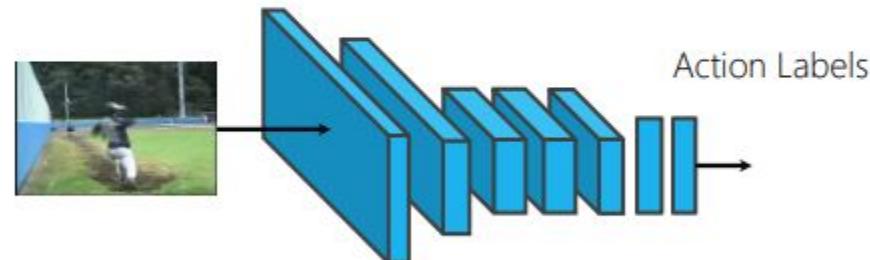


Temporal coherence

Self-supervised Pre-train



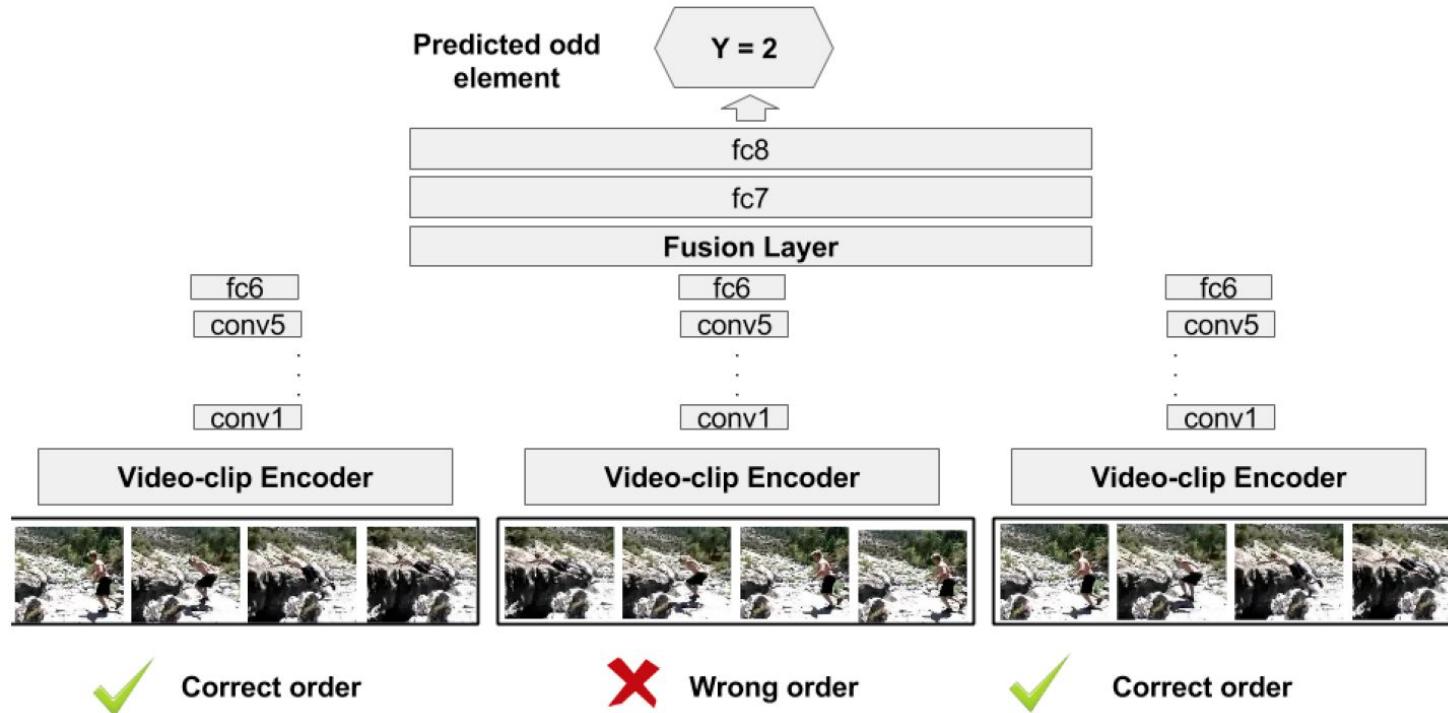
Test -> Finetune



#**shuffle&learn** Misra, Ishan, C. Lawrence Zitnick, and Martial Hebert. "[Shuffle and learn: unsupervised learning using temporal order verification.](#)" ECCV 2016. [[code](#)] ([Slides](#) by Xunyu Lin):

Temporal verification

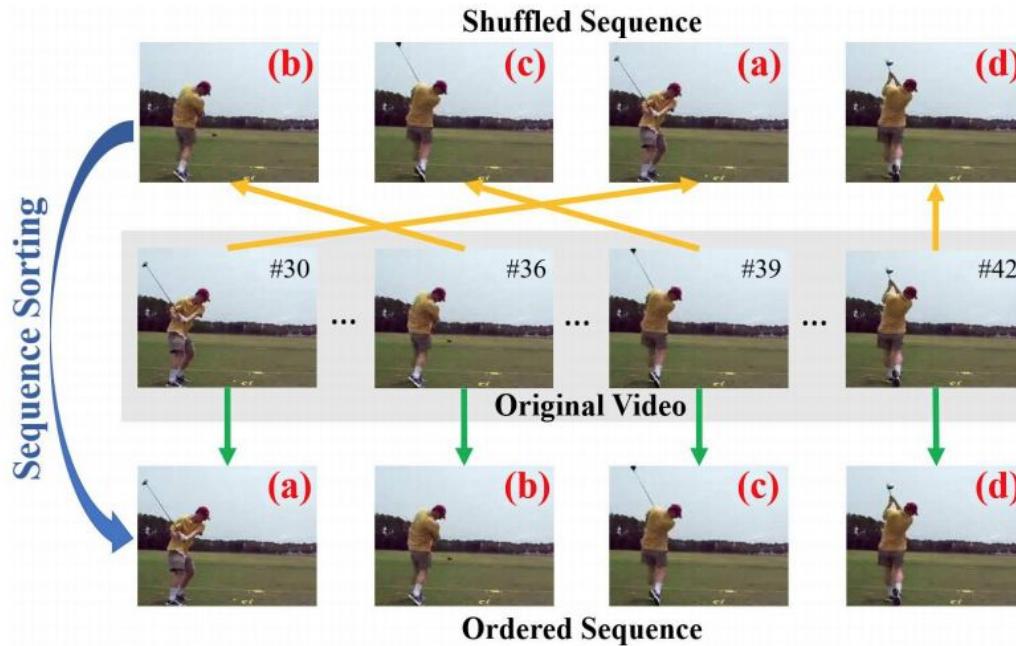
Train a network to detect which of the video sequences contains frames in the wrong order.



#Odd-one-out Fernando, Basura, Hakan Bilen, Efstratios Gavves, and Stephen Gould. "[Self-supervised video representation learning with odd-one-out networks.](#)" ICCV 2017

Temporal sorting

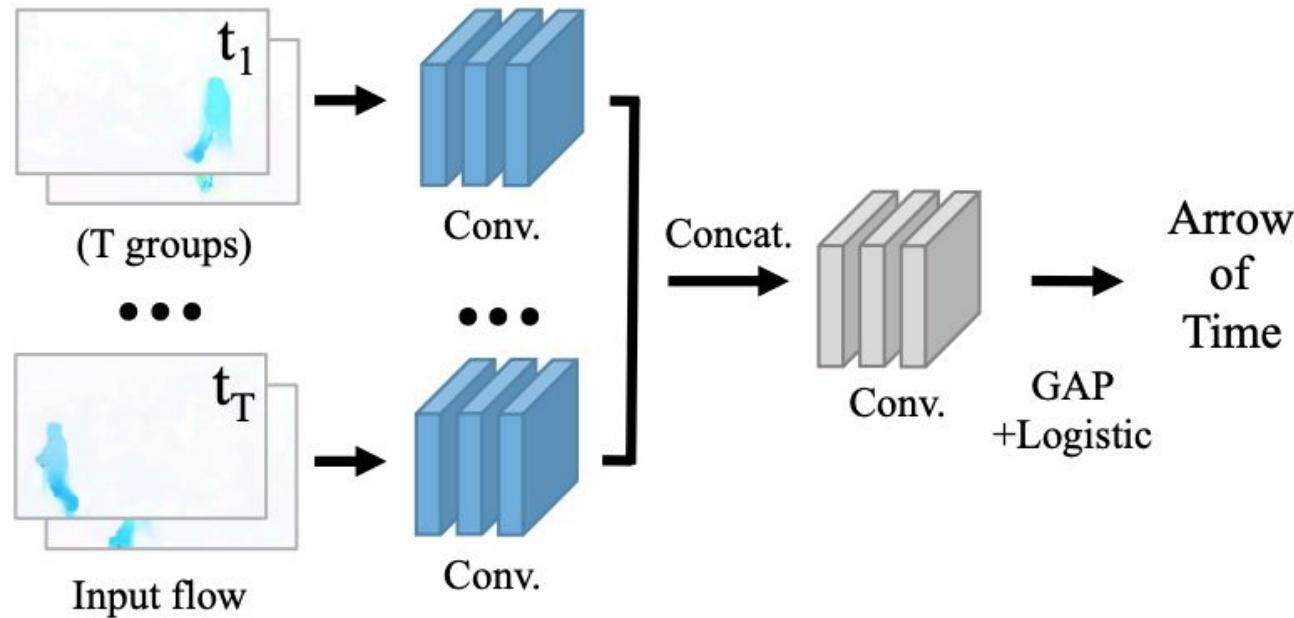
Sort the sequence of frames.



Lee, Hsin-Ying, Jia-Bin Huang, Maneesh Singh, and Ming-Hsuan Yang. "[Unsupervised representation learning by sorting sequences.](#)" ICCV 2017.

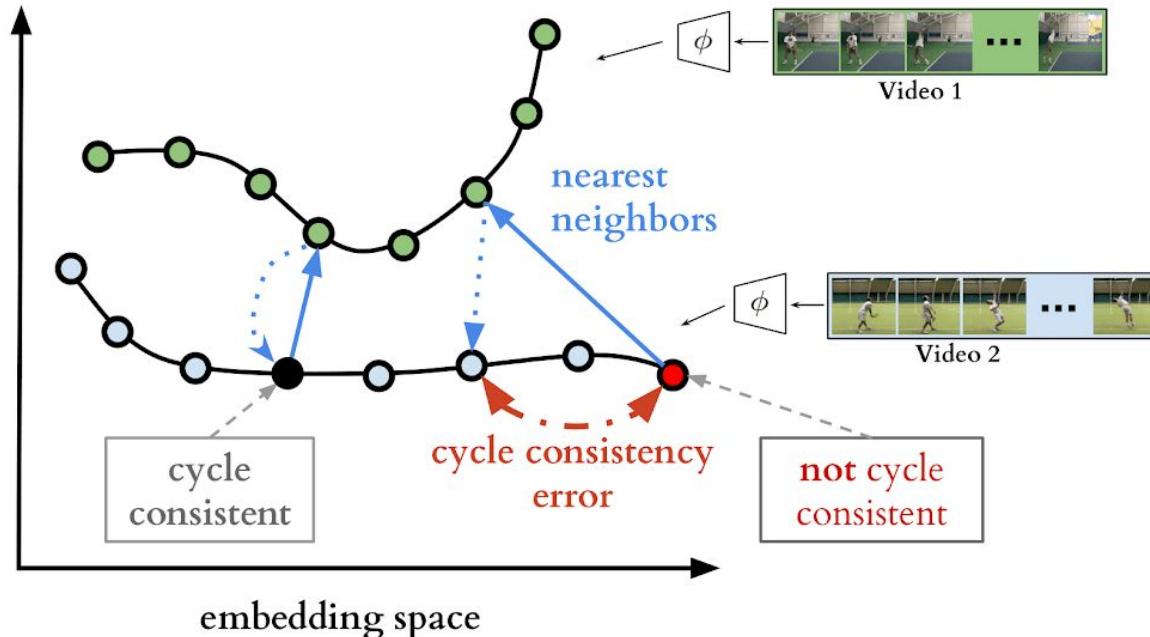
Arrow of time

Predict whether the video moves forward or backward.



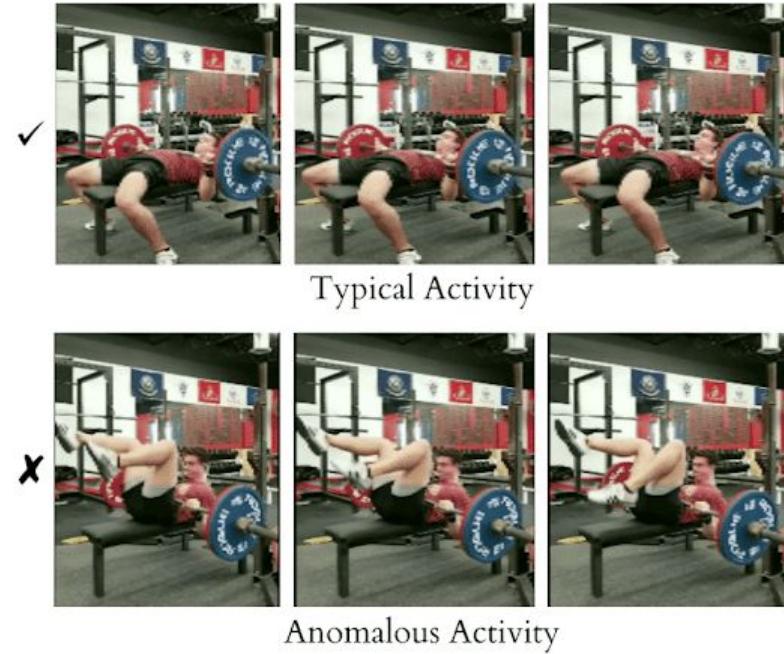
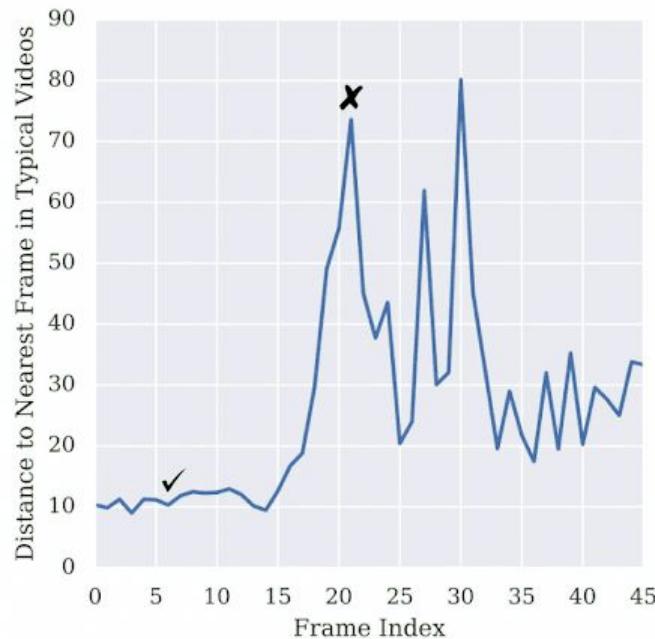
Temporal coherence

Training: Learn video frame embeddings that must satisfy a temporal alignment between videos depicting the same activity.



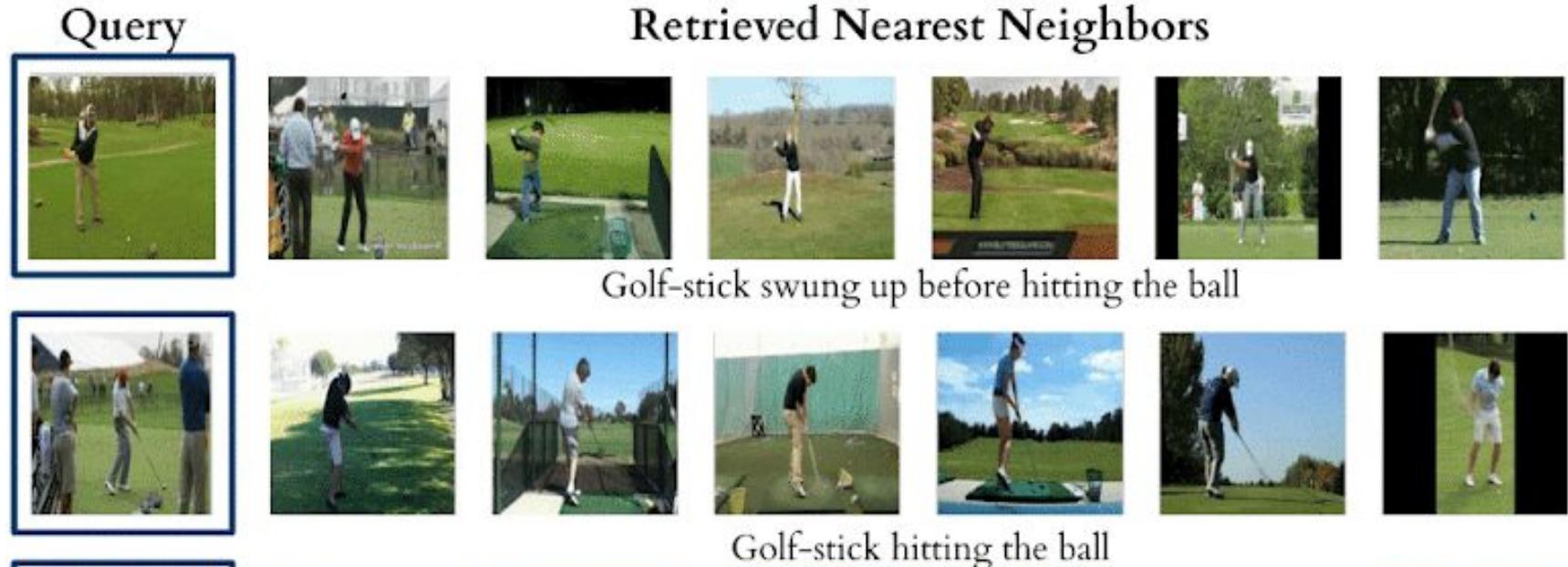
Temporal coherence

Application: Anomaly detection.



Temporal coherence

Application: Fine-grained retrieval



Outline

1. Unsupervised Learning
2. Self-supervised Learning
 - a. Autoencoder
 - b. Temporal regularisations
 - c. Temporal verification
 - d. Predictive Learning**
 - e. Miscellaneous

Predictive Learning

■ "Pure" Reinforcement Learning (cherry)

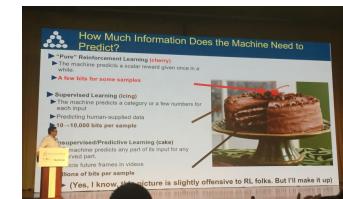
- ▶ The machine predicts a scalar reward given once in a while.
- ▶ **A few bits for some samples**

■ Supervised Learning (icing)

- ▶ The machine predicts a category or a few numbers for each input
- ▶ Predicting human-supplied data
- ▶ **10→10,000 bits per sample**

■ Unsupervised/Predictive Learning (cake)

- ▶ The machine predicts any part of its input for any observed part.
- ▶ **Predicts future frames in videos**
- ▶ **Millions of bits per sample**



Slide credit:
Yann LeCun

■ (Yes, I know, this picture is slightly offensive to RL folks. But I'll make it up)

Next word prediction (language model)

Self-supervised learning

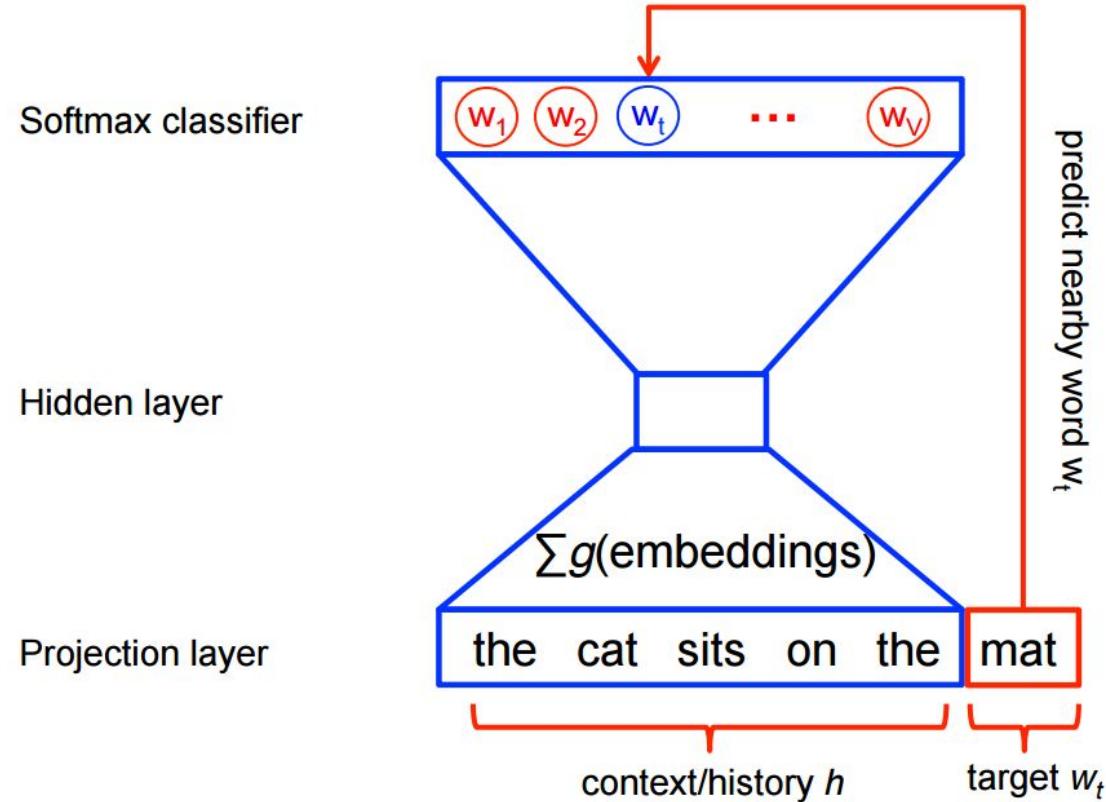
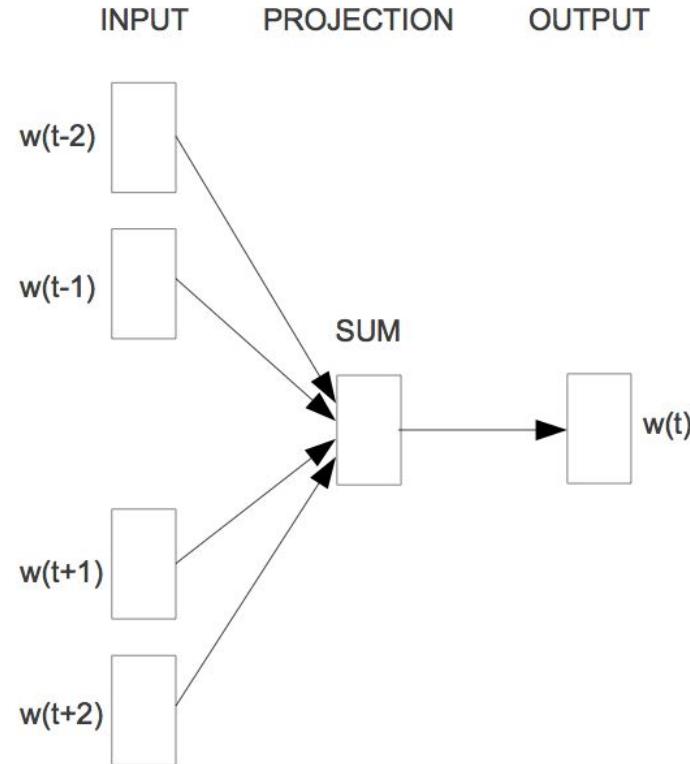


Figure:
[TensorFlow tutorial](#)

Bengio, Yoshua, Réjean Ducharme, Pascal Vincent, and Christian Jauvin. "[A neural probabilistic language model.](#)" Journal of machine learning research 3, no. Feb (2003): 1137-1155.

Missing word prediction (language model)



Self-supervised learning

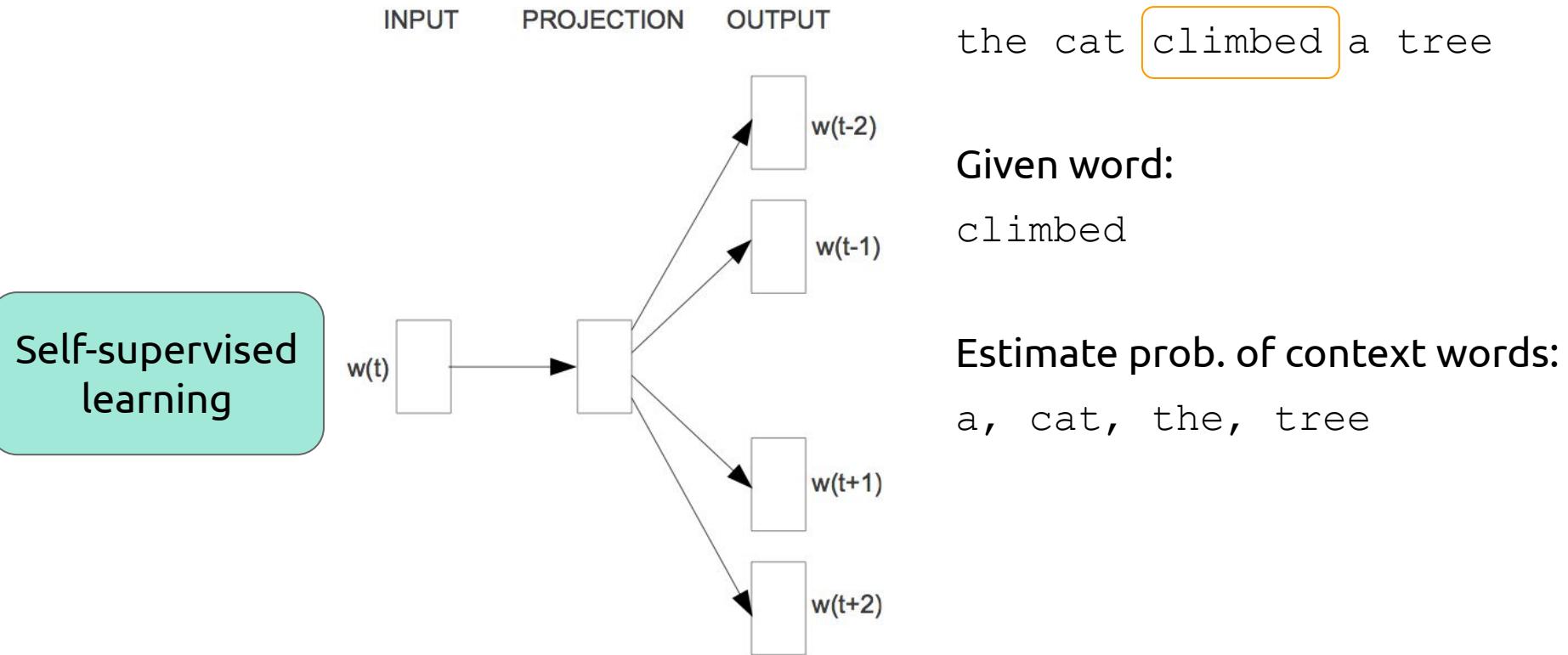
the cat climbed a tree

Given context:

a, cat, the, tree

Estimate prob. of
climbed

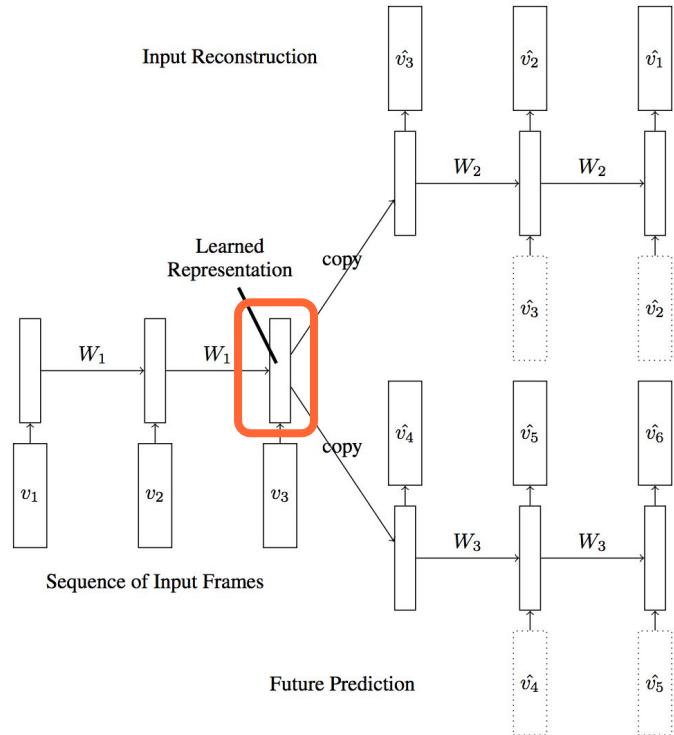
Context Words Prediction (language model)



#word2vec #skipgram Mikolov, Tomas, Ilya Sutskever, Kai Chen, Greg S. Corrado, and Jeff Dean. ["Distributed representations of words and phrases and their compositionality."](#) NIPS 2013

Frame Prediction

Learning video representations (features) by...

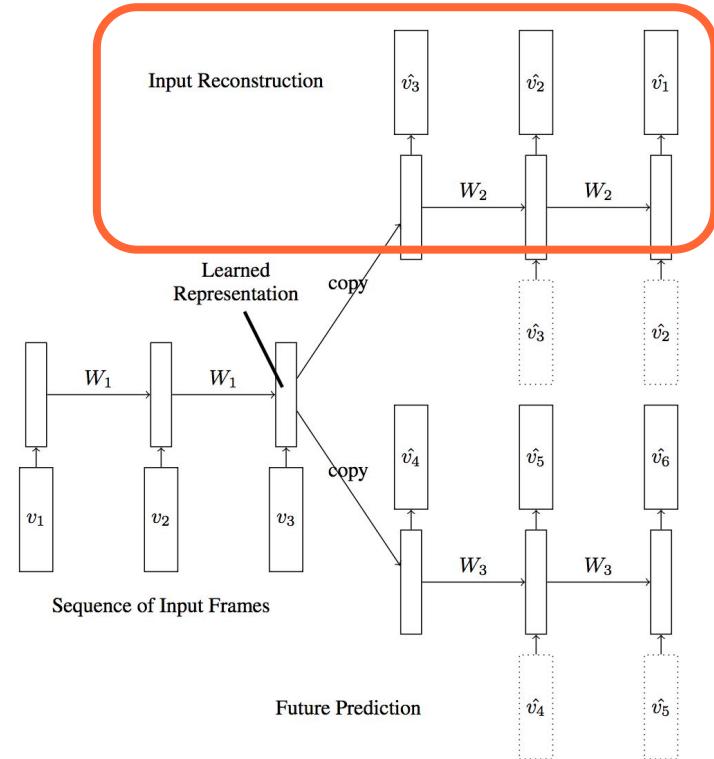
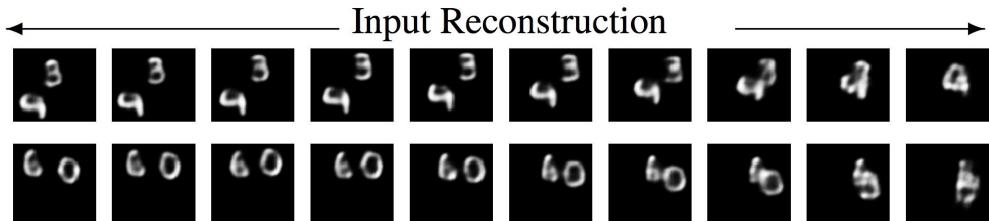
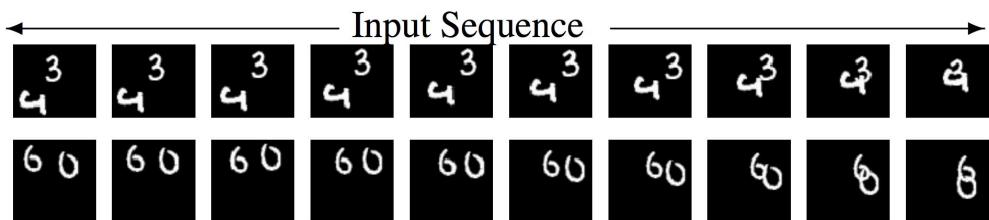


Srivastava, Nitish, Elman Mansimov, and Ruslan Salakhutdinov. "[Unsupervised Learning of Video Representations using LSTMs.](#)" In ICML 2015. [\[Github\]](#)

Frame Prediction

Learning video representations (features) by...

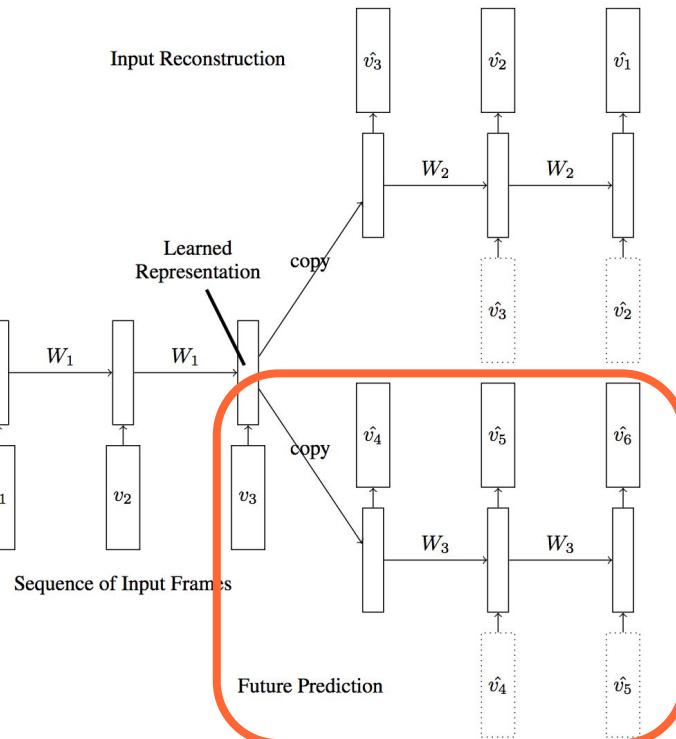
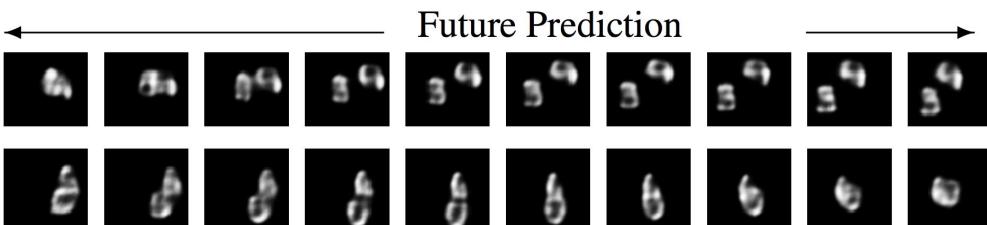
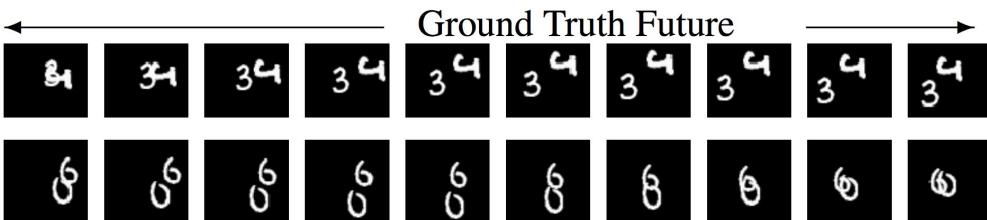
(1) frame reconstruction (AE):



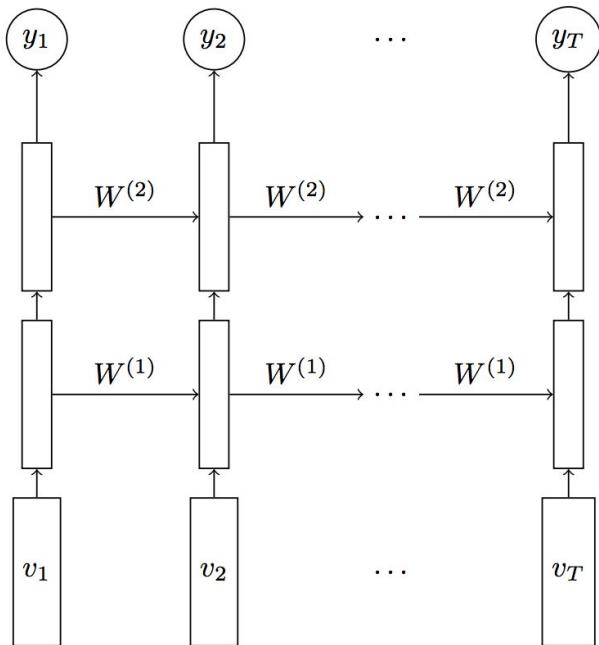
Frame Prediction

Learning video representations (features) by...

(2) frame prediction



Frame Prediction



Unsupervised learned features (lots of data) are fine-tuned for activity recognition (small data).

Model	UCF-101	UCF-101	HMDB-51
	RGB	1-frame flow	RGB
Single Frame	72.2	72.2	40.1
LSTM classifier	74.5	74.3	42.8
Composite LSTM	75.8	74.9	44.1
Model + Finetuning			

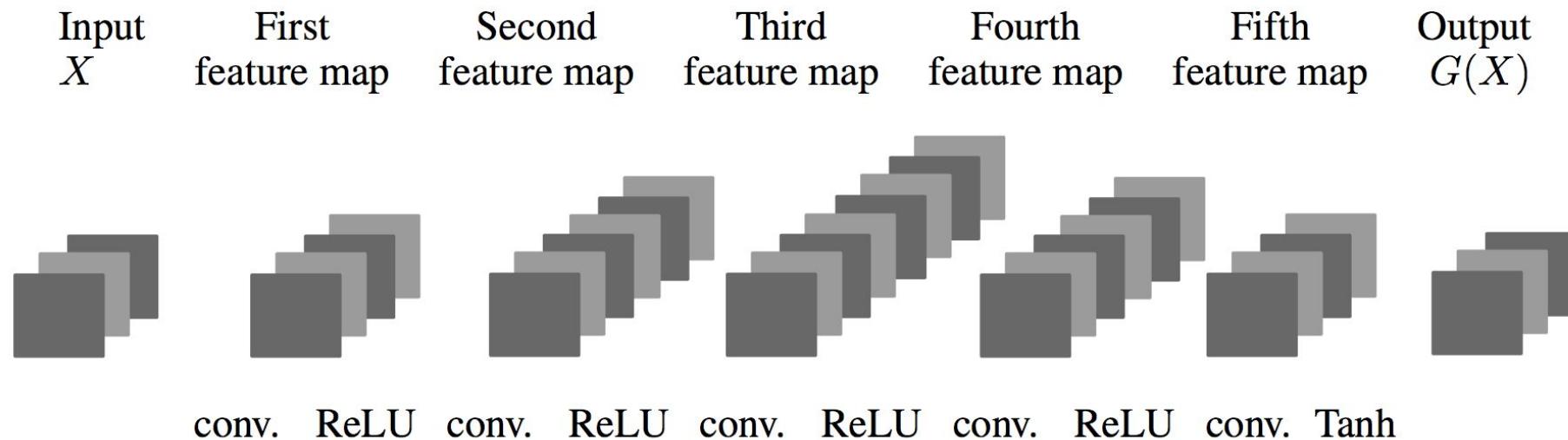
Table 1. Summary of Results on Action Recognition.

Figure 6. LSTM Classifier.

Frame Prediction

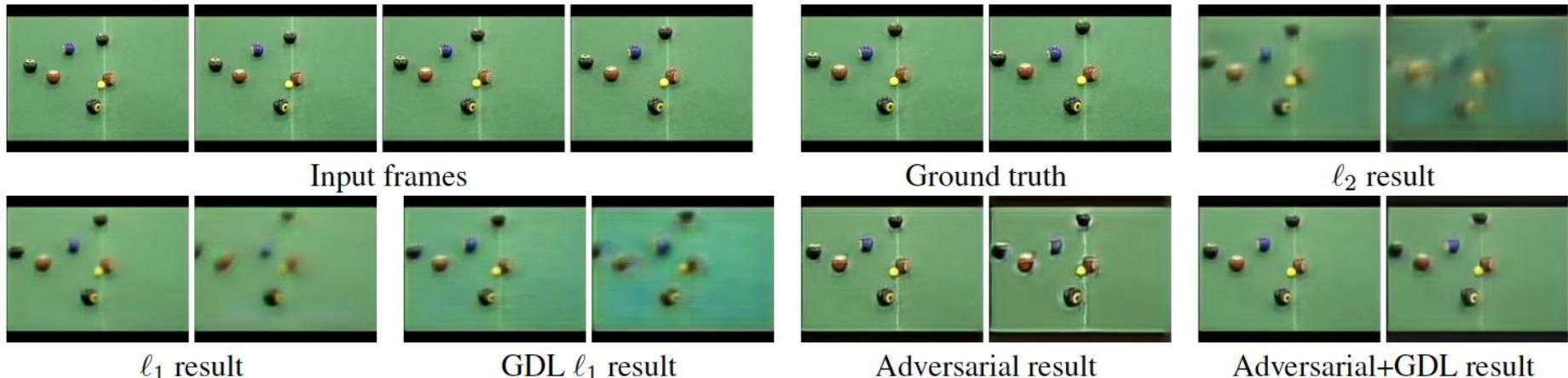
Video frame prediction with a ConvNet.

Figure 1: A basic next frame prediction convnet



Frame Prediction

The blurry predictions from MSE (L1 loss) are improved with multi-scale architecture, **adversarial training** and an image gradient difference loss (GDL) function.



Frame Prediction



Mathieu, Michael, Camille Couprie, and Yann LeCun. "[Deep multi-scale video prediction beyond mean square error.](#)"
ICLR 2016 [\[project\]](#) [\[code\]](#)

Frame Prediction + Disentangled features

The model learns to disentangle (“separate”) the visual features that correspond to the:

Object Content
(class)



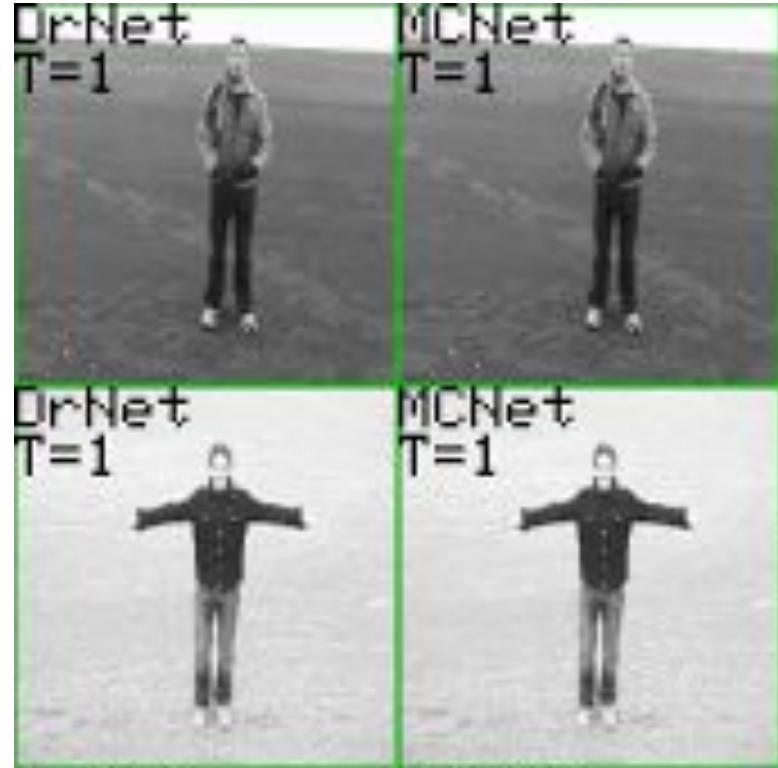
Object Pose
(wrt the camera)



Frame Prediction + Disentangled features

100 step video generation on KTH dataset where:

- green frames indicate reconstructions
- red frames indicate frame predictions.



#DrNet Denton, Emily L. "[Unsupervised learning of disentangled representations from video.](#)" NIPS 2017.

#MCNet R. Villegas, J. Yang, S. Hong, X. Lin, and H. Lee. Decomposing motion and content for natural video sequence prediction. In ICLR, 2017

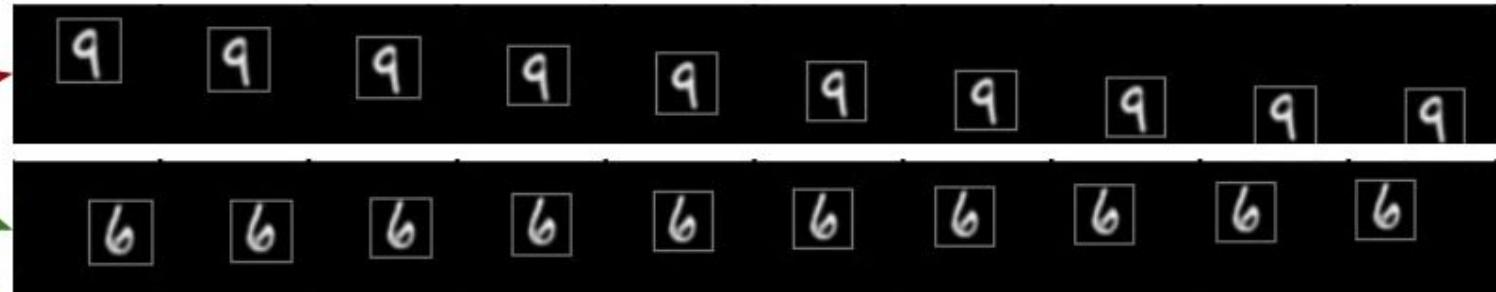
Frame Prediction + Disentangled features

Decompositional Disentangled Predictive Auto-Encoder (DDPAE):

Ground truth



Decompose

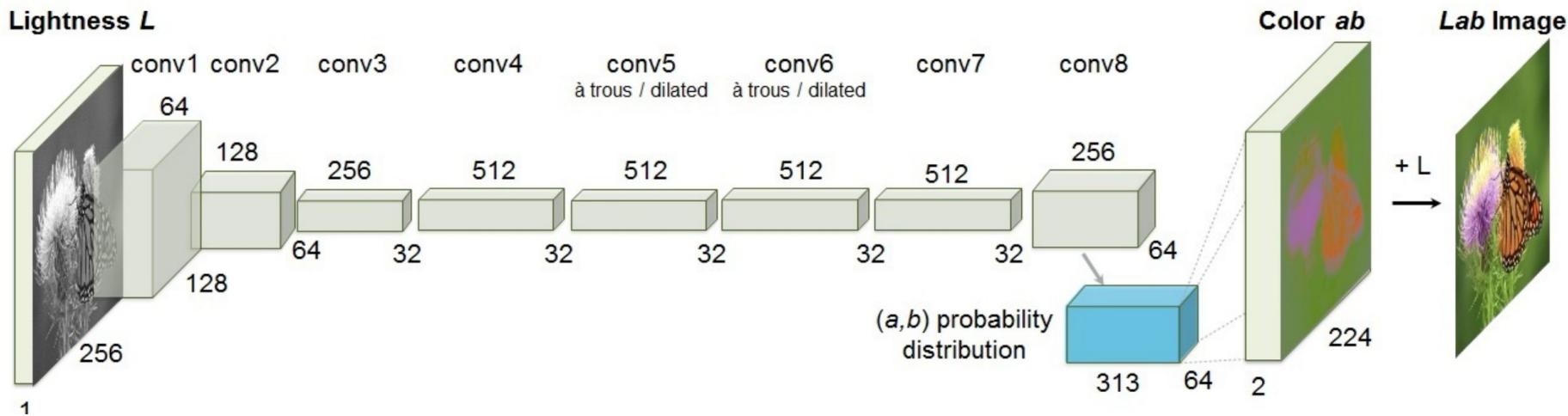


Outline

1. Unsupervised Learning
2. Self-supervised Learning
 - a. Autoencoder
 - b. Temporal regularisations
 - c. Temporal verification
 - d. Frame Prediction
 - e. Miscellaneous**

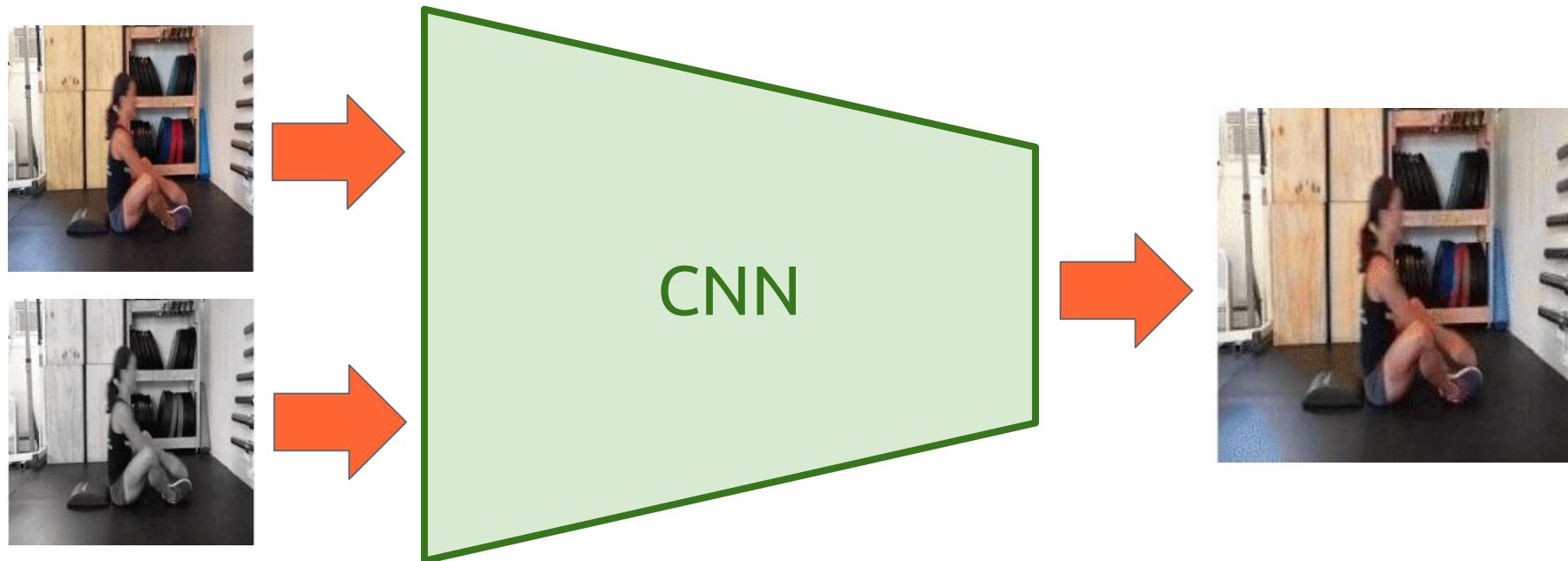
Colorization (still image)

A NN is trained to colorize a gray scale image..



Colorization (video)

A NN is trained to colorize a video frame, given the color of the first frame of the video sequence.



Vondrick, Carl, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. "[Tracking emerges by colorizing videos.](#)" ECCV 2018. [\[blog\]](#)

Colorization (video)

Training: A NN is trained to colorize a video frame, given the color of the first frame of the video sequence.

Reference Frame



What color is this?



Vondrick, Carl, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. "[Tracking emerges by colorizing videos.](#)" ECCV 2018. [\[blog\]](#)

Colorization (video)

Inference: Video Object tracking with the learned pixel embeddings.



Vondrick, Carl, Abhinav Shrivastava, Alireza Fathi, Sergio Guadarrama, and Kevin Murphy. "[Tracking emerges by colorizing videos.](#)" ECCV 2018. [\[blog\]](#)

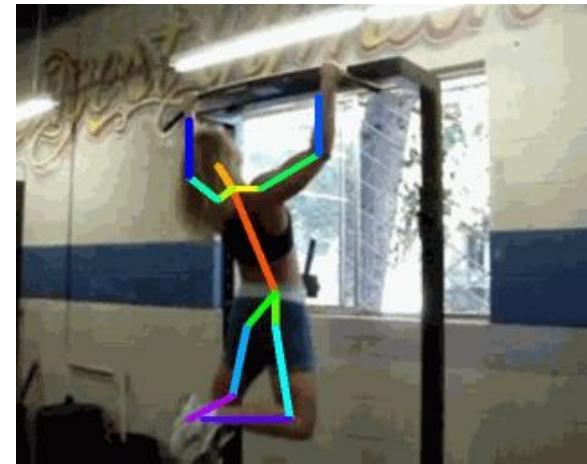
Cycle-consistency

Training: Tracking backward and then forward to learn from the divergence of the predictions.



Cycle-consistency

Inference: Video Object tracking with the learned pixel embeddings.

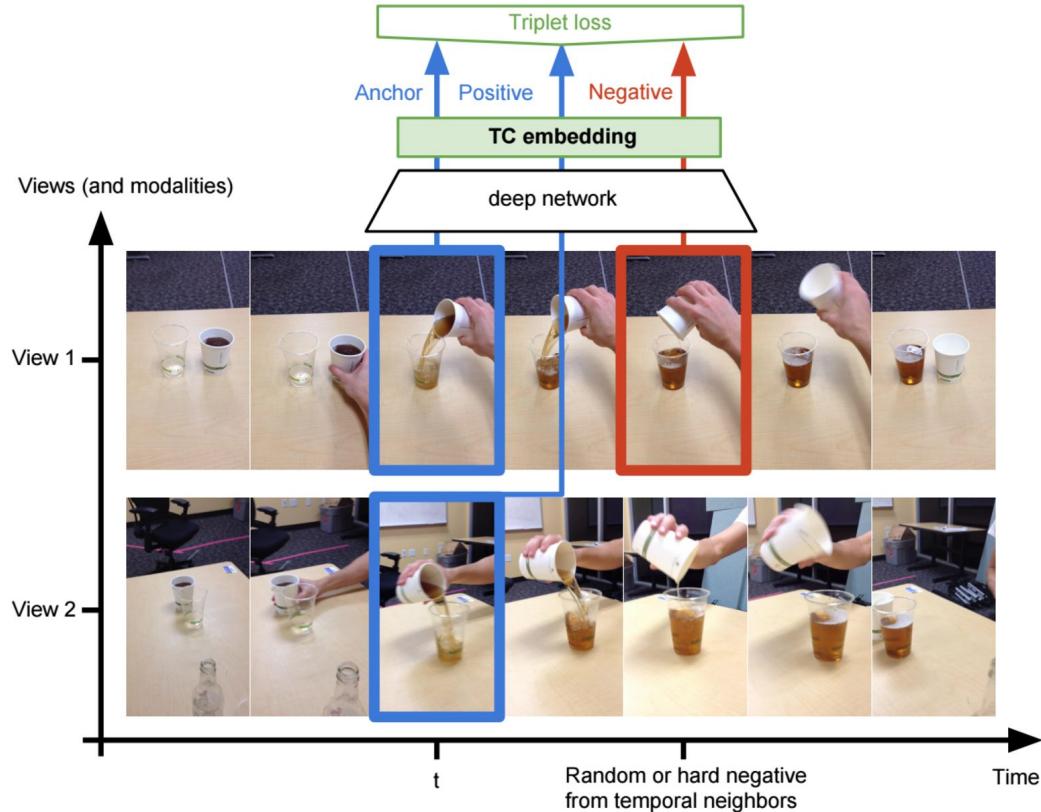


Multiview + Time

Training: Learn features that

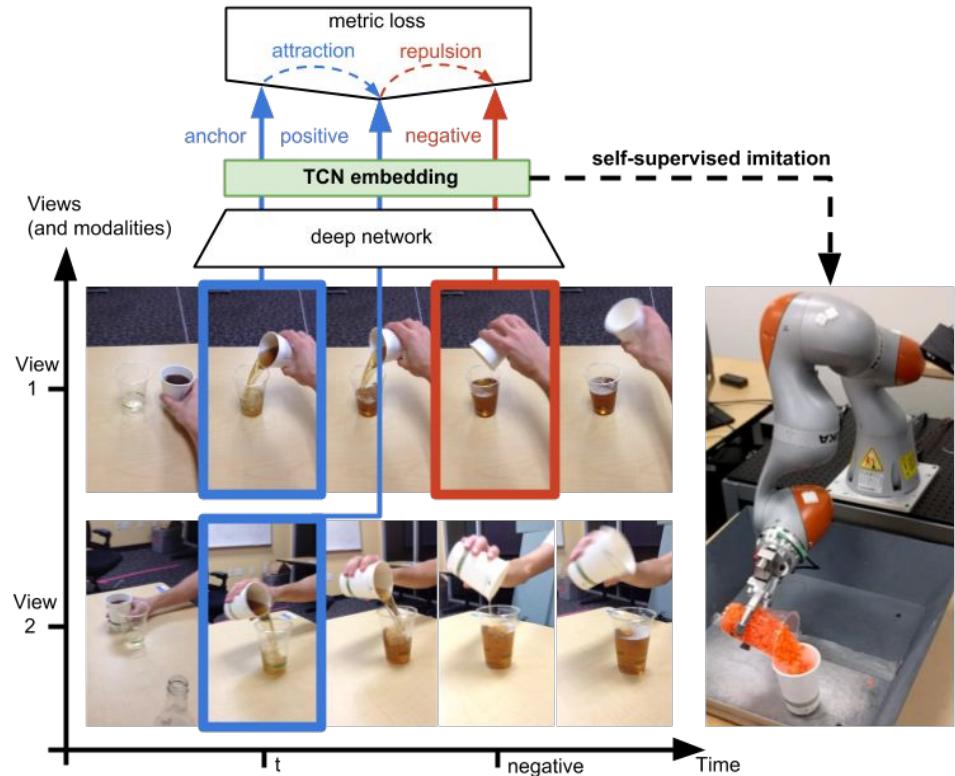
are:

- Invariant to viewpoint.
- Very sensitive to temporal ordering



Multiview + Time

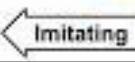
Inference: Imitation learning for a robotic arm.



Learning to imitate, from video, without supervision



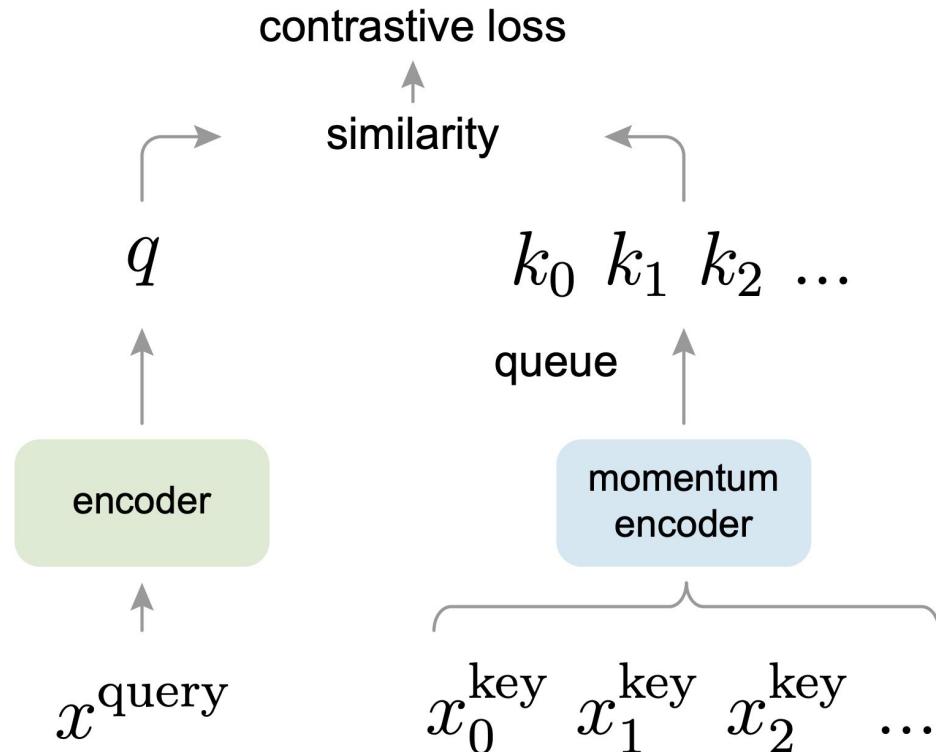
3rd-person observation



Learned policy

Momentum Contrast (MoCo)

“This suggests that the gap between unsupervised and supervised representation learning has been largely closed in many vision tasks.!“



Outline

1. Unsupervised Learning
2. Self-supervised Learning
 - a. Autoencoder
 - b. Temporal regularisations
 - c. Temporal verification
 - d. Frame Prediction
 - e. Miscellaneous

Bonus track



@DocXavi



Master in
Computer Vision
Barcelona

[\[http://pagines.uab.cat/mcv/\]](http://pagines.uab.cat/mcv/)





Xavier Giro-i-Nieto
xavier.giro@upc.edu

Associate Professor
Universitat Politècnica de Catalunya



Module 6 - Day 9 - Lecture 2
**Self-supervised
Audiovisual Learning**
4th April 2019

Self-supervised Audiovisual Learning [[slides](#)] [[video](#)]

Suggested reading



Alex Graves, Kelly Clancy, "[Unsupervised learning: The curious pupil](#)". Deepmind blog 2019.

Suggested talks



Demis Hassabis (Deepmind)



Aloysha Efros (UC Berkeley)

Use audio and visual features

facebook research

What can be learnt by watching and listening to videos?

- Good representations
 - Visual features
 - Audio features
- Intra- and cross-modal retrieval
 - Aligned audio and visual embeddings
- "What is making the sound?"
 - Learn to localize objects that sound

Diagram: Two parallel processing paths. The top path takes a 'single frame' as input and processes it through a 'visual subnetwork' to produce a '200D' vector. The bottom path takes a '1 s' audio signal as input and processes it through an 'audio subnetwork' to produce a '200D' vector. These two vectors are then compared by a 'correspond' module to output a 'yes/no' response.

"Objects that Sound", Arandjelović and Zisserman, ICCV 2017 & ECCV 2018

Andrew Zisserman (Oxford/Deepmind)



Yann LeCun (Facebook AI)

Questions ?

Undergradese

What undergrads ask vs. what they're REALLY asking

"Is it going to be an open book exam?"

Translation: "I don't have to actually memorize anything, do I?"

"Hmm, what do you mean by that?"

Translation: "What's the answer so we can all go home."

"Are you going to have office hours today?"

Translation: "Can I do my homework in your office?"

"Can i get an extension?"

Translation: "Can you re-arrange your life around mine?"

"Is this going to be on the test?"

Translation: "Tell us what's going to be on the test."

"Is grading going to be curved?"

Translation: "Can I do a mediocre job and still get an A?"

