

# Temporal Saliency Adaption in Egocentric Videos



Panagiotis Linardos



Eva Mohedano



Monica Cherto

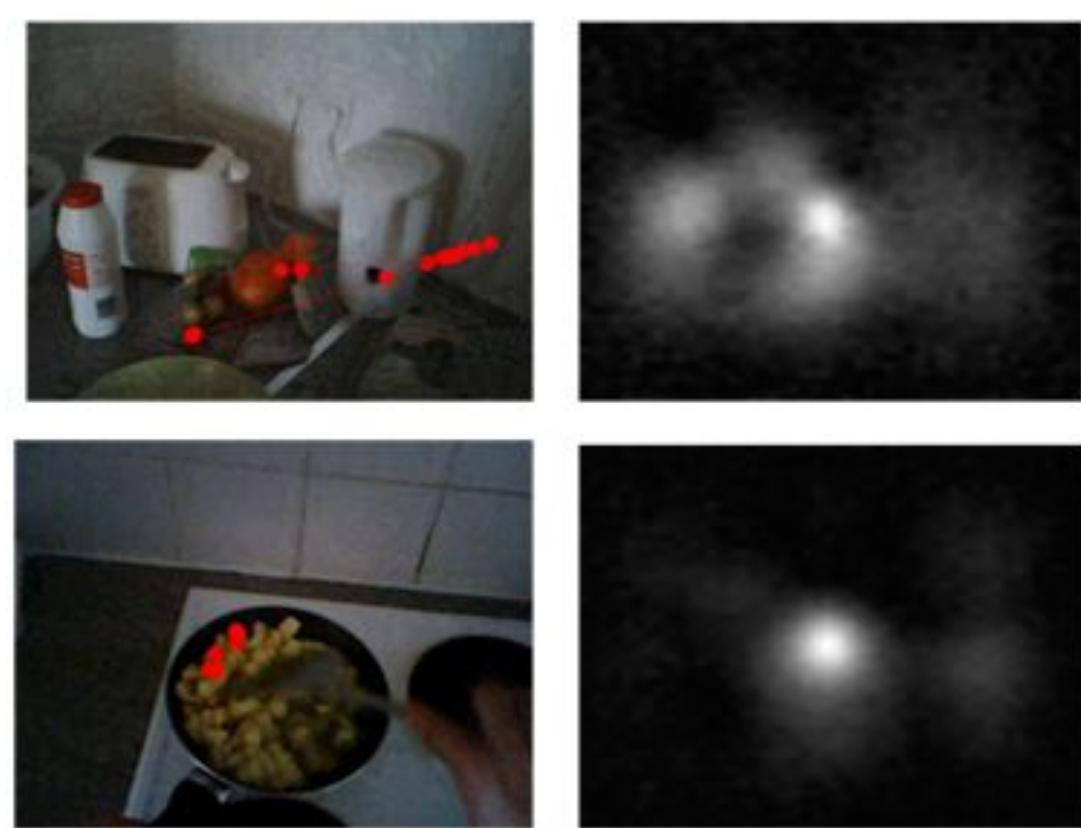


Cathal Gurrin



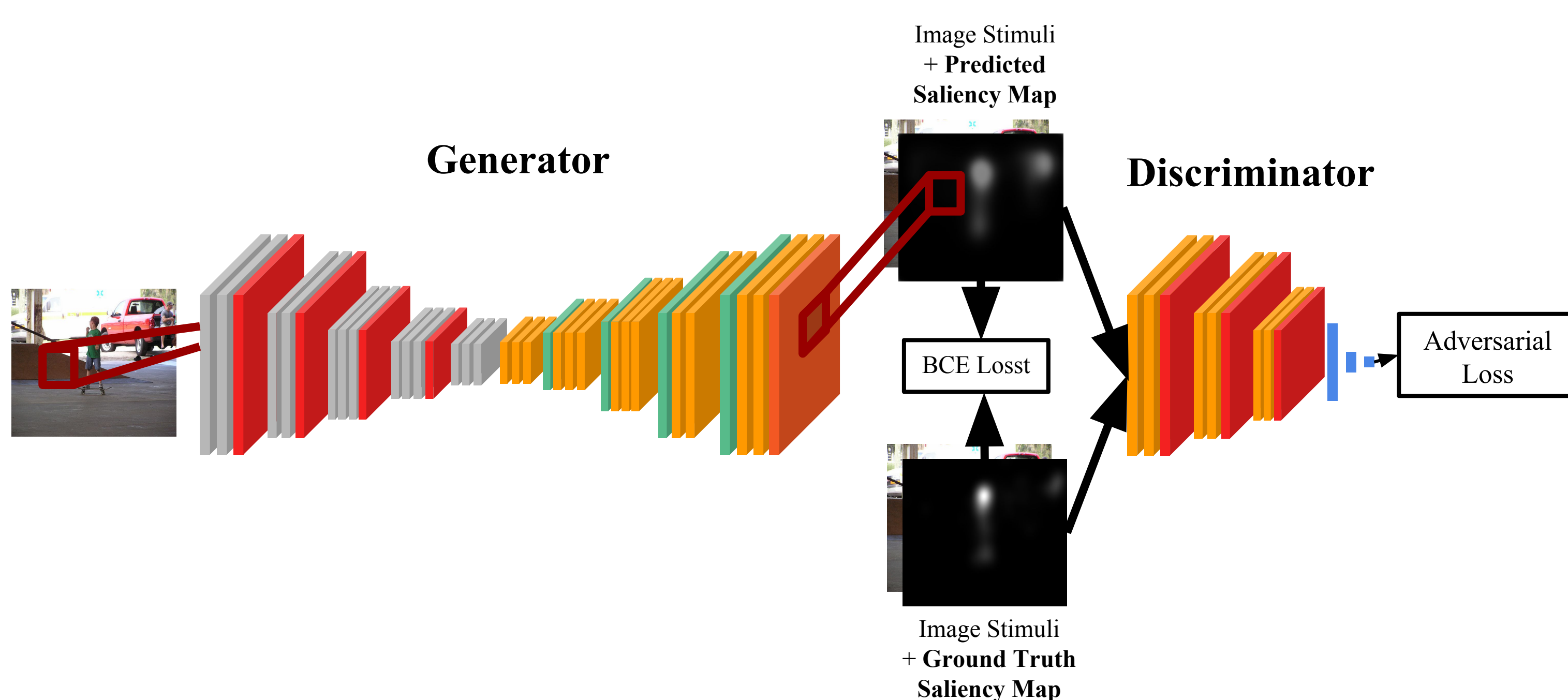
Xavier Giro-i-Nieto

## Motivation

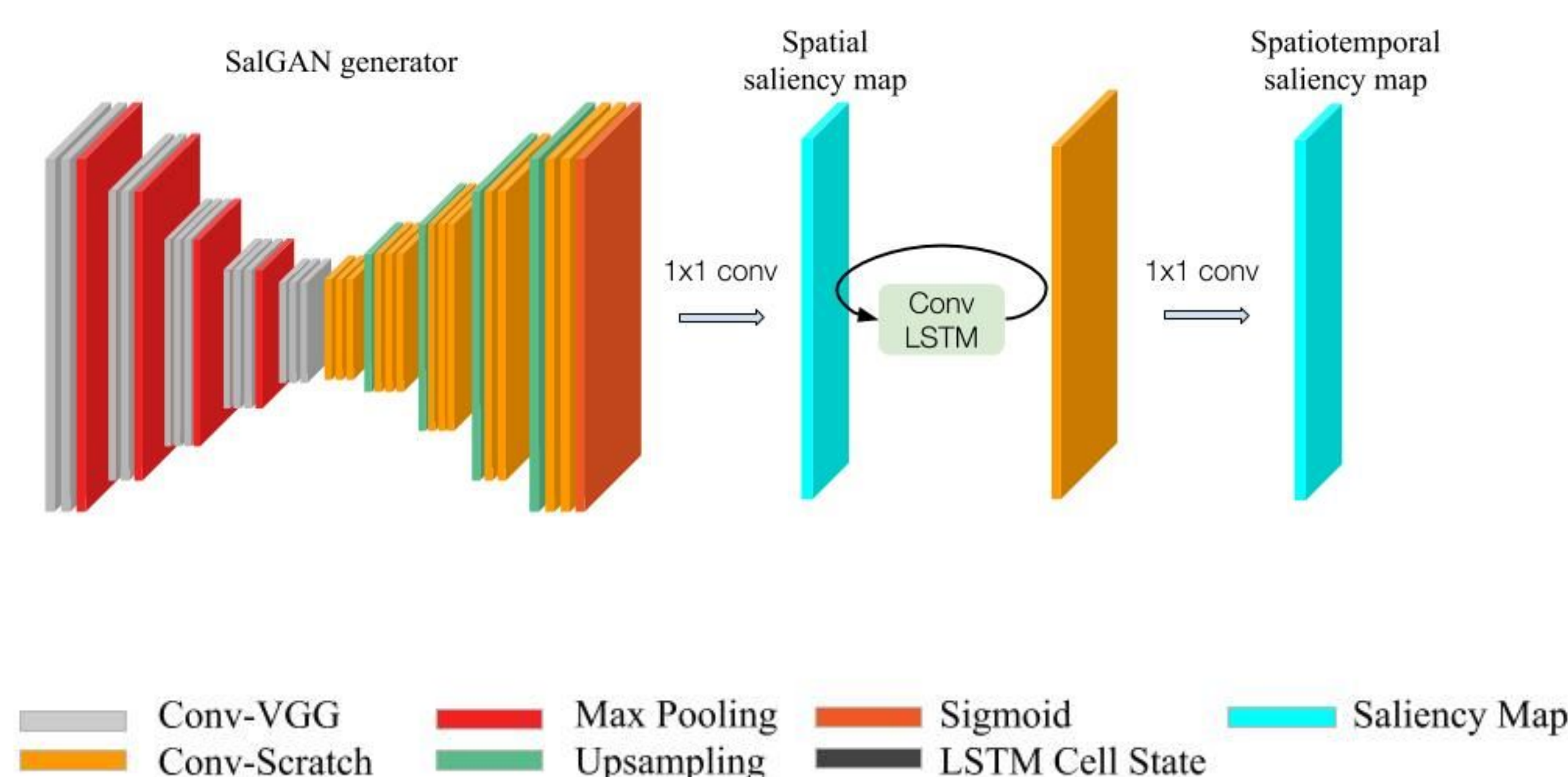


- Saliency prediction refers to the task of estimating which regions of an image have a higher probability of being observed by a viewer.
- It has been shown that this information can be helpful to improve tasks such as activity recognition or object detection.
- In this work, we focus in the task of **predicting saliency in egocentric videos**. For that, we extend a state of the art saliency model to the temporal domain.

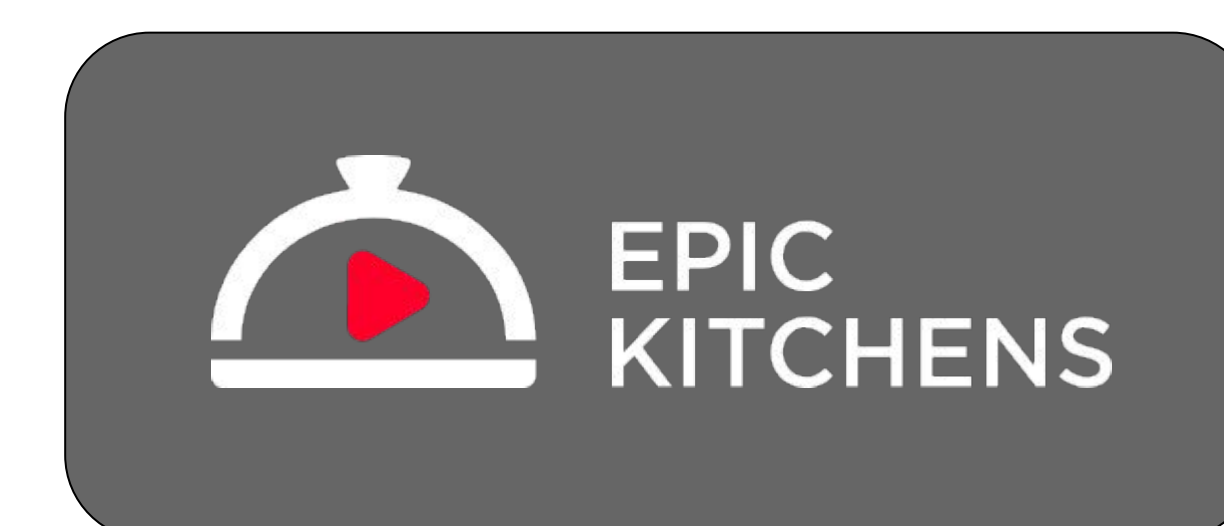
## SalGAN [1]



## SalGAN [1] + convLSTM



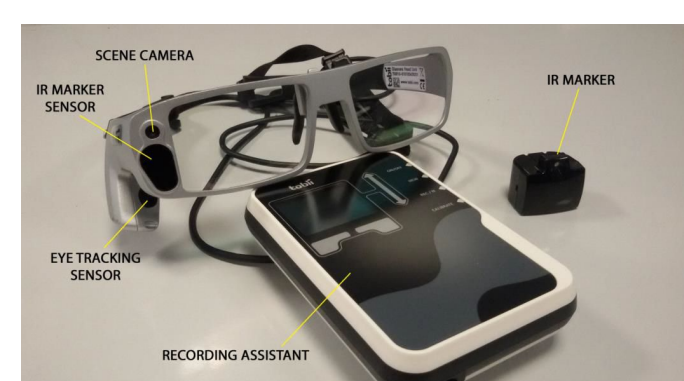
Saliency maps from both configurations available for [2]:



## EgoMon Gaze & Video Dataset

New egocentric dataset

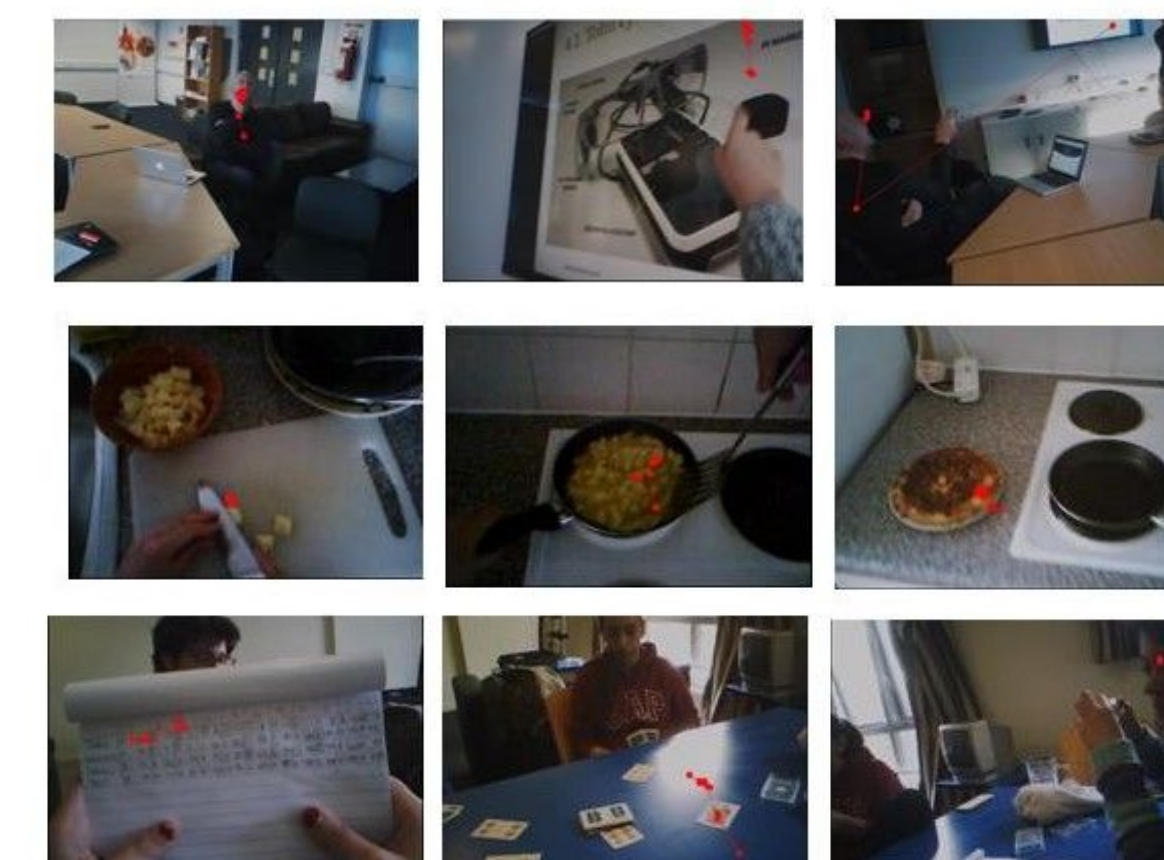
- 7 sequences: 4 free-viewing + 3 task-oriented
- Average length of 30 minutes
- 3 different users wearing Tobii glasses



### Free-viewing activities



### Task-oriented activities



...and also images from Narrative clip for one sequence..



## Results

### Performance on the DHF1K dataset

	AUC-J	sAUC	NSS	CC	SIM
DHF1K SoA	0.885	0.553	2.259	<b>0.415</b>	0.311
SalGAN	<b>0.930</b>	<b>0.834</b>	<b>2.468</b>	0.372	<b>0.264</b>
+ conv	0.743	0.723	2.208	0.303	0.261
+ convLST	0.744	0.722	2.246	0.302	0.260

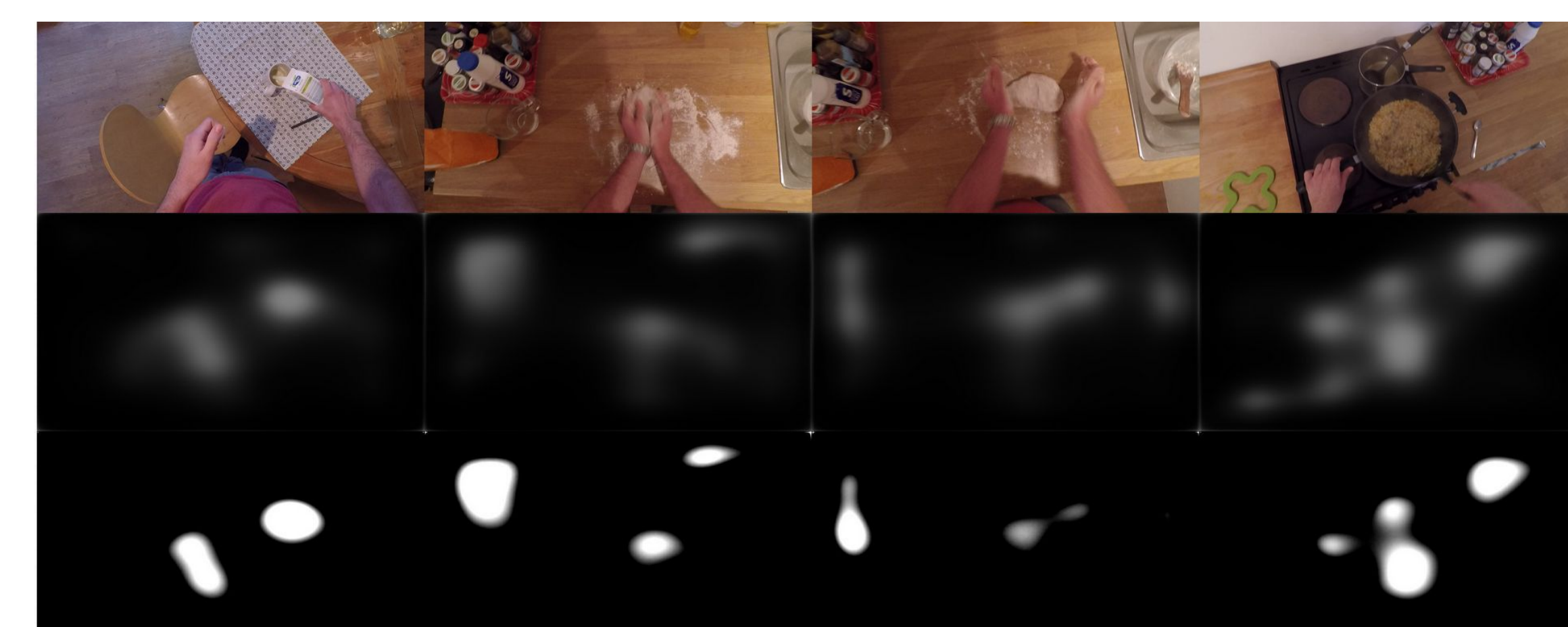
### Performance on EgoMon tasks (NSS metric)

	Free-Viewing	Task-oriented	Total
SalGAN	<b>2.652</b>	1.313	<b>2.079</b>
+ conv	0.805	1.694	1.249
+ convLST	0.904	<b>1.705</b>	1.247

Video frames

SalGAN

+ convLSTM



Results indicate that SalGAN represents a strong baseline for video-saliency prediction. However, our strategy of including time information based directly on saliency predictions ignores any semantic information of the scene, so we will extend the work by working with lower SalGAN layers.

