

CERTH CENTRE FOR ESEARCH & TECHNOLOGY HELLAS



### Motivation

- Object recognition constitutes an open research challenge Scope: Hand and object segmentation ✓ Almost exclusively use of static appearance features • Output: a) colorized object depth maps, b) colorized hand depth maps, and c) ✓ Challenges: object appearance variance, occlusions, deformations, colorized hand 3D flow magnitude fields illumination variation Findings in cognitive neuroscience  $\checkmark$  Object perception is based on the fusion of sensory (object appearance) 3D Volume o and motor (human-object interaction) information Interest Cropped  $\checkmark$  Object affrdances: the types of actions that humans typically perform when interacting with them 300x300 Proposed approach 512x424 Accumulated evidence on ventral and dorsal stream interaction Depth Map (512x424) **Generalized Template Matching (GTM) architecture** Jorsal stream Object Appearance interaction decision: Information stream Scope: Estimation of co-Object appearance mplex multi-level affor-Colorized Ventral stream Object Depth dance-related patterns Novel contributions on sensorimotor object recognition along the spatial dimen-✓ Neurobiologically and neurophysiologically grounded Neural Network architectures sions Evaluation of multiple recent neuro-scientific findings **Colorized Hand 3D** Flow Based solely on the use ✓ Large number of complex affordance types Magnitude ✓ Large-scale public RGB-D object recognition dataset of CNNs Affordan Introduced SOR3D dataset Figure Top: appearance CNN 14 object categories for object recognition and 13 affordance types affordance CNN (single-stream Bottom: Detailed model). • 54 object-affordance GTM topology the of combinations architecture for: a) late fusion at FC layer, b) late fusion at 105 subjects last CONV layer, c) single-level • Over 20,800 instances slow fusion and d) multi-level slow fusion. Appearance information stream



- 3 synchronized Kinect II
- http://sor3d.vcl.iti.gr



# **Deep Affordance-grounded Sensorimotor Object Recognition**

## Spyridon Thermos, Georgios Th. Papadopoulos, Petros Daras, Gerasimos Potamianos

Information Technologies Institute, Centre for Research and Technology Hellas, Greece <sup>2</sup> Department of Electrical and Computer Engineering, University of Thessaly, Greece

### Visual front-end

Colorized depth map





### **Generalized Spatio-Temporal (GST) architecture**

- Scope: Encoding of the time-evolving procedures of the performed human actions
- Composite CNN-LSTM NN considered
- Support of asynchronous fusion

Figure Top: affordance CNN-LSTM (single-stream model). Bottom: Detailed GST architecture for: a) late fusion and b) slow fusion.

- Up to 29% relative erro reduction, compared to the baseline model
- Recognition performance all supported object Of types boosted

### **Conclusions & future work**

- probabilistic approaches of the literature.



#### **Experimental results**

	Fusion architecture	Accuracy (%)
or e	Baseline CNN (appearance only)	85.12
	GTM late	88.24
	GTM slow single level	88.13
	GTM slow multi-level	89.43
е	GST late	86.50
ct	GST slow single level	79.65
	Product Rule	73.45
	SVM	83.43
	Bayes	75.86

• The proposed NN-based sensorimotor approach outperforms similar

• Future work will investigate the modeling of the human-object interactions in more details and the application of the proposed methodology to "in the wild" object recognition scenarios.