# Learning to Sense Robustly and Act Effectively

**Benjamin Kuipers, P.I. and Silvio Savarese, co-P.I.**
**EECS, University of Michigan, Ann Arbor, Michigan**

## Introduction

The physical environment of a cyber-physical system is unboundedly complex, changing continuously in time and space. An embodied cyber-physical system, embedded in the physical world, receives a high bandwidth stream of sensory information, and sends continuous control signals. Traditional embedded systems restrict the environment or the attributes considered relevant, and depend on human supervision.
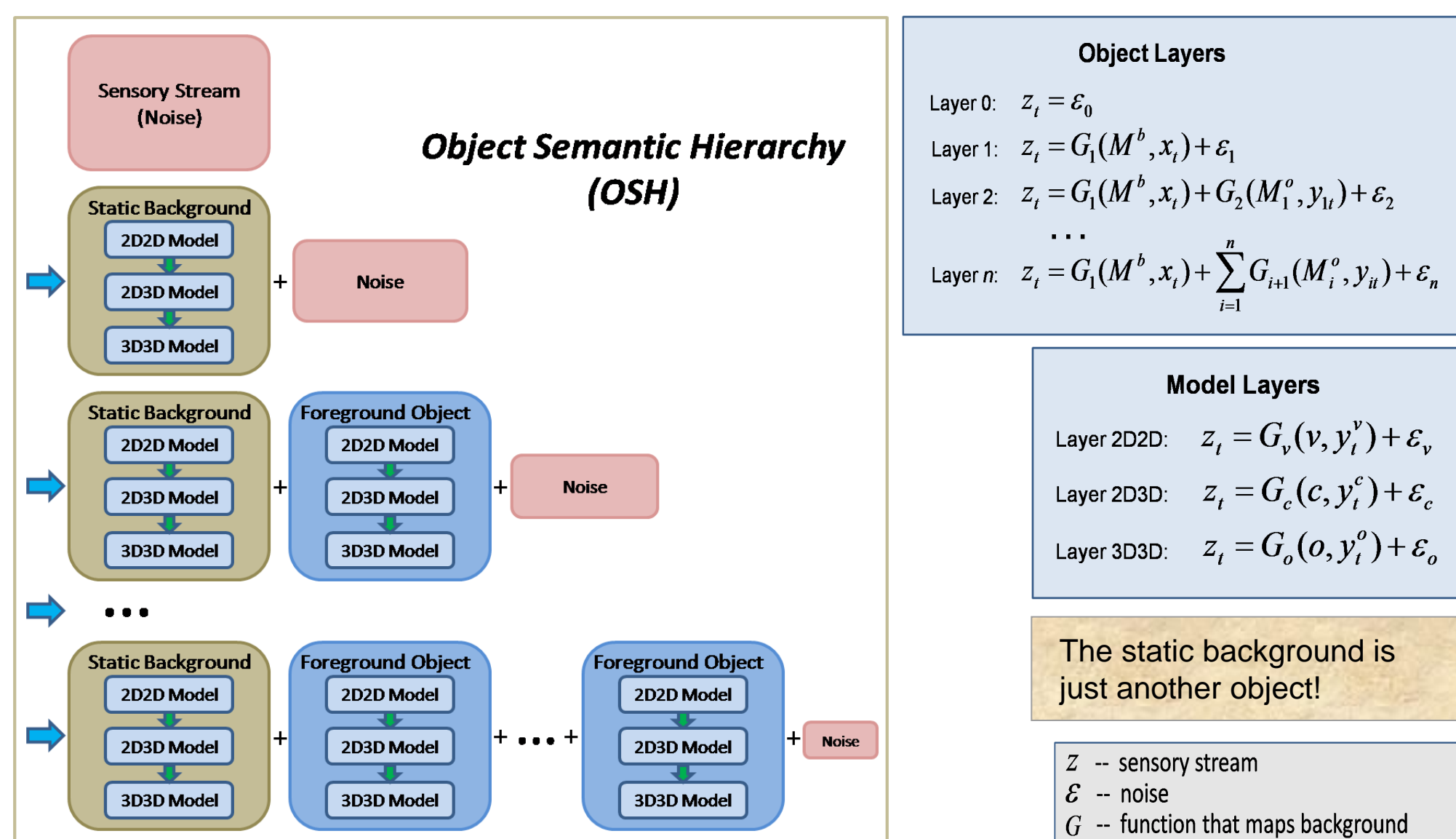
To handle the complexity of unrestricted environments, future cyber-physical systems will need to be learning agents, learning the properties of sensors, effectors, and environment from their own experience, and adapting over time. Foundational concepts such as **Space**, **Object**, **Action**, etc., will be essential for abstracting and controlling the complexity of its world.

Our previous work on the **Spatial Semantic Hierarchy (SSH)** [Kuipers, AIJ, 2000; Beeson, et al, IJRR, 2010] shows how multiple representations of space can bridge the gap between continuous interaction with the physical environment, and discrete symbolic descriptions that support effective planning.

We are developing robot agents that use vision and manipulation to learn models of objects and actions at multiple levels of representation:
 (1) learning to perceive objects in its environment;
 (2) joint optimization of semantic constraints in vision;
 (3) learning a hierarchy of increasingly skilled actions.

The **Object Semantic Hierarchy (OSH)** [Xu & Kuipers, ICDL, 2010] shows how a learning agent can create a hierarchy of representations for visual perception of objects it interacts with. The OSH "object abstraction" factors uncertainty in the sensor stream into object models and object trajectories.



**Object Semantic Hierarchy (OSH)**

The static background is just another object!

**Object Layers**

Layer 0: $z_i = \varepsilon_0$
Layer 1: $z_i = G_1(M^b, x_i) + \varepsilon_1$
Layer 2: $z_i = G_1(M^b, x_i) + G_2(M_1^o, y_{i1}) + \varepsilon_2$
...
Layer n: $z_i = G_1(M^b, x_i) + \sum_{i=1}^{n} G_{i+1}(M_i^o, y_{i1}) + \varepsilon_n$

**Model Layers**

Layer 2D2D: $z_i = G_v(v, y_i^r) + \varepsilon_v$
Layer 2D3D: $z_i = G_c(c, y_i^r) + \varepsilon_c$
Layer 3D3D: $z_i = G_o(o, y_i^r) + \varepsilon_o$

$Z$ -- sensory stream
$\mathcal{E}$ -- noise
$G$ -- function that maps background and/or object models to a 2D image given background and/or object poses
$M^b$ -- static background model
$X$ -- robot pose
$M^o$ -- object model
$\mathcal{Y}$ -- object pose
$V$ -- 2D view model
$C$ -- 2D component model
$O$ -- 3D object model

The uncertainty in the agent's sensory stream is factored into a collection of relatively compact representations:
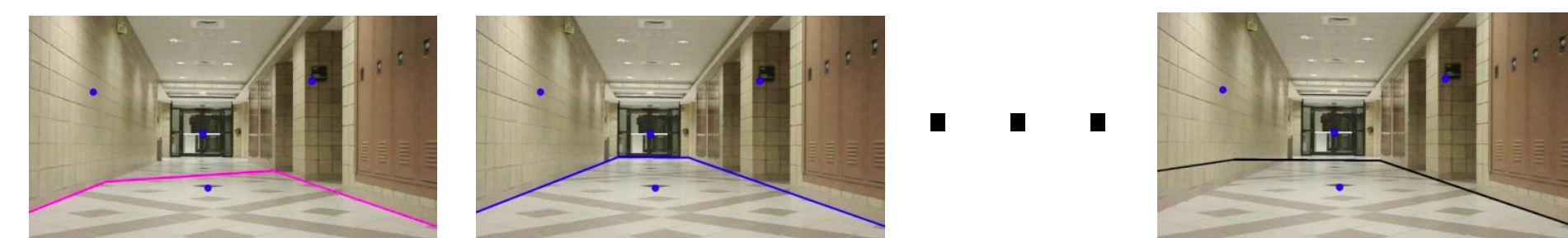 ➤ static background model
 ➤ pose trajectory of the agent
 ➤ constant foreground object models
 ➤ pose trajectories of foreground objects
 ➤ any remaining noise
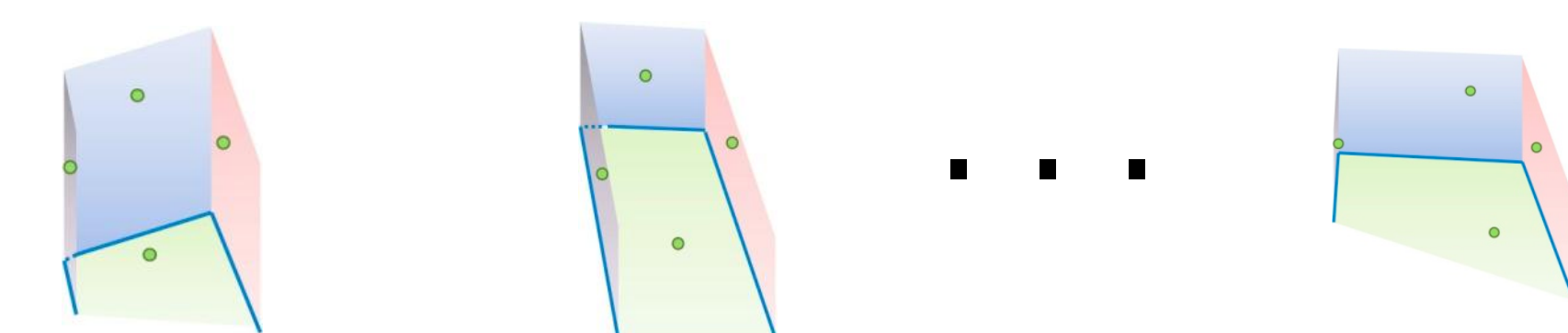
## Indoor Scene Understanding

Building on the OSH, and treating the surrounding environment as an "object", Tsai, Xu, Liu & Kuipers [ICCV, 2011] present a new method whereby an embodied agent using visual perception can efficiently create a model of a local indoor environment from its experience moving within it.

Our method uses a single-image analysis, not to attempt to identify a single accurate model, but to propose a set of plausible hypotheses about the structure of the environment from an initial frame. We then use data from subsequent frames to update a Bayesian posterior probability distribution over the set of hypotheses. The likelihood function is efficiently computable by comparing the predicted location of point features on the environment model to their actual tracked locations in the image stream.
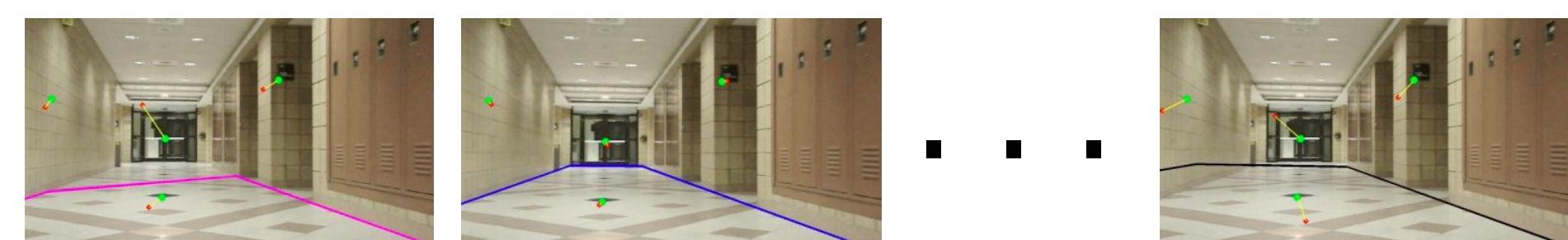
### Generate hypotheses



### Reconstruct 3D planar model



### Estimate camera pose



### Bayesian filtering

$$p\left(H_i | \mathbf{O}_1, \mathbf{O}_2, ..., \mathbf{O}_m\right) \propto p\left(H_i\right) \prod_{j=1...m} p\left(\mathbf{O}_j | H_i\right)$$

$$p(\mathbf{O}_j | H_i) \propto \prod_{o_k \in \mathbf{O}_j} \exp\left(\frac{-||\hat{\mathbf{L}}(o_k) - \mathbf{L}(o_k)||^2}{2\sigma^2}\right)$$

Confidence

$Hypothesis_1$
$Hypothesis_2$
...
$Hypothesis_n$

Time

Our method runs in real time, and avoids the need for extensive prior training and the Manhattan-world assumption, which makes it more practical and efficient for an intelligent robot to understand its surroundings compared to most previous scene understanding methods. Experimental results on a collection of indoor videos suggest that our method is capable of an unprecedented combination of accuracy and efficiency.
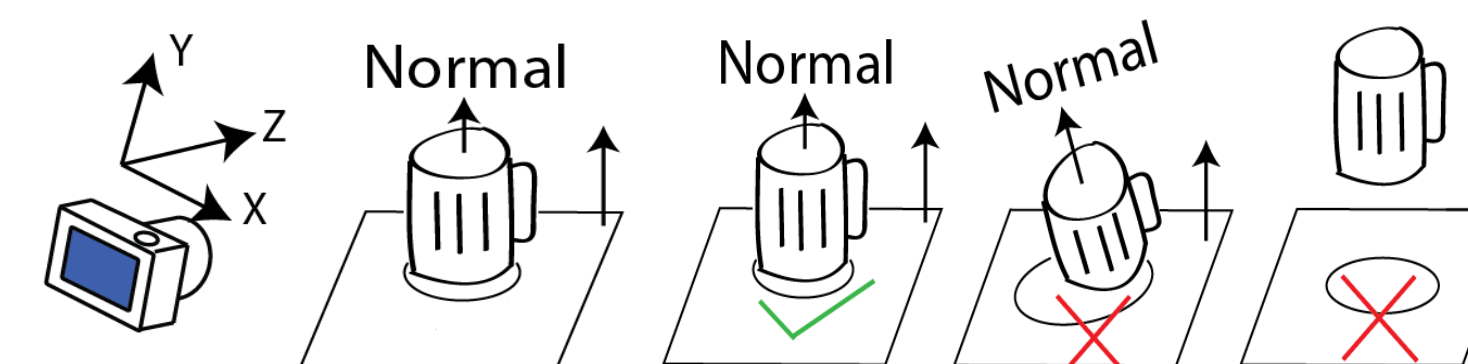
## Semantic Constraints in Vision

• **Overview**
  • Single un-calibrated image
  • Improve object detection's accuracy
  • Estimate camera pose and focal length
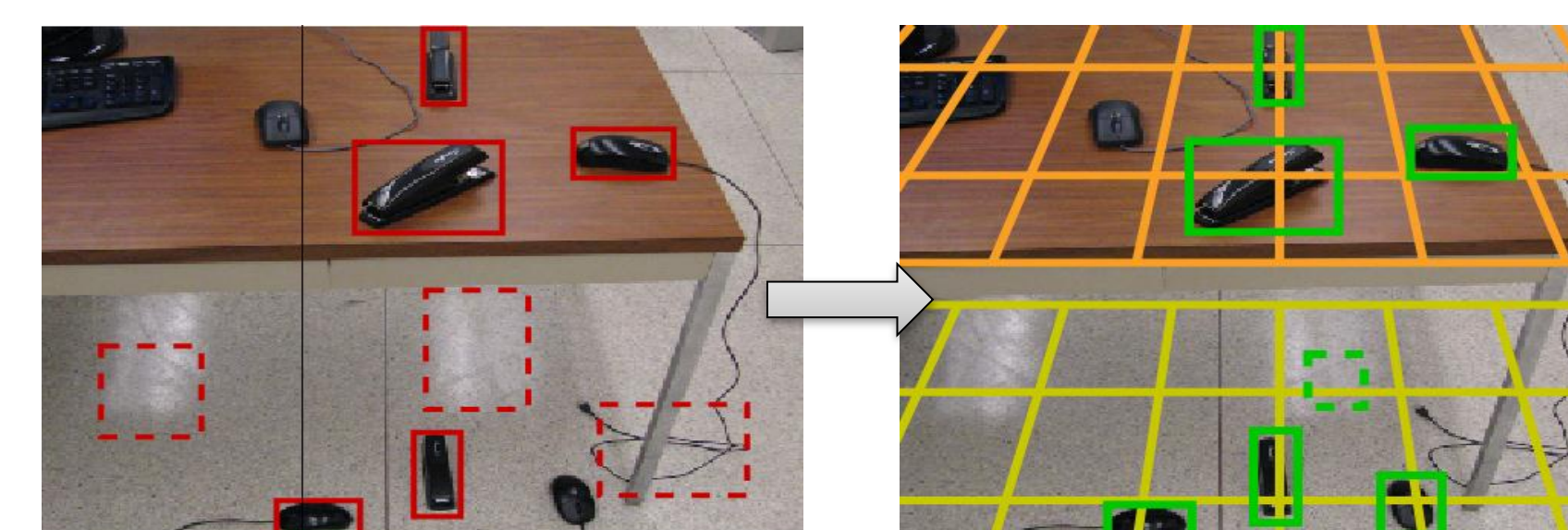  • Recover 3D supporting planes
  • Locate object in 3D space



Baseline Object Detector

**Joint** Probabilistic Optimization

• **Tools**
  • Novel relationship between object's pose & location, and supporting plane
  • Layout priors



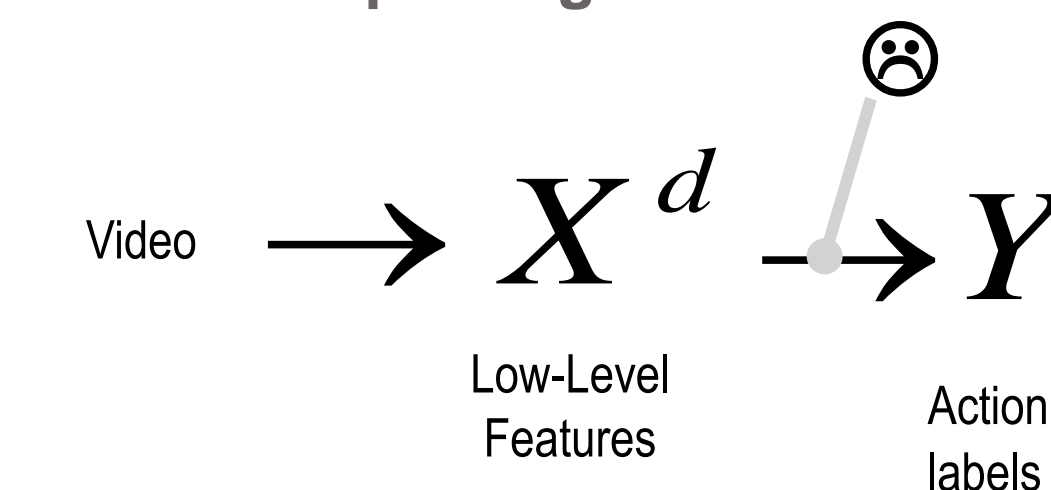Normal   Normal   Normal

• Experimental Results



### References

• J. Liu, B. Kuipers, S. Savarese, *Recognizing Human Actions by Attributes*, Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2011
• J. Liu, M. Shah, B. Kuipers, S. Savarese, *Cross-View Action Recognition via View Knowledge Transfer*, Proceedings of the IEEE International Conference on Computer Vision and Pattern Recognition (CVPR), 2011.
• Y. Bao, M. Sun & S. Savarese, Toward coherent object detection and scene layout understanding. CVPR, 2010.
• M. Sun, G. Bradsky, B. Xu & S. Savarese, Depth-encoded Hough voting for joint object detection and shape recovery. ECCV, 2010.
• G. Tsai, C. Xu, J. Liu & B. Kuipers, Real-time indoor scene understanding using Bayesian filtering with motion cues. ICCV, 2011.
• C. Xu & B. Kuipers, Towards the Object Semantic Hierarchy. ICDL, 2010.

## Learning Human Actions by Attributes

**Traditional paradigm**

Video $\rightarrow X^d \rightarrow Y$
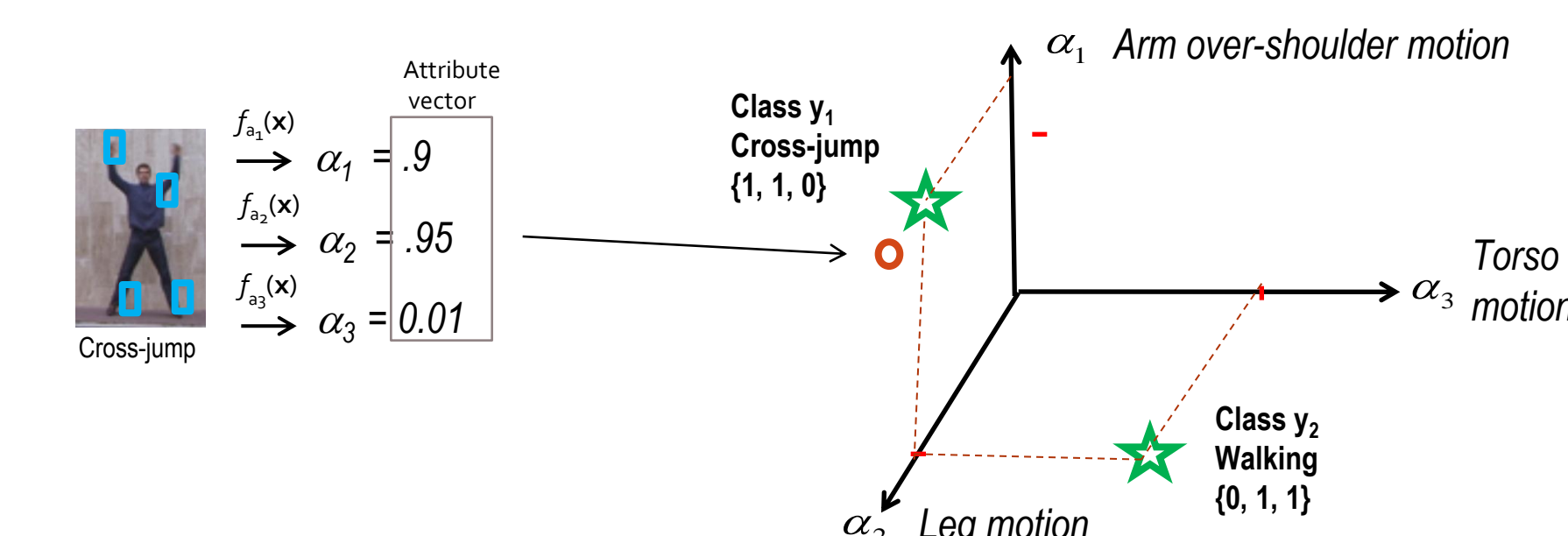
Low-Level Features      Action class labels

• Rich visual temporal-spatial structures cannot be well characterized by a single class label
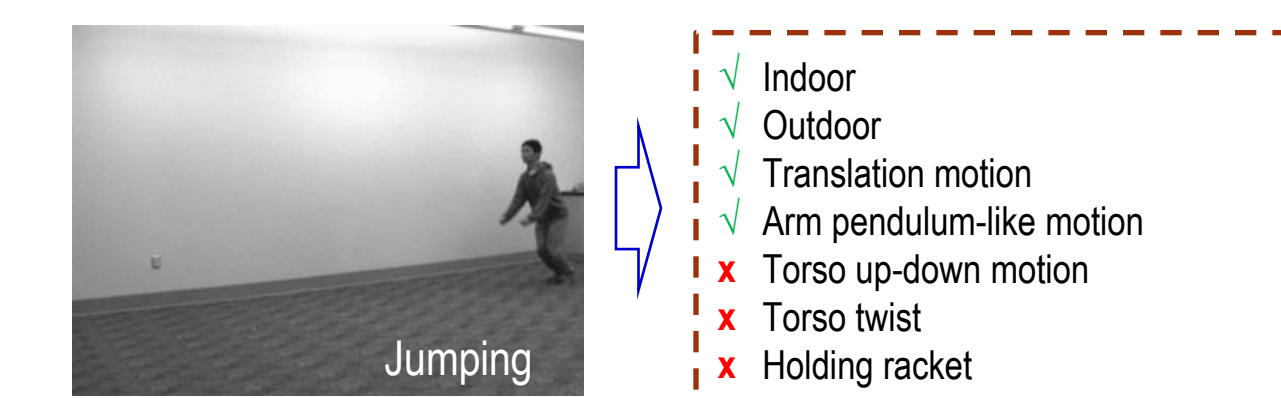• For complex activities this process is too restrictive and reductive

**Proposed paradigm**

• The action classifier $F: X^d \rightarrow Y$ can be decomposed into:

$$S: X^d \rightarrow A^m \qquad L: A^m \rightarrow Y$$



Attribute vector

$f_{\alpha_1}(x) \rightarrow \alpha_1 = .9$
$f_{\alpha_2}(x) \rightarrow \alpha_2 = .95$
$f_{\alpha_3}(x) \rightarrow \alpha_3 = 0.01$

Cross-jump

Class $y_1$ Cross-jump {1, 1, 0}

Class $y_2$ Walking {0, 1, 1}

$\alpha_1$ Arm over-shoulder motion
$\alpha_2$ Leg motion
$\alpha_3$ Torso motion

• **Experimental Results**



Jumping

Indoor
Outdoor
Translation motion
Arm pendulum-like motion
✗ Torso up-down motion
✗ Torso twist
✗ Holding racket

| | Average Accuracy (%) |
|---|---|
| raw-feature | 51.83 |
| specified attributes | 60.48 |
| raw-feature + specified attributes | 63.60 |
| data-driven attributes | 45.31 |
| raw-feature + all attributes | 65.09 |

Activity data set Niebles et al. 2010

### Acknowledgements