

**VISUAL ANALYSIS &** 

PERCEPTION LAB



# Visual and Multimodal Learning Systems for Human Behavior Understanding



UNIVERSITAT DE BARCELONA Dr. Sergio Escalera

www.sergioescalera.com

sergio@maia.ub.es



ELLIS Fellow, IAPR Fellow, AAIA Fellow

### University of Barcelona





# University of Barcelona – Faculty of Mathematics and Informatics



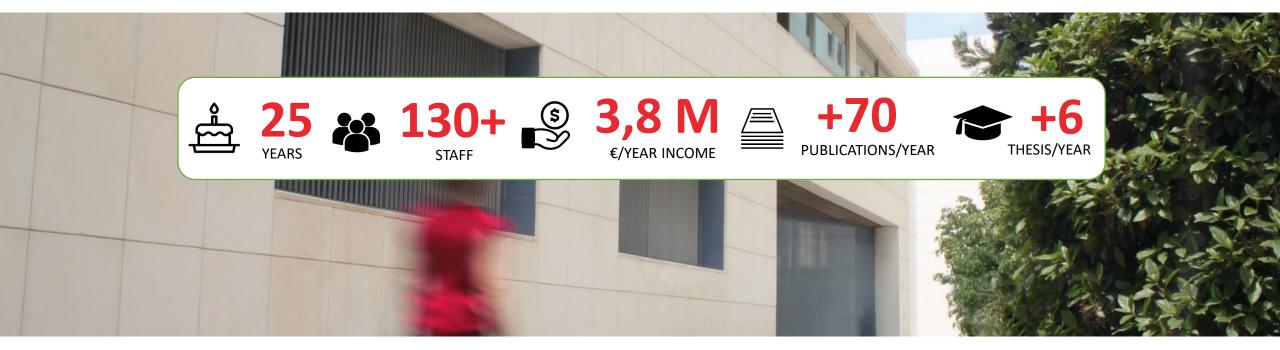


- Informatics degree, Prof.
   Sergio Escalera head of Informatics
- Mathematics degree
- Master involvement (interuniversitary):
  - Al
  - Data Science
  - Computer Vision
  - Behavioral Data Science

Computer Vision Center (CVC-UAB)

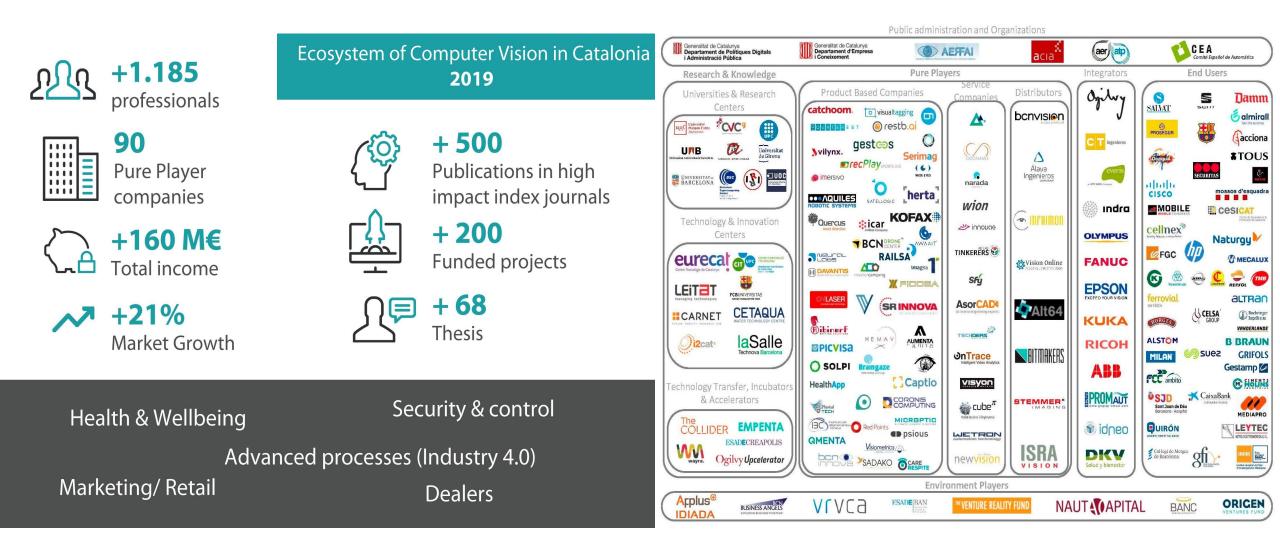
Only Center in Europe fully devoted to Computer Vision

- CVC is a legally independent non-profit institution founded in 1995, belonging to the Catalan CERCA network. Located in Bellaterra (Barcelona). Dedicated to research, technology transfer, training, and outreach.
- 38 senior researchers + 52 students (2019)
- 45 JCR indexed journals, 65 international conference papers, 12 thesis (2019).
- HuPBA, head Prof. Sergio Escalera, 1 of the 8 strategic research lines





## Computer Vision ecosystem en Catalonia





esearch fields Computer Vision Machine Learning Social Signal Processing Affective Computing Personality Computing	<ul> <li>Application domains:</li> <li><u>eHealth and well-being</u></li> <li>Security</li> <li>Smart cities</li> <li>Leisure</li> </ul>	LARS VEARS VEARS 13 VEARS 15+ MEMBER MEMBER
<ul> <li>Research lines</li> <li>Deep Learning</li> <li>Domain Adaptation</li> <li>Bias and fairness</li> <li>Explainability and interpretability</li> <li>Spatio-temporal modeling &amp; video understanding</li> <li>Multi-modality, multi-view &amp; multi-task learning</li> <li>Attention mechanisms</li> </ul>	<ul> <li>Main collaborators</li> <li>Medical researchers: <ul> <li>Psychologists</li> <li>Psychiatrists</li> <li>Neurologists</li> </ul> </li> <li>National (UAB, UOC), and international universities (<u>AAU</u>, Berkeley, Boston)</li> <li>Companies (Google, Microsoft, Disney, Amazon, NVIDIA, and Facebook)</li> </ul>	TOTAL PUBLICATIONS

# Examples of transfered technology





Visual and multimodal monitoring in neurorehabilitation



Risk event monitoring



ADHD and mental health monitoring

Coaching for the elder



SLR



Physiotherapy and fitness



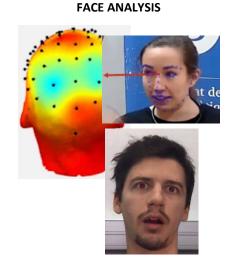


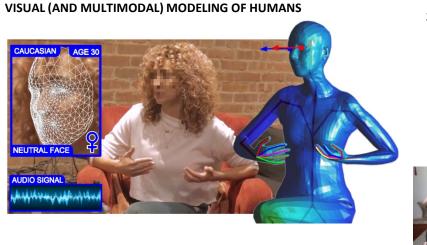
Virtual try ons

#### 7



## Overview of current research lines in LAP





3D (& 4D) POSE, SHAPE, TEXTURE (IN 3D AND FROM 2D)



**BEHAVIOR ANALYSIS** 

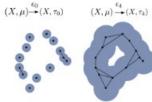


UNDERSTANDING AND EXPLAINING HUMAN BEHAVIOR (Affective & Personality Computing) -INTERPRETABILITY & EXPLAINABILITY

-FAIRNESS



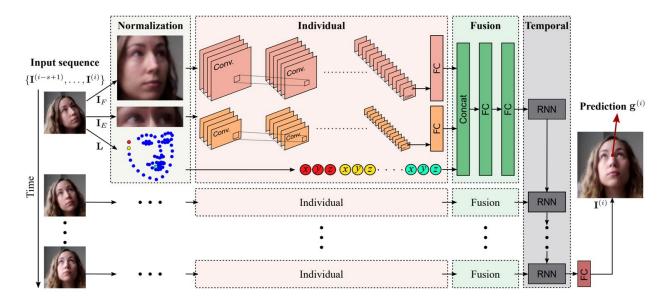
BIAS ANALYSIS VISUALIZATION



INTERPRETING AND EXPLAINING LEARNING

# Face analysis

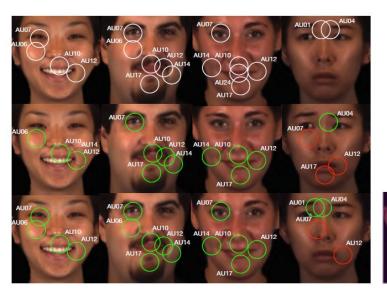
- Gaze
- Facial Action Units
- Multimodal app.



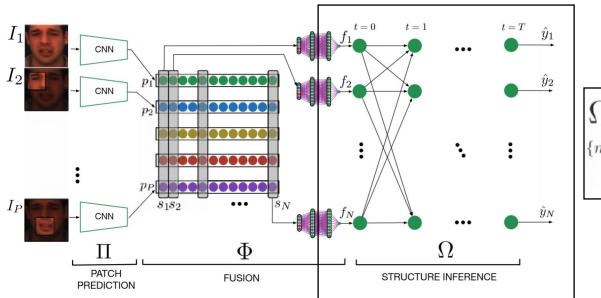


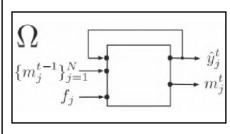
# Face analysis

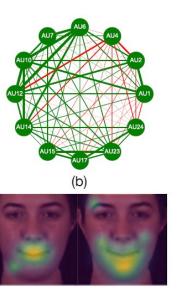
- Gaze
- Facial Action Units
- Multimodal app.

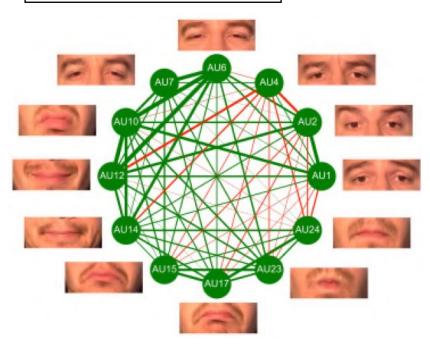


Corneanu et.al. ECCV 2018



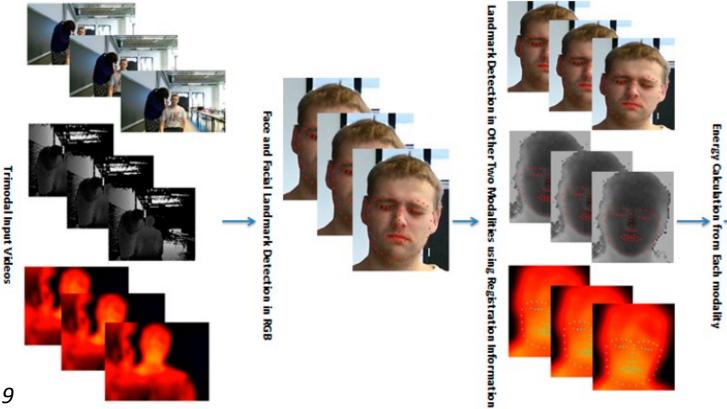






# Face analysis

- Gaze
- Facial Action Units
- Multimodal app.



Faces and Gestures, 2019

Fusing the

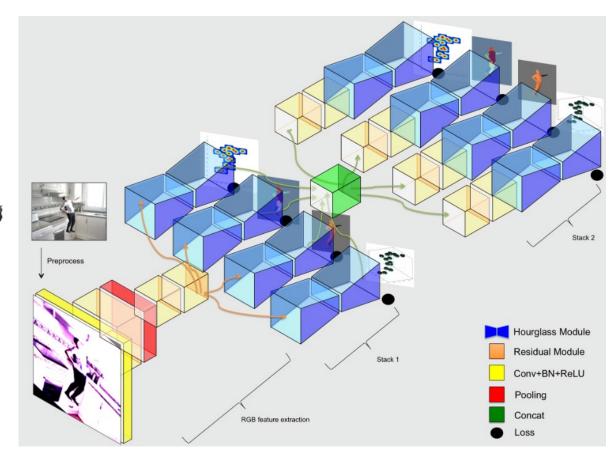
ults

for

dal Pain Recognition

# Body Analysis • Posture and multimodal H,RGB

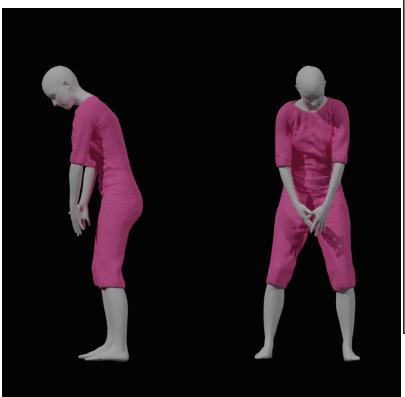
Sanchez et.al. FG 2019



Shanxin Yuan et.al. CVPR 2018

# Body Analysis

• Posture and multimodal



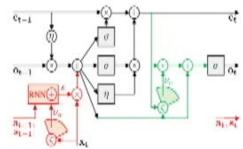
Madadi et.al. PR 2020 Bertiche et.al. ICCV 2021, SIGGRAPH ASIA 2021 & 2022

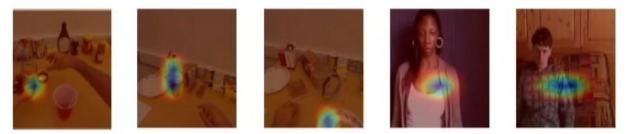
## Behavior



## Long Short-Term Attention

State-of-the-art recognition performance

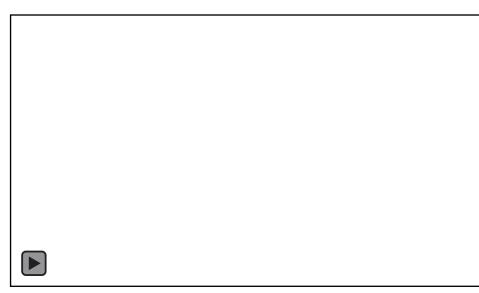




Sudhakaran et.al. CVPR 2019, CVPR 2020 EPIC Kitchens 2019, 2020, Top-3 winning solutions

# Behavior

- Moving towards affective HCI and assistive technology
- Regressing self-reported personality from dyadic multimodal video data



Cross-attention Cross-attention encoder encoder Cross-subject attention QA KA VA  $Q \not\models K \not\models V \not\models$ Self-attention Self-attention encoder encoder Q 4 K 4 V 4  $Q \downarrow K \downarrow V \downarrow$ Cross-attention Cross-attention encoder encoder Cross-modal attention  $Q \downarrow K \downarrow V$ Q**↓** K**≜** V4 Self-attention Self-attention encoder encoder Positional  $Q \downarrow K \downarrow V$  $Q \downarrow K \downarrow V \downarrow$ Metadata encoding  $\bigcirc \rightarrow \bigcirc \checkmark \bigcirc$  $\bigcirc \bullet \bullet \bullet \bigstar$ Video Audio Video Audio embeddinas embeddings embeddings embeddings Video chunks Audio chunks Video chunks Audio chunks

Participant A OCEAN

Feed-forward

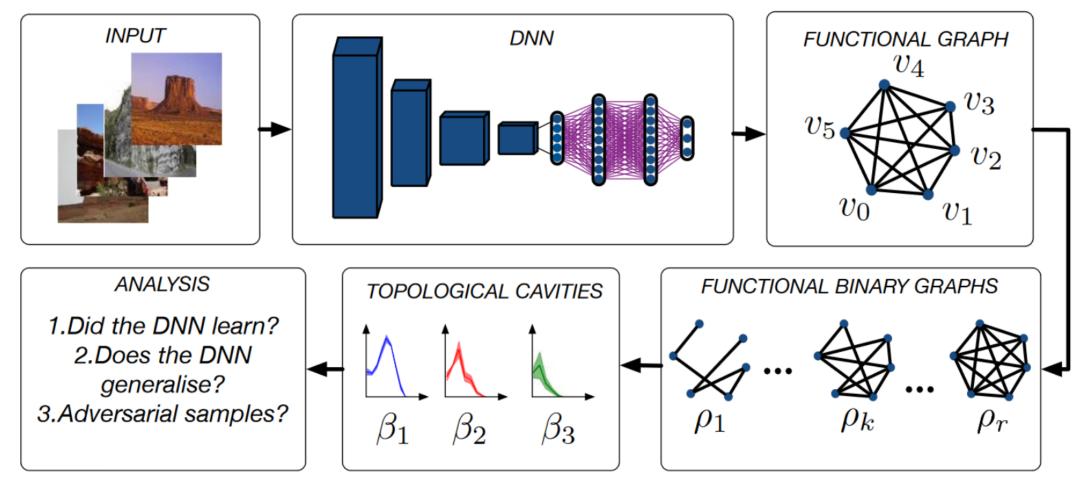
**Participant B** 

OCEAN

Feed-forward

Dyadformer David Curto, Albert Clapés, ICCV 2021

## Interpretability

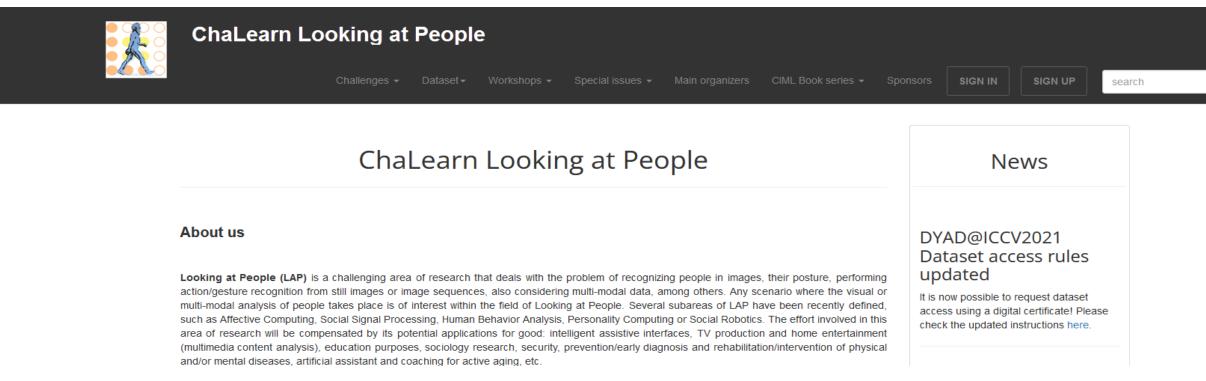


Corneanu et.al. CVPR 2019, CVPR 2020

CVPR best paper award nominee

## ChaLearn

- Non-profit organization. Berkeley. We organize challenges to stimulate research in this field. The web sites of past challenges remain open for post-challenge submission as ever-going benchmarks Promoting open data, educational materials, and challenge organization. Link with ChaLearn and Codalab initiatives.
- President: Isabelle Guyon, Google, Université Paris-Saclay, France
- Vice-president: Sergio Escalera, University of Barcelona, Spain
- Involved in the implementation of Neurips Competition track and Neurips Benchmarking track

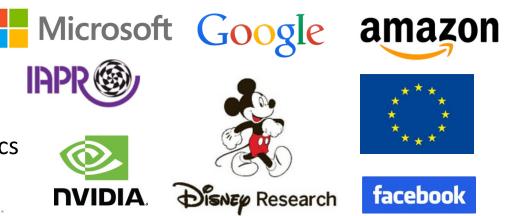






### Main sponsors

- > 25 new datasets
- > 25 organized challenges at CVPR, ICCV, ECCV, NeurIPS, ...
- > 25 organized workshops at CVPR, ICCV, ECCV, NeurIPS, ...
- > 10 organized Special Issues to related workshop/challenge topics



### The Springer Series on Challenges in Machine Learning ≌ Springer

### The Springer Series on Challenges in Machine Learning

Series Editors: **Escalante**, Hugo Jair, **Guyon**, Isabelle, **Escalera**, Sergio ISSN: 2520-131X

ABOUT THIS SERIES

~10 edited volumes up to date

The books in this innovative series collect papers written in the context of successful competitions in machine learning. They also include analyses of the challenges, tutorial material, dataset descriptions, and pointers to data and software. Together with the websites of the challenge competitions, they offer a complete teaching toolkit and a valuable resource for engineers and scientists.

# Calab

#### Accelerating reproducible computational research.

#### The CodaLab Team



Percy Liang is an assistant professor of Computer Science at Stanford University. His primary research areas are machine learning and natural language processing. He leads the development of CodaLab in close collaboration with Microsoft Research and the rest of the community.



Isabelle Guyon is full professor at UPSud University Paris-Saclay and president of ChaLearn a non-profit organization dedicated to running machine learning competitions. Her research interested include automatic machine learning, transfer learning, and causal discovery. Isabelle served as an advisor in the development of the CodaLab competition platform and pioneered the implementation of several challenges on Codalab.



Evelyne Viegas is a Director at Microsoft Research responsible for the outreach artificial intelligence program. She leads the CodaLab project working in collaboration with Isabelle Guyon, Percy Liang and the machine learning and artificial intelligence communities.



Sergio Escalera is adjunct professor at Universitat Oberta de Catalunya, Aalborg University, and Dalhousie University and a member of the Visual and Computational Learning consolidated research group of Catalonia and a member of the Computer Vision Center at UAB. He is series editor of The Springer Series on Challenges in Machine Learning. He is Editor-in-Chief of American Journal of Intelligent Systems and editorial board member of more than 5 international journals. He is vice-president of ChaLearn Challenges in Machine Learning, leading ChaLearn Looking at People events.

♥ Why GitHub? ∨ Tear	n Enterprise Explore ~	Marketplac	e Pricing $\sim$			Search		Sign in (	Sign up
🗟 codalab / codalab-	competitions Public					Q Notifications	itar 359	প্ট Fork	110
<> Code  • Issues 46	î Pull requests 4	Actions	미 Projects 1	🛱 Wiki	Security	🗠 Insights			
Project About	Codalab								

#### Isabelle Guyon edited this page on 9 Feb · 31 revisions

#### About CodaLab

Pages 78

CodaLab is an open-source platform that provides an ecosystem for conducting computational research in a more efficient, reproducible, and collaborative manner. There are two aspects of CodaLab: worksheets and competitions.

- Preferred open source competition platform by the community.
- More than thousand organized competitions



## Research interests

- Main interest in image, video, and multimodal data analysis for LAP
- Promoting explainability and interpretability for transparency
- Bias detection and mitigation for fairness, context, personalization
- Multimodal learning and with noise and asynchronous data
- Self-supervised learning to manage huge amount of data and reduce the need of annotated data
- Domain adaptation
- Uncertainty estimation and human in the loop
- LAP challenges design and implementation
- Real applications for good!

Dr. Sergio Escalera, <u>www.sergioescalera.com</u>, <u>sergio@maia.ub.es</u>

# **Video Transformers**

### Javier Selva Castelló

javierselva.github.io







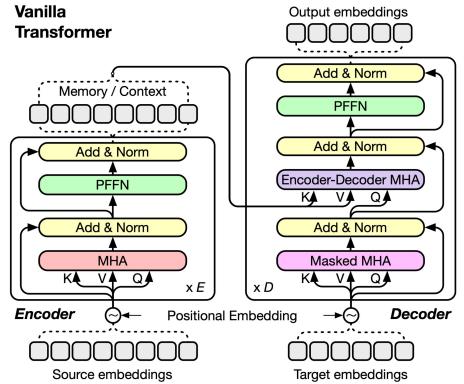
## **Outline**

- Transformers.
- Video pre-processing.
- Architectures for video.
- Training Video Transformers.
- Conclusions.

# Transformers

## **Transformer Overview**

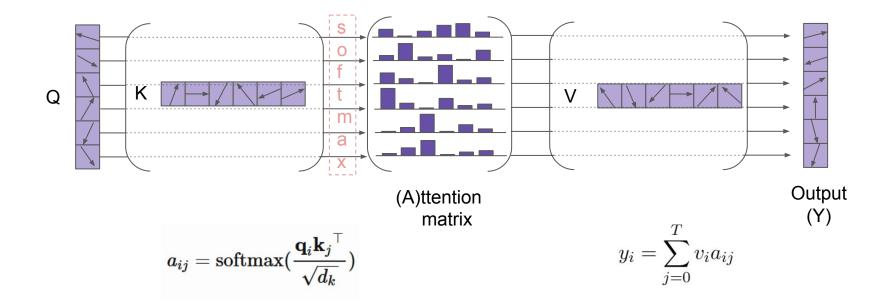
- Sequence modeling (tokens).
- Multi-head Attention.
  - Non-local interactions.
- Position-wise FeedForward.
  - Evolve token representations.



## **Non-local Token Mixing**

Queries, Keys and Values are linear transformations of the input.

 $\mathbf{Q} = \mathbf{X}\mathbf{W}_Q, \mathbf{K} = \mathbf{X}\mathbf{W}_{\mathrm{K}}, \mathbf{V} = \mathbf{X}\mathbf{W}_{\mathrm{V}}$ 



## Long-term Modeling for Video

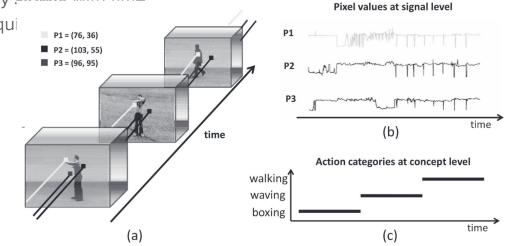
- Video is inherently a sequence.
- Opposed to RNNs or CNNs, allows to model long-term interactions in a single operation.
- Benefit motion modeling.





Video:

- Highly dimensional.
- Highly redundant:
  - Appearance-based semantics vary slowly with time
  - Capturing fine-grained motion requi



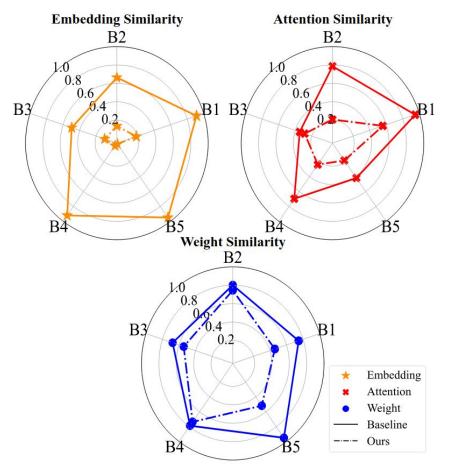
## **Challenges**

Video:

- Highly dimensional.
- Highly redundant.

Transformers:

- Token smoothness



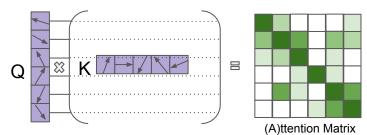


Video:

- Highly dimensional.
- Highly redundant.

Transformers:

- Token smoothness
- Quadratic complexity O(N<sup>2</sup>)



Pair-wise token affinity

For scale, 25 frames at 256 x 256 resolution, tokenized as 1 x 16 x 16 sized patches, results in ~**41M** elements in A.



Video:

- Highly dimensional.
- Highly redundant.

Transformers:

- Token smoothness
- Quadratic complexity **O(N<sup>2</sup>)**
- Lack of inductive biases
  - Versatile architecture.
  - Requires large datasets!

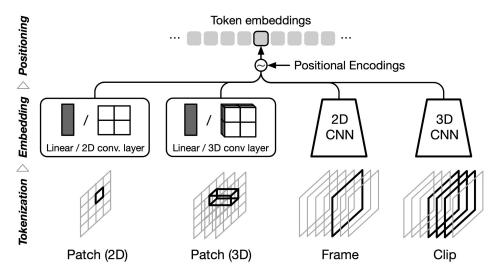
# Video pre-processing

## **Pre-processing**

**Tokenization**: Turned into a sequence of elements.

Embedding: Map raw pixel values into continuous vector representation.

Positional information: Transformers are agnostic to position.





Defines granularity at which interactions can be learned.

More token receptive field, less fine-grained feature modeling.

Defines **complexity** of the self-attention computation.

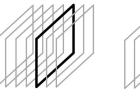
Smaller tokens = More tokens = More compute!



Patch (2D)



Patch (3D)



Frame



Clip

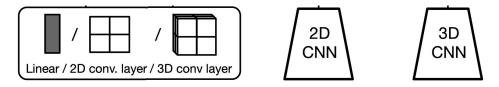


Minimal Embeddings: Single layer

Large Embedding Networks: Complete SOTA architectures

Dimensionality reduction.

CNNs will alleviate Transformer training.



Minimal Embeddings

Large Embedding Networks

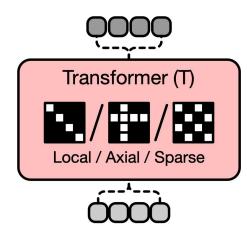
# Architectural designs

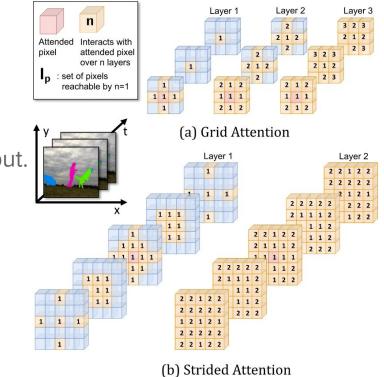
## **Restricted Approaches**

General techniques to reduce complexity!

Limit scope of attention operation.

Depend on **stacking** to account for complete input.







Leverage a reduced set of tokens.

May help reduce **redundancy**.

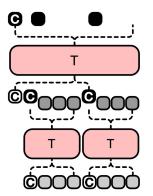
Needs to be done carefully:

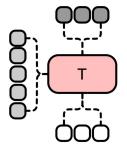
- May lose ability to model some dependencies later on!

Hierarchy

VS.

**Query-driven Compression** 

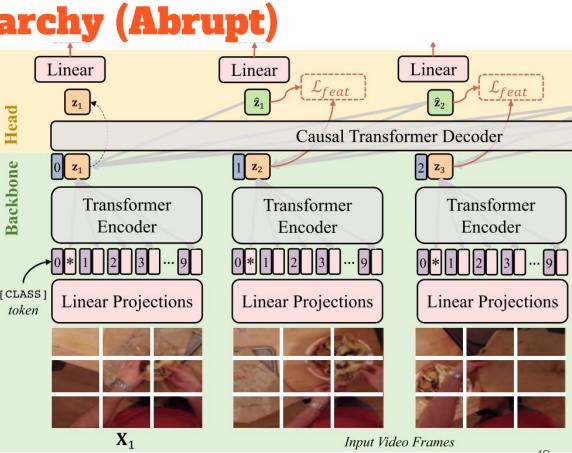




## **Aggregation: Hierarchy (Abrupt)**

<u>Abrupt</u> approaches:

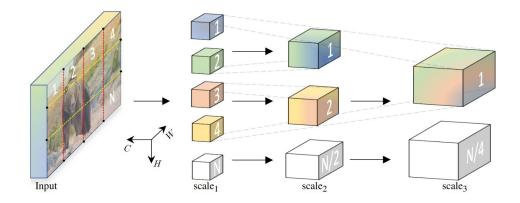
- Early aggregation: miss fine-grained motion
- Local neighborhoods
- CLS token



## **Aggregation: Hierarchy (Progressive)**

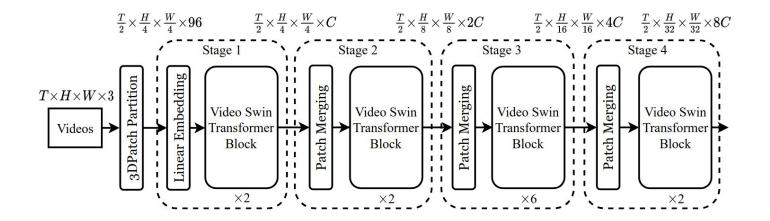
Progressive approaches:

- Gradual aggregation.
- Global interactions.
- Local learned **pooling.**



## **Aggregation: Hierarchy (Progressive)**

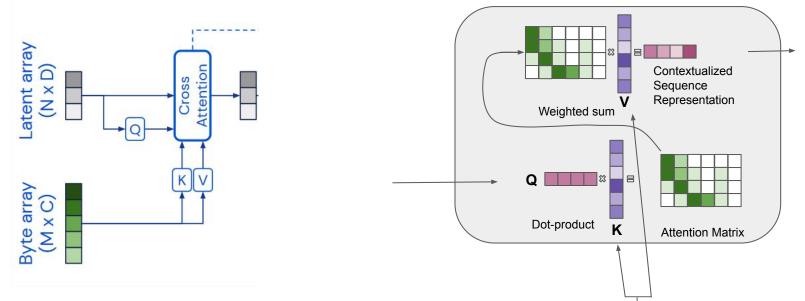
- Maintain temporal fidelity.
- Feature dimensionality increase with depth.
- Best performance/cost ratio!



## **Aggregation: Query-driven Compression**

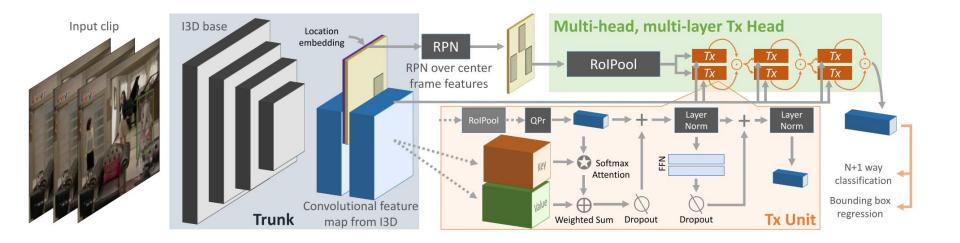
General technique, allows many different applications.

Compress full sequence into fewer tokens. Reduce complexity!



## **Aggregation: Query-driven Compression**

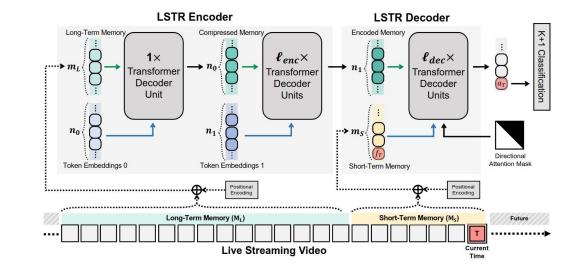
#### Data driven



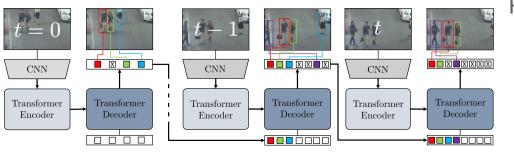
## **Long-term modeling**

#### Memory

- Stores several past observations.
- Can be accessed:
  - Through Cross-Attention.
  - By concatenation.
- Generally **compressed**:
  - At storing time.
  - Sparse sampling.
  - Query-driven Compression.
- **Discard** old elements:
  - FIFO.
  - Relevance (attention).



## **Long-term modeling**



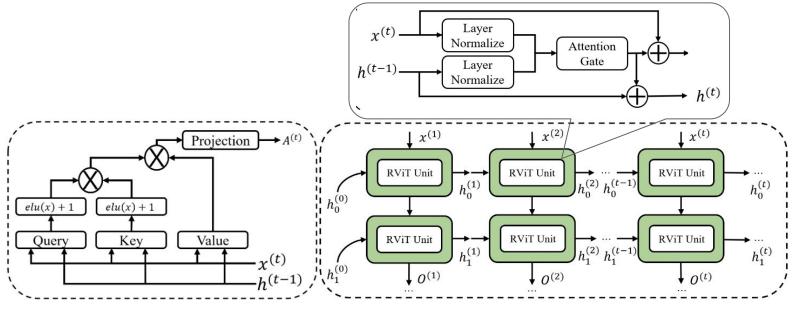


Images source: https://github.com/timmeinhardt/trackformer

#### Recurrence

- **Collapses** past observations in a single recurrent state.
- Send information to next time-step:
  - Within the Attention Operation.
  - Forwarding Output.
- Used when forwarding **local temporal information** is enough.
- Maintain access to **full resolution** at each new observation.

## **Long-term modeling: Recurrence**



Attention Gate (Aggregates **current input with recurrent state**) Traditional Recurrence Unrolling (Replaces **RNN Units** with a **Transformer**)

Yang, J. et al. "Recurring the Transformer for Video Action Recognition." In CVPR (2022).

# Training the Transformer

## **Training Regime**

Solving lack of **inductive biases**:

Large datasets.

Self-supervised Learning.

**Computational** limitations!

**End-to-end** training of embedding layers.

Large batches, videos, architecture...

Multiple training stages.

- Minimal Embedding Networks:
  - End-to-end
  - Require large-datasets or SSL.
- End-to-end with Large Emb. Net.:
  - Small Transformer.
  - **Efficient** designs.
  - CNN-Transformer tandem.
- Frozen embeddings:
  - Most common.
  - Transformer boost!
  - Limited by the pre-trained features.
- Image vs Video pre-training
  - Appearance variability
  - Motion modeling.

## Self-Supervised Learning (SSL)

Allows to leverage un-annotated data.

Provides data-specific biases.

Traditional Instance-based Learning.

- Spatial data augmentation.
- Align views through InfoNCE [1].
- Learn:
  - **invariance** to perturbations.
  - instance-based representations.

$$\mathcal{L} = -\sum_{i,k} \left[ \log \frac{\exp(\hat{z}_{i,k}^{\top} \cdot z_{i,k})}{\sum_{j,m} \exp(\hat{z}_{i,k}^{\top} \cdot z_{j,m})} \right]$$

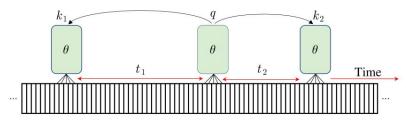
[1] Oord, A. et al. "Representation Learning with Contrastive Predictive Coding." In ArXiv (2018).

[2] Chen, T. et al. "A Simple Framework for Contrastive Learning of Visual Representations". In ICML (2020).

[3] Purushwalkam, S. et al. "Demystifying Contrastive Self-Supervised Learning: Invariances, Augmentations and Dataset Biases." In NeurIPS (2019).

[4] Feichtenhofer, C. et al. "A Large-Scale Study on Unsupervised Spatiotemporal Representation Learning." In CVPR (2021).





May learn invariance to temporal deformations!





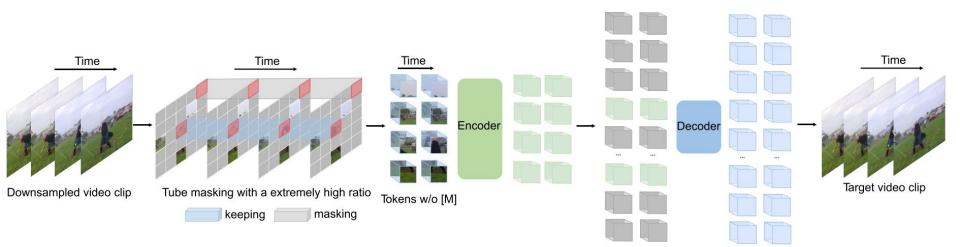
(h) Gaussian noise (i) Gaussian blur (j) Sobel filtering

## **SSL: Masked Token Modeling**

Force to leverage global appearance and motion semantics...

... to solve **local token-dependent** prediction tasks.

Crucial: High masking ratios, mask entire tubes!



## Conclusions

### **Conclusions**

Complexity:

Reduce appearance redundancy.

Careful! Need to keep fine-grained motion.

Progressive aggregation.

Memory models (variable length!).

MTM exploits Transformer abilities.

