

Context Recognition for the Application of Visual Privacy

ESR 14

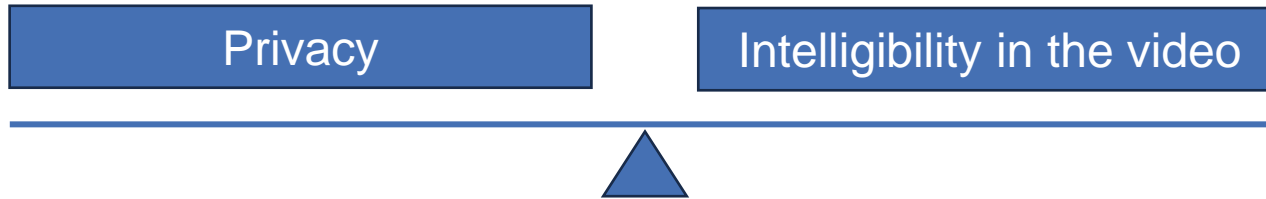
University of Alicante

20/11/2023

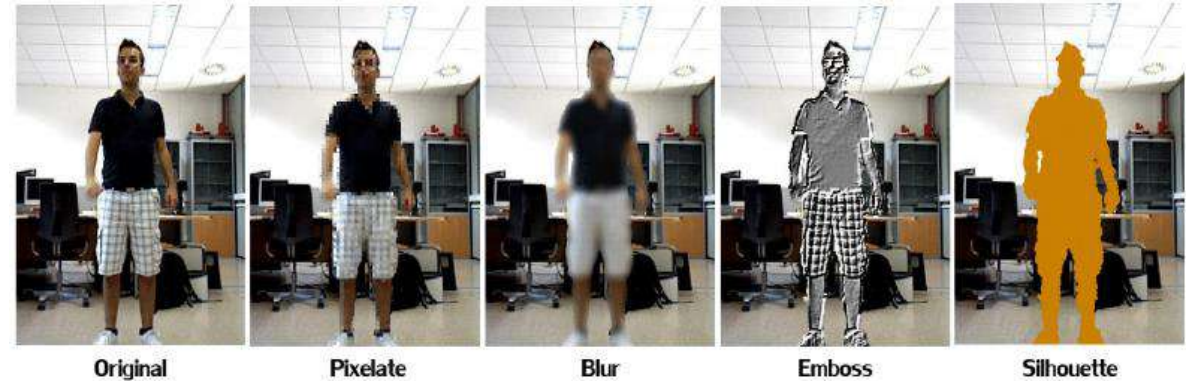
Kooshan Hashemifard

- Demographic changes
- Burden to care personnel and facilities
- Damage to autonomy, self-esteem and spirit
- Ambient-assisted living (AAL) and sensors
- Video-based technology
- The most directed and natural way to record events
 - Pros: Provide richer information
 - Cons: Easy to interpret by unauthorized viewers and privacy issues





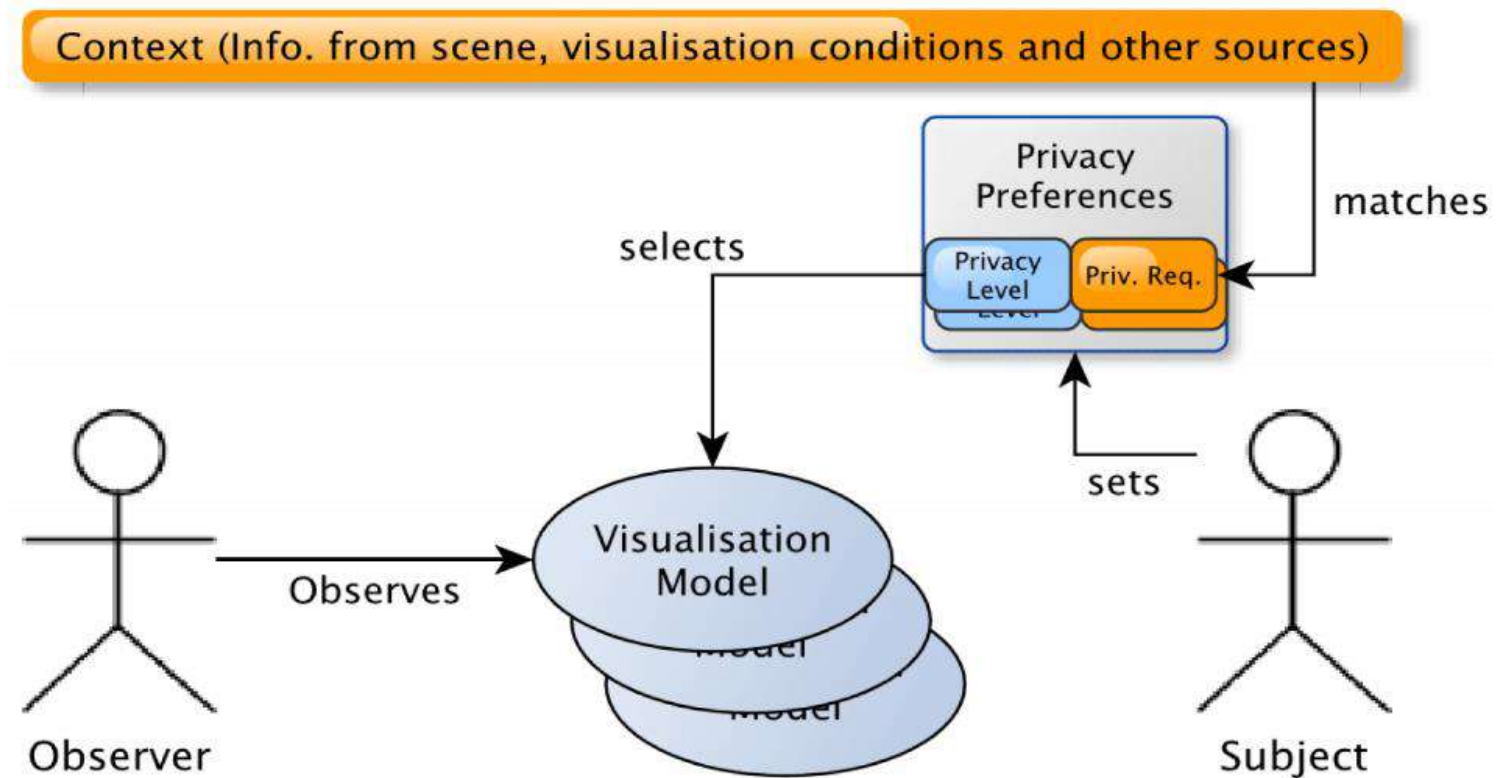
- Find a visualization method to understand what is happening and at the same time preserve privacy
- Previous work introduced privacy-by-context:
 - Level-based visualization
 - Selected according to the context



Objective: Automatic selection of visualization

Context is modeled by the following variables:

- Activity
- Incident
- Appearance
- Observer
- Place, ...



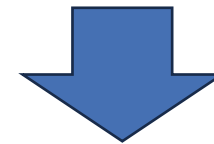
- This leads to the following research questions, in line with the objective:
 - - Can we automatically estimate the context that may affect visualization preference?
 - - Can we estimate relevant activities?
 - - Can we estimate relevant events, such a fall?
 - - Can we estimate degree of nudity?
 - - Can we integrate them into a privacy-by-context approach so that visual privacy adapts in real time?



① **Activity**



② **Incident**



③ **Appearance**

1

Activity

Human Activity Recognition (HAR)

Action recognition

involves identifying and classifying specific actions or movements.

- Examples: walking, jumping, specific gestures like waving or handshaking
- Applications: human-computer interaction, sports analysis, and robotics

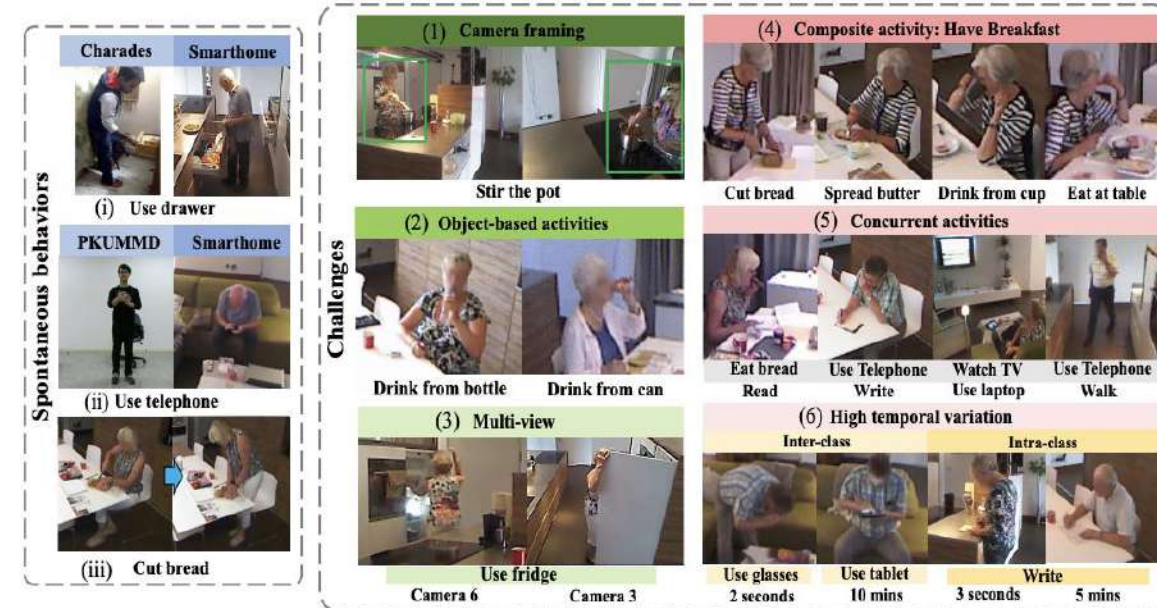
Activity recognition

is a broader concept that involves identifying and understanding a sequence of actions or interactions over a certain period to infer the overall activity or behavior.

- Examples: cooking, playing basketball, working at a desk
- Applications: healthcare monitoring, smart homes, and context-aware computing

Toyota Smarthome Dataset

- The subjects are senior people in the age range 60-80 years old.
 - 35 Activities from daily living:
 - Composite Activities: cooking, cleaning, making breakfast
 - Elementary Activities: laydown, watch tv, use laptop, reading, phone call, take pill
 - Object-based Activities: drink from bottle/can
 - Data Modalities available: RGB, depth, skeleton
 - Data Modalities use:
 - RGB for scene details, depth maps for 3D structural information, skeletal data for 3D joint locations.
- Multi-modal approaches integrate different modalities for a comprehensive understanding of complex actions.

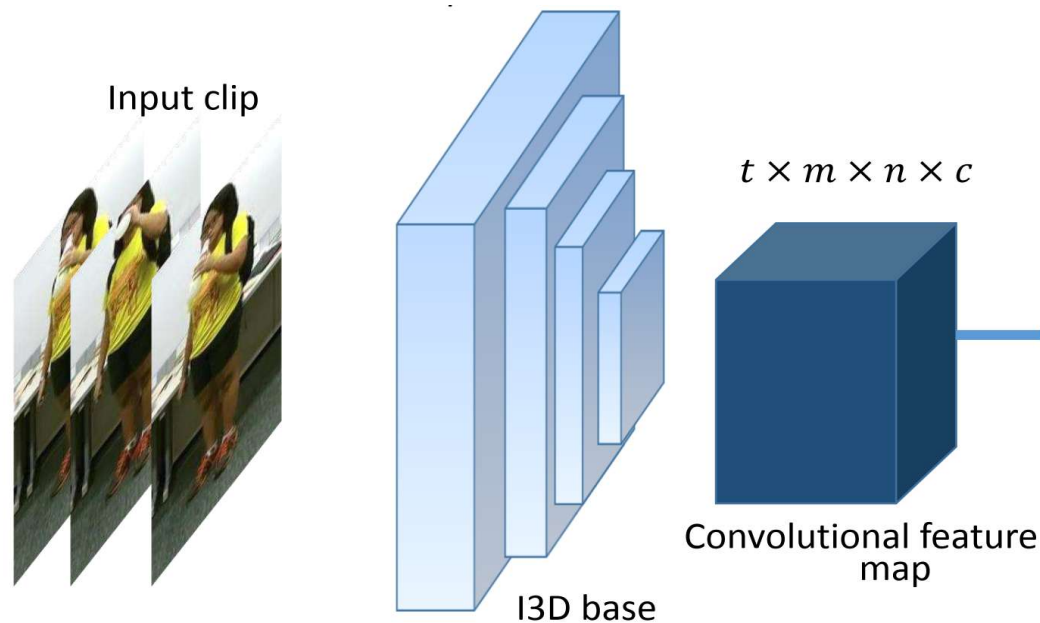


- **Spatial Information Encoding:**
 - Utilization of CNN. CNNs' ability to extract useful and discriminative features
 - **Temporal Information Encoding:**
 - Existing deep architectures encode temporal information with limited solutions.
 - Challenges in acquiring both local and global variations of temporal features.
 - RNNs and time series models, higher dimension CNNs, newer Transformer architecture.
-
- **Promising Solution: Transformers in HAR**
 - New encoder-decoder architecture using attention mechanism.
 - Success in natural language processing tasks; now applied to images and video recognition.
 - Video as a sequence of images, akin to language processing (image frames as words).
 - Not restricted to sequential processing; attention mechanism provides context for any position.
 - Relatively new approach, Increasing research focus on transformers for action recognition in recent years.

1- Choosing Backbone

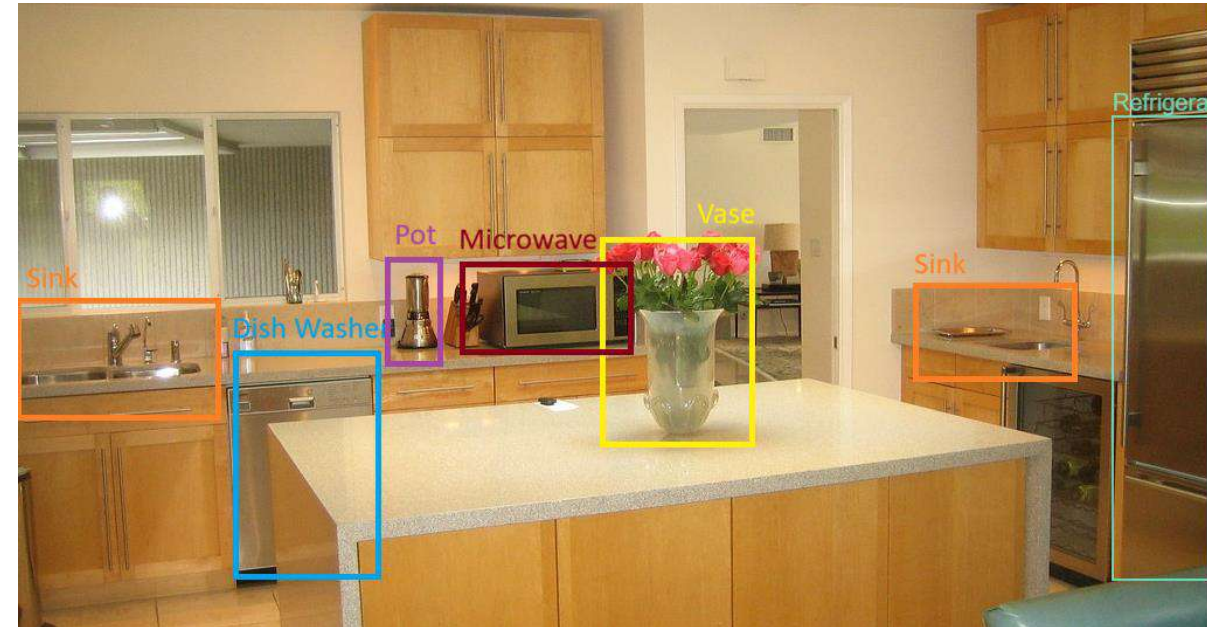
CNN-based Network for rich feature extraction

- 3D CNN Networks: I3D Net
- Time-distributed 2D CNN Networks: EfficientNet, DenseNet, InceptionNet, ...



2- Reinforcing Feature extraction with Object Detection

- Auxiliary information about Daily Activities
 - House Objects can play an import role for a given activity
 - Location information
- Advances in Object Detection
- House objects datasets

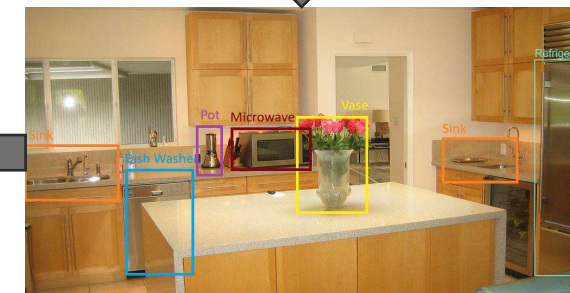


Object Groups of Interests

- **Person**
- **Kitchen Furniture:** Stove, Refrigerator, Oven, Microwave, Sink
- **Living room Furniture:** TV, Sofa, Table
- **Cutlery:** Spoon, Knife, Dish, Glass
- **Food**
- **Bathroom Furniture:** Toilet, Shower, Tube
- **Electronics:** Laptop, Cell phone, Tablet
- **Bedroom Furniture:** Bed, Wardrobe

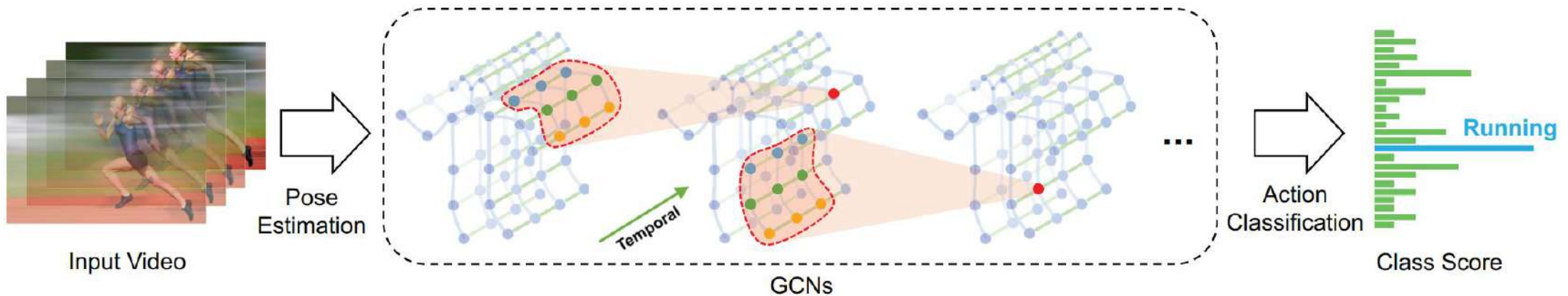


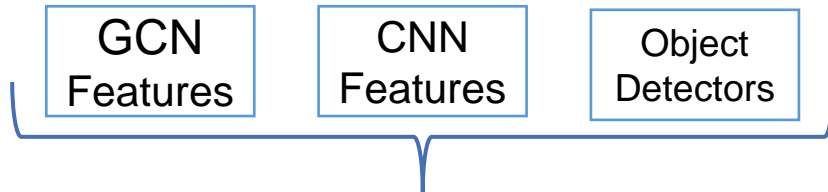
Object Detector
YOLO v8



3- Pose Information

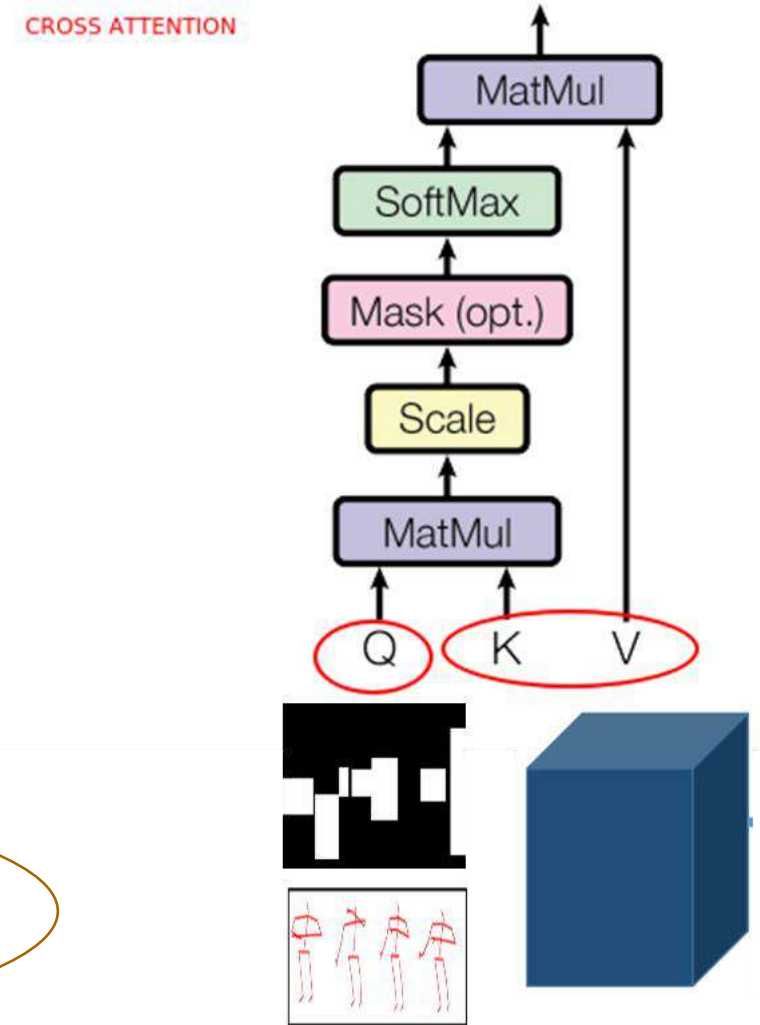
- The RGB image includes rich scene and object information, but does not encode long-range temporal information
- Pose info is more robust to changes in appearance, clothing, or lighting condition
- View invariance: Relationship between body joints coordinates can tackle Cross-view (CV) problem.
- Pose as Graph: joints as Vertexes and bones as Edges
- Approach: GCN and Body Joints Graph



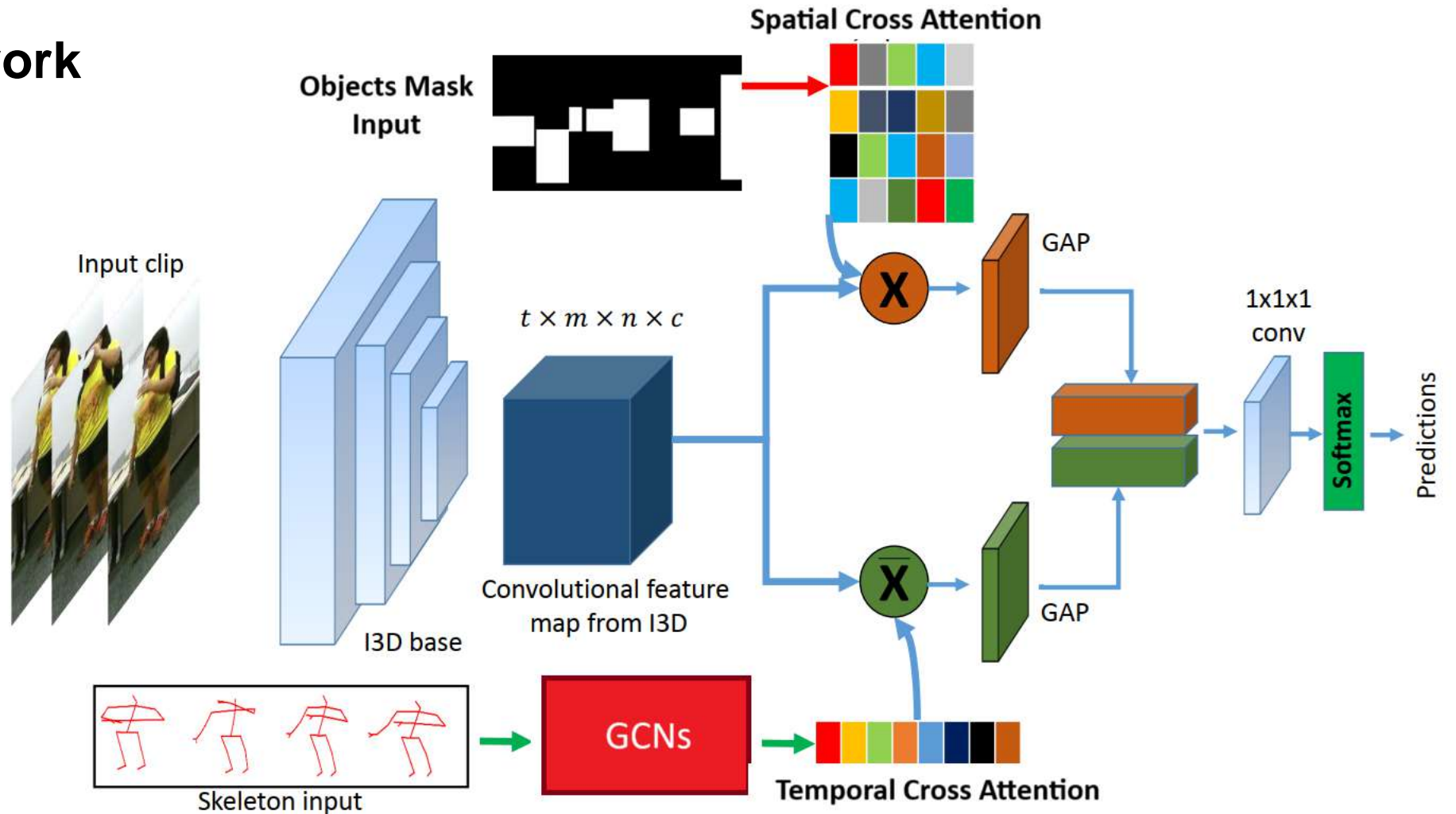


Multi-Source Information Integration Methods:

- **Feature Concatenation:** Combining features by concatenating them along a specific dimension.
- **Feature Fusion:** Integrating features through a fusion process, which can involve mathematical operations like addition, averaging, etc.
- **Cross Attention:** Focusing on relevant parts of input data through attention mechanisms, allowing the model to emphasize important information during merging.



Proposed Network



2

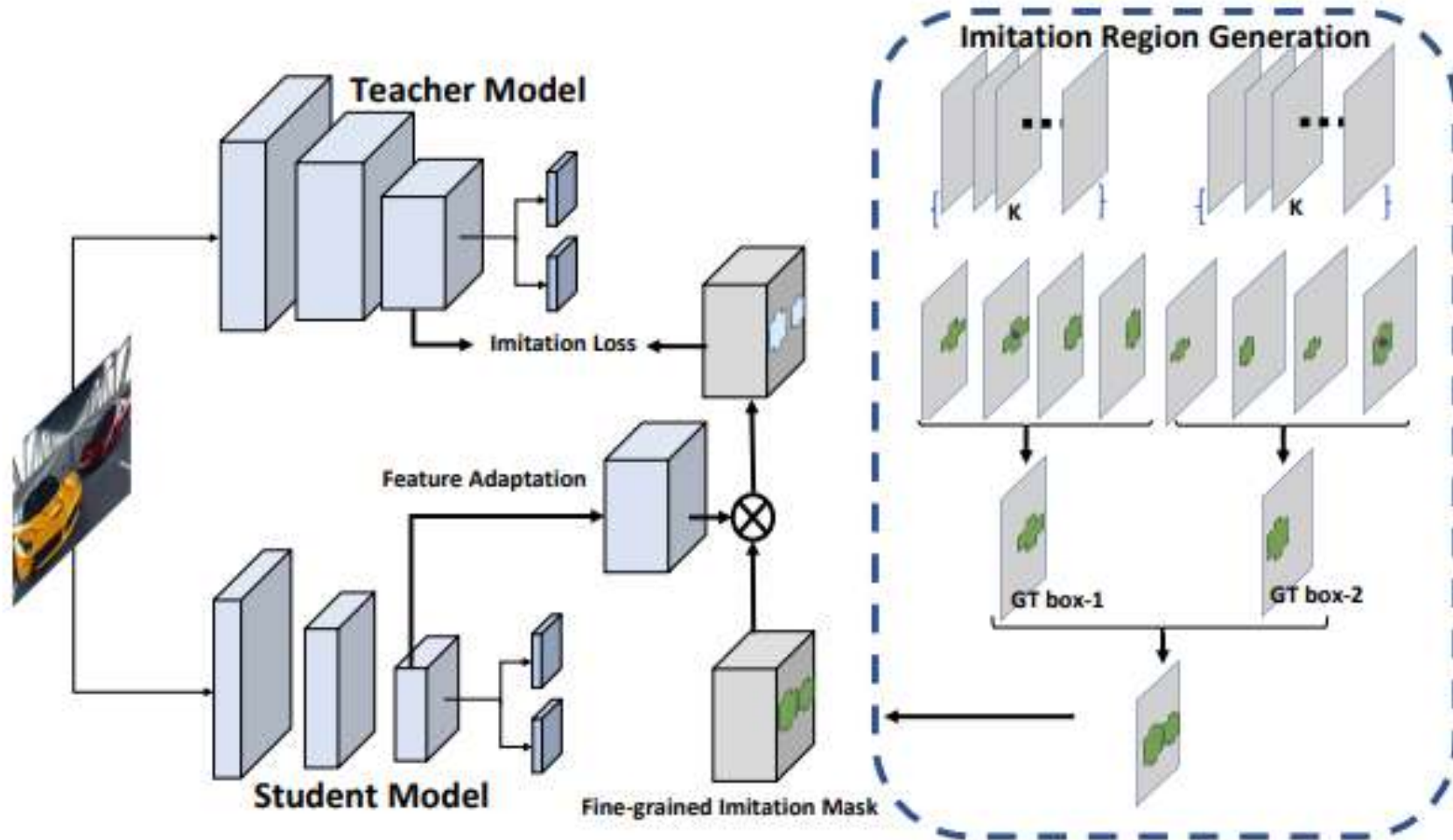
Incident

- End-to-end **Object Detection** approach with 2 classes: fallen and upright individuals
- Top layer of the object detector network was replaced with two output neurons for each status
- Pre-trained parameters retained, except for the last layer
- Fine-tune on E-FPDS dataset, contains 6982 images captured in indoor environments, with 5023 instances of falls and 2275 instances of non-falls in various scenarios, including variations in pose, size, occlusions, and lighting



- state-of-the-art CNN-based networks can be computationally expensive and difficult to deploy on smaller devices, especially in real-time and multi-tasking scenarios
- knowledge distillation has emerged to directly learn compact models by transferring knowledge from a large model (teacher network) to a smaller one (student network)
- Reducing computational costs without sacrificing validity
- Fine-grained Feature Imitation method
 - Local features in the object region and near its anchor location contain important information and are more crucial for the detector and how the teacher model tends to generalize
 - The smaller student detector is trained by using both ground truth supervision and feature response imitation on object anchor locations from the teacher networks

Knowledge Distillation



Evaluation In-the-wild with Luxonis

19

- model's performance in real-world settings after conversions
- videos, captured using the Luxonis OAK-D camera
- 38 videos recorded from various angles and labeled



③ Appearance

Appearance Detection in Literature

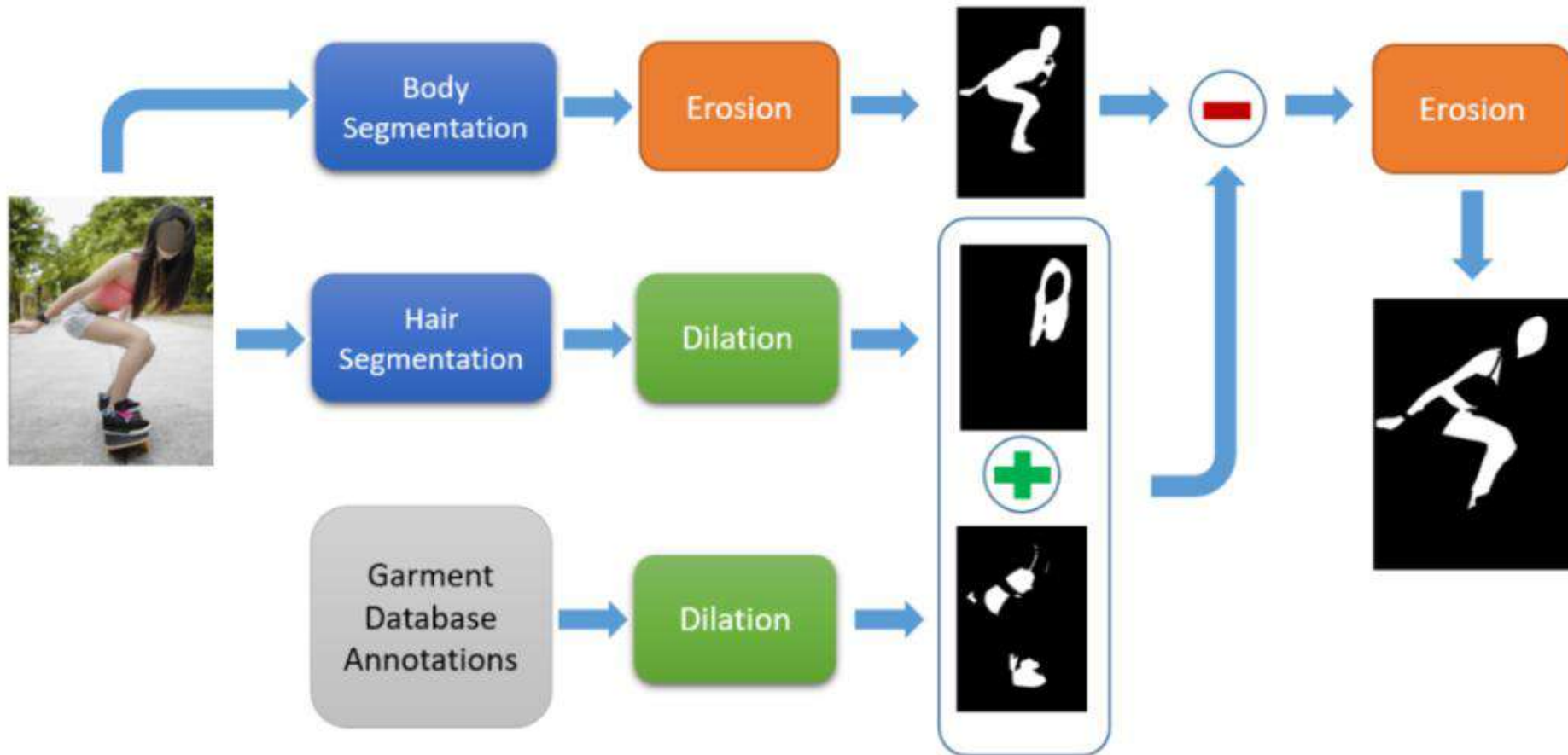
- Obscene image Classification
 - Naïve approach
 - No level-based classification
 - Subjective
 - Fail to generalize the problem
 - Limited standard datasets
- Garment Detection
 - Detect different clothing and wearables
 - Market analysis and finding trends
 - Does not necessarily represent the coverage level
 - Cannot cover all types of garments



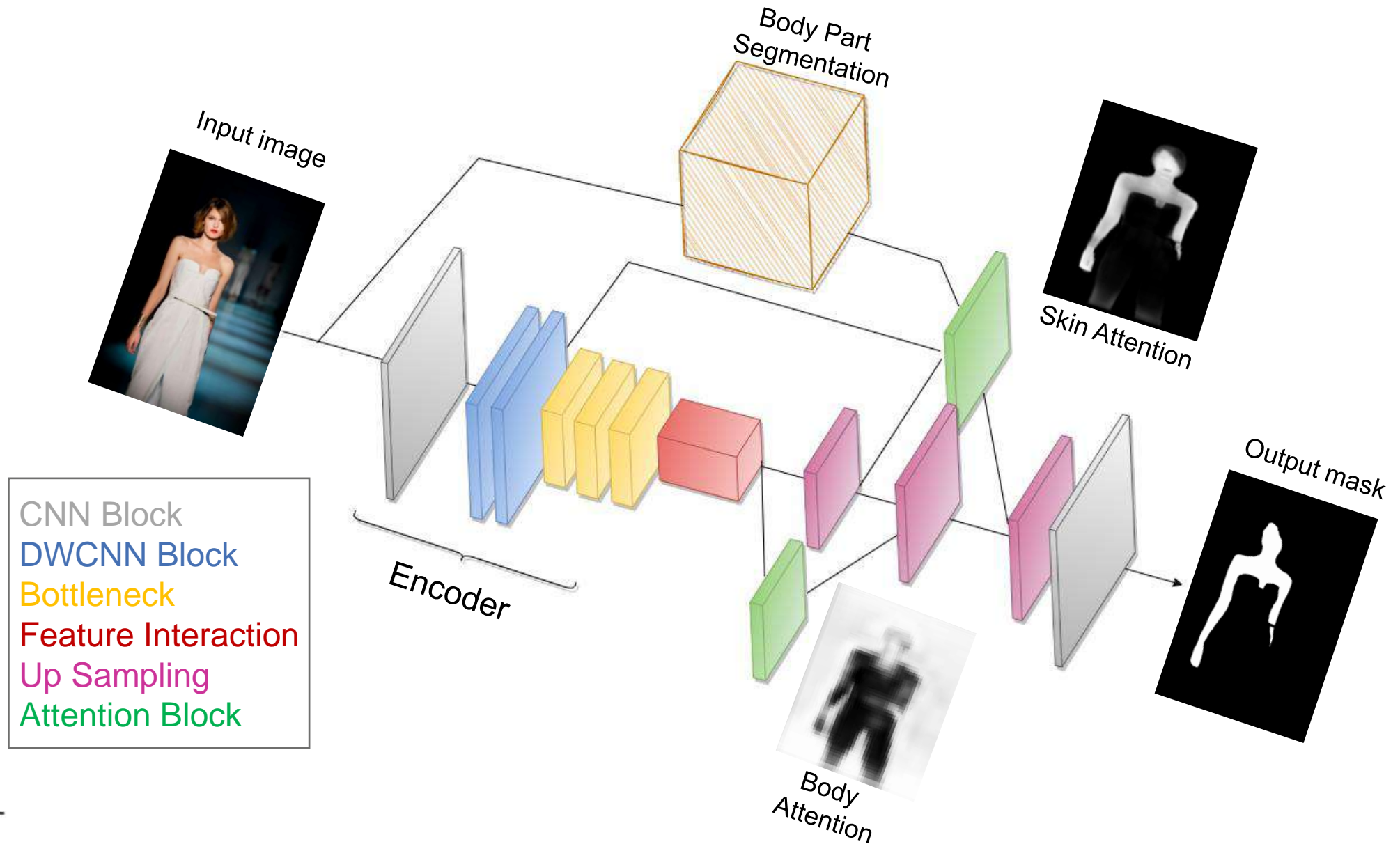
Nudity and AAL

- Serious concerns about person appearance recorded in a video
- Appearance detection can be interpreted as nudity detection in the context of AAL
- Skin Exposure Detection plays a crucial role in nudity detection
- Estimate the degree of nudity

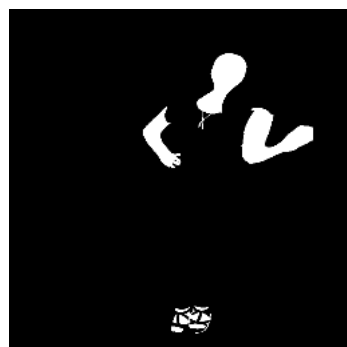
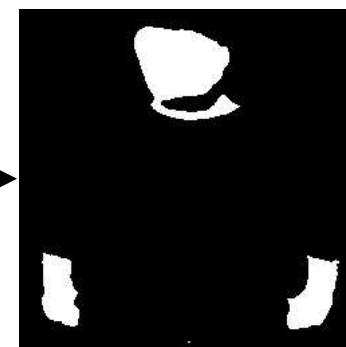
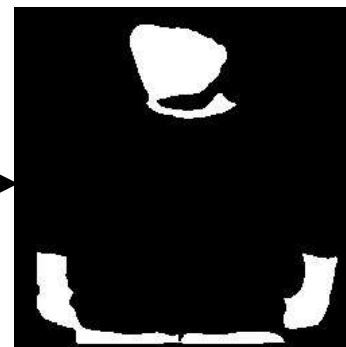
1- Dataset



2- Architecture

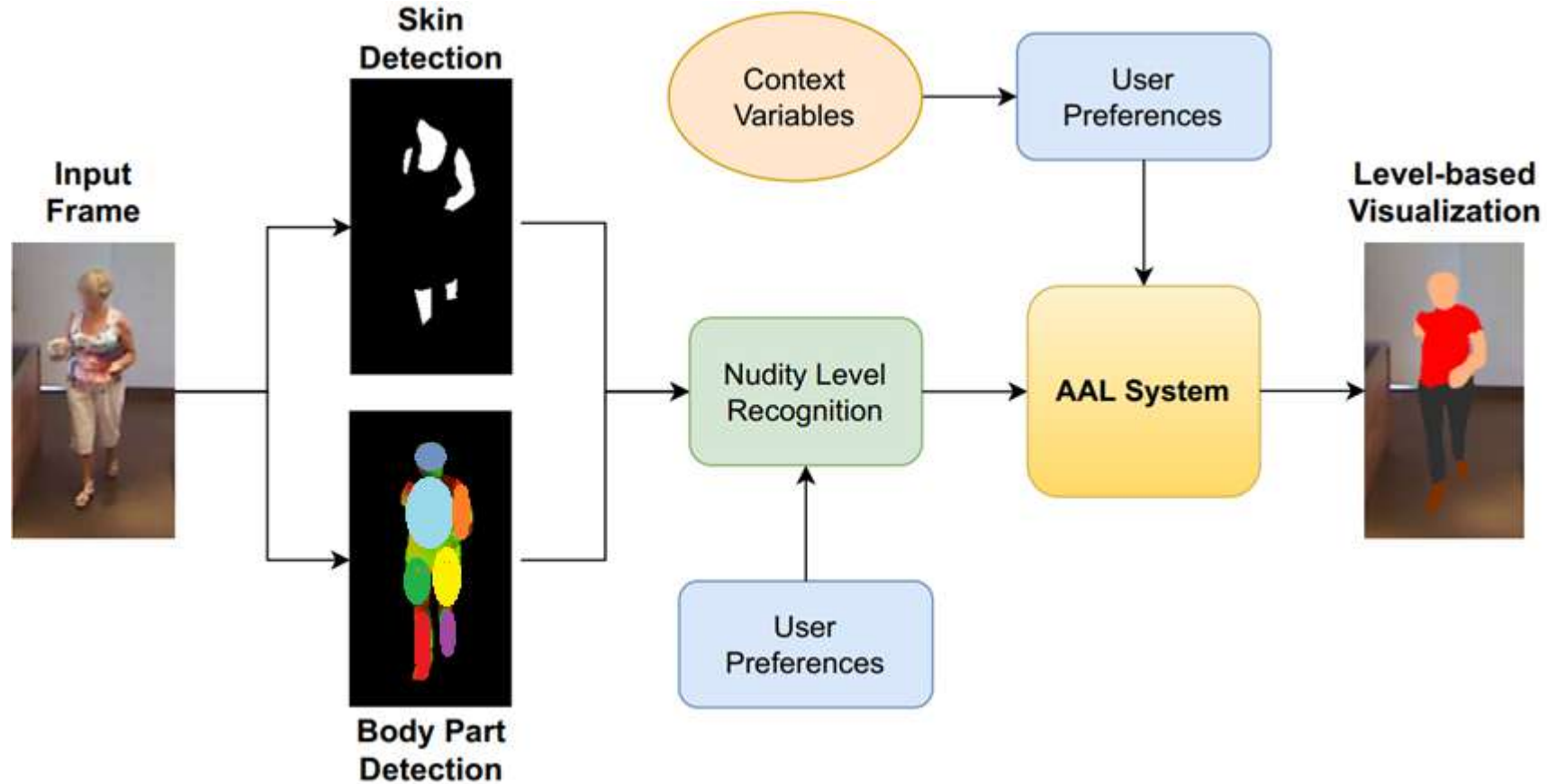


3- Training

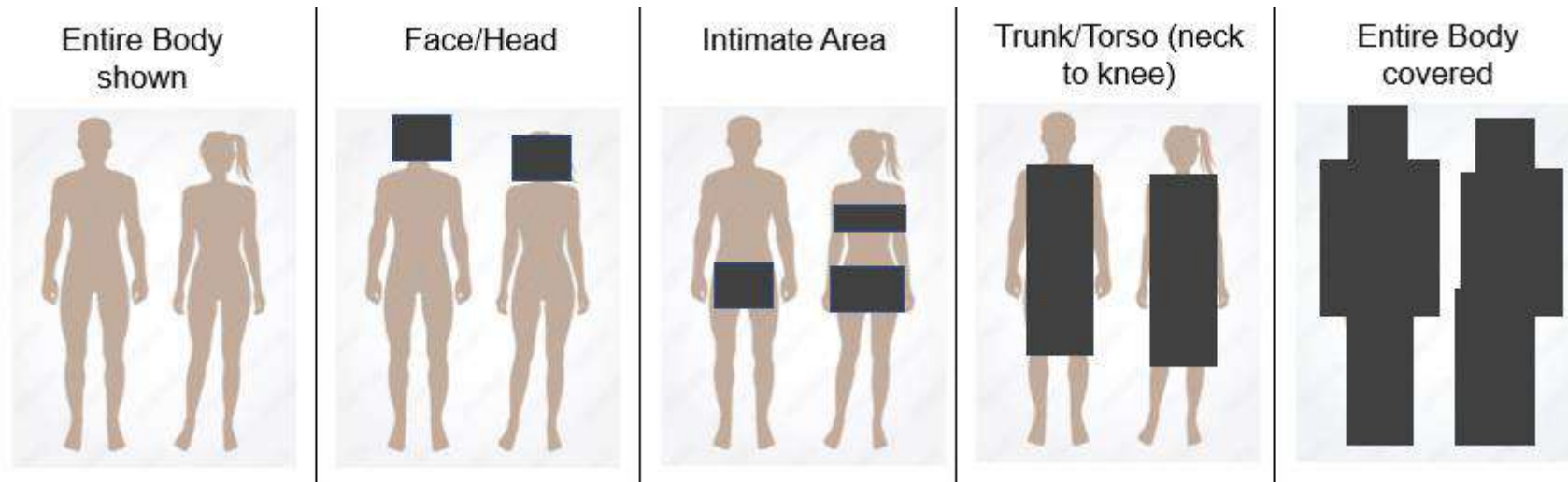


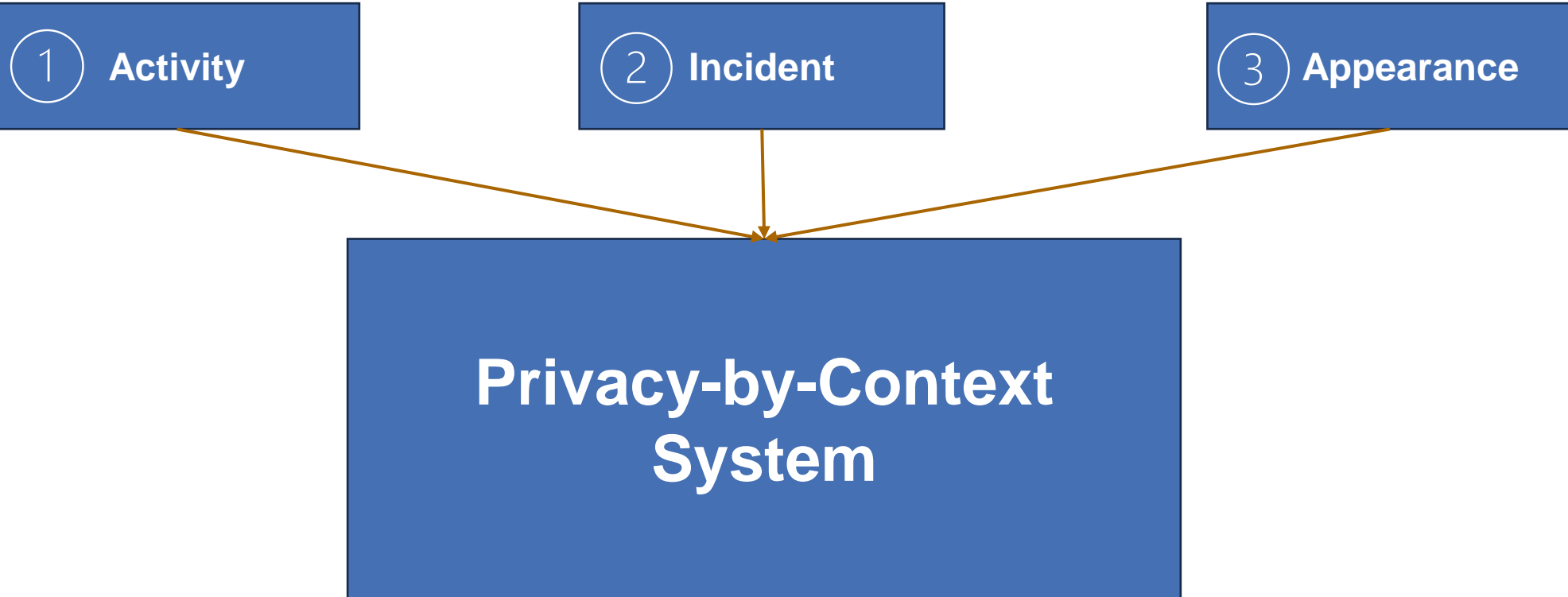
4- Nudity Recognition

Proposed Appearance Detection in AAL



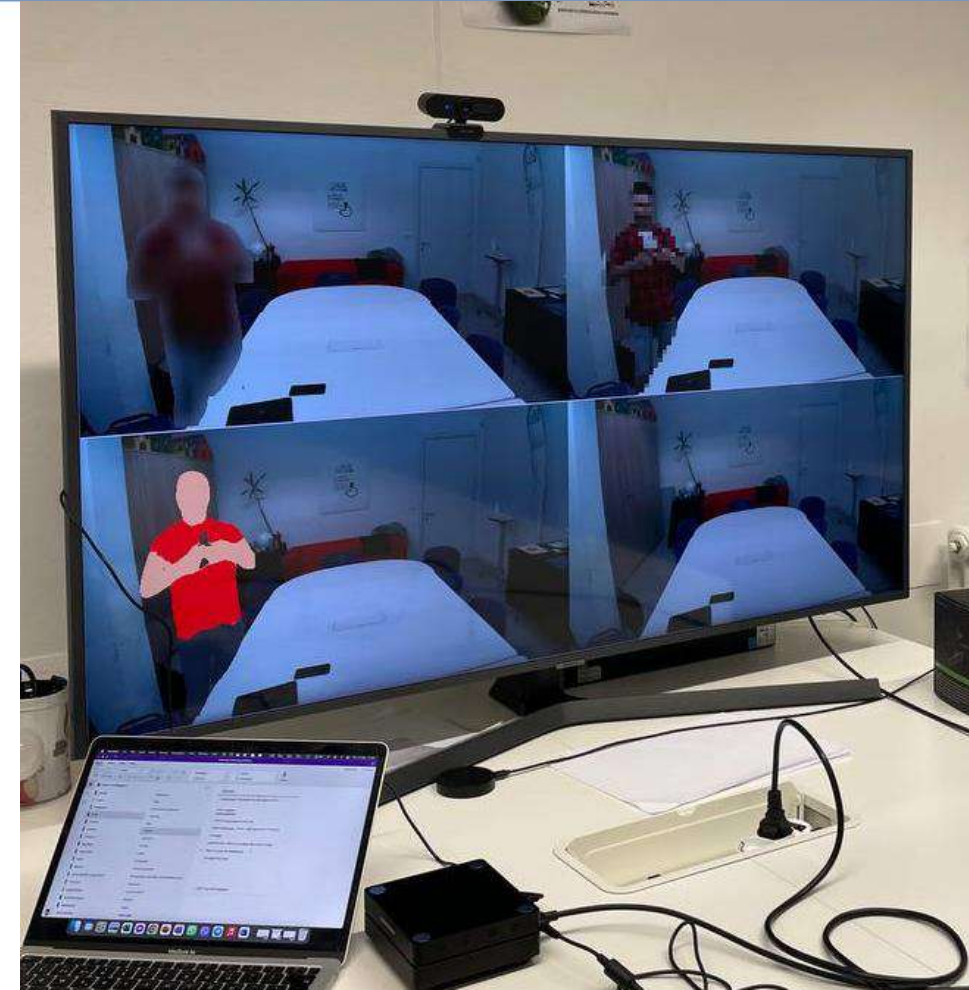
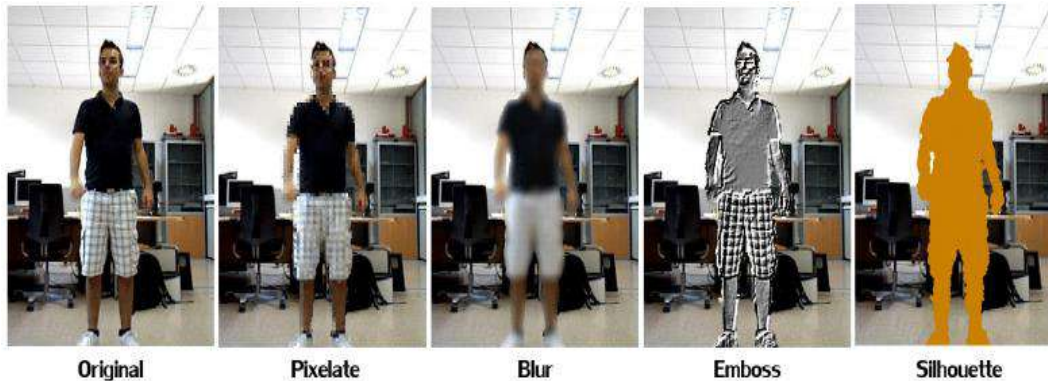
Nudity Level	Description
1	<i>Completely covered</i>
2	<i>Covered torso (neck to knee)</i>
3	<i>Covered intimate areas</i>
4	<i>Covered faces</i>
5	<i>Full body or exposed intimate areas</i>





Visualization Demo:

- Real-time illustration of different visualization
- Deployed on Jetson Xavier NX



A Step Closer to Privacy by Context System

30

- Activity Detection ✓
- Incident Detection ✓
- Appearance Detection ✓
- Visualization Filters ✓





Remaining work:

1. Finalizing research paper(s) in HAR
2. Applying the method to the top view ODIN dataset
3. Finalizing Appearance Recognition in practice
4. Proposing the whole Privacy by Context monitoring system
5. Finalizing Thesis

Pursuing Research Position

- Post Doc position in Academia
- Research Engineer in Industry

Fields of Interest

Computer vision

Implementation on Hardware and
devices

Working with Robots

Work closer to product

- **Journal Paper**

- K. Hashemifard, P. Climent-Perez, and F. Florez-Revuelta, “Weakly supervised human skin segmentation using guidance attention mechanisms,” *Multimedia Tools and Applications*, 2023.

- **Conference Papers**

- K. Hashemifard, F. Flórez-Revuelta, and G. Lacey, “A fallen person detector with a privacy-preserving edge-ai camera,” 2023.
- S. Ravi, P. Climent-Perez, T. Morales, C. Huesca-Spairani, K. Hashemifard, and F. Flórez-Revuelta, “Odin: An omnidirectional indoor dataset capturing activities of daily living from multiple synchronized modalities,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 6487–6496.
- K. Hashemifard and F. Florez-Revuelta, “From garment to skin: The visual skin segmentation dataset,” in *International Conference on Image Analysis and Processing*, Springer, 2022, pp. 59–70.
- C. Maidhof, K. Hashemifard, J. Offermann, M. Ziefle, and F. Florez-Revuelta, “Underneath your clothes: A social and technological perspective on nudity in the context of aal technology,” in *Proceedings of the 15th International Conference on Pervasive Technologies Related to Assistive Environments*, 2022, pp. 439–445.

Thank you!

Kooshan Hashemifard

University of Alicante

k.hashemifard@ua.es

Evaluation Protocols

- **Cross-Subject (CS):** same camera view, performed by different person
- **Cross-View (CV):** different camera view, performed by the same person

