

# DEEP LEARNING FOR COMPUTER VISION

Summer School at UPC TelecomBCN Barcelona. June 28-July 4, 2018



## Instructors



Organized by



UNIVERSITAT POLITÈCNICA  
DE CATALUNYA  
BARCELONATECH



Supported by



Co-funded by the  
Erasmus+ Programme  
of the European Union

GitHub Education

Google Cloud Platform

+ info: <http://bit.ly/dlcv2018>

<http://bit.ly/dlcv2018>



Day 3 Lecture 5

## Saliency Prediction



Kevin McGuinness

[kevin.mcguinness@dcu.ie](mailto:kevin.mcguinness@dcu.ie)

Assistant Professor  
School of Electronic Engineering  
Dublin City University



Insight  
Centre for Data Analytics

# The importance of visual attention







# The importance of visual attention



# The importance of visual attention







# The importance of visual attention



# Why don't we see the changes?

We don't really see the whole image

We only focus on small specific regions: the **salient** parts

Human beings reliably attend to the same regions of images  
when shown



# What we perceive



# Where we look



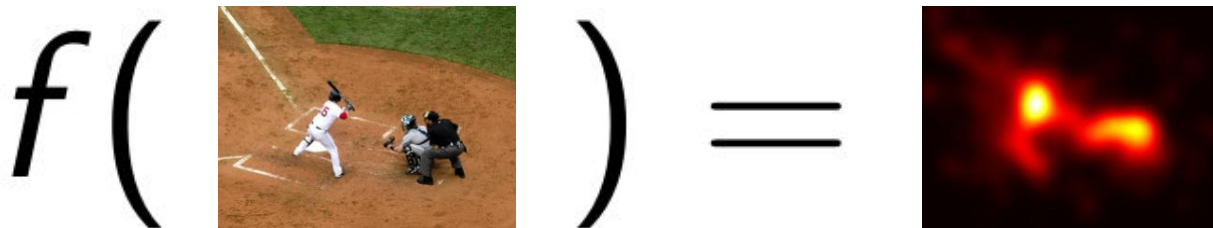
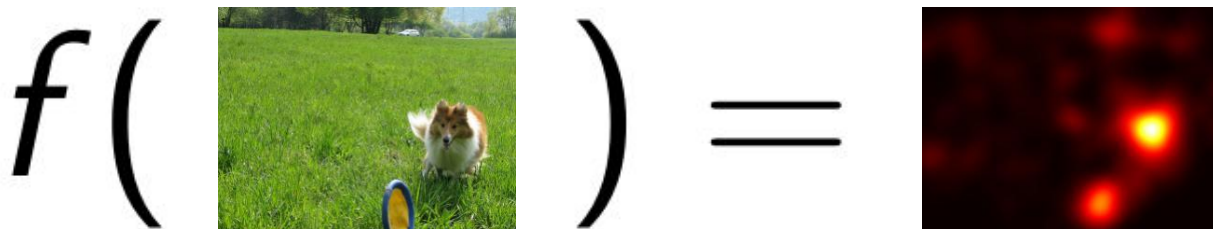
# What we actually see



# Saliency prediction

Produce a **computational model of visual attention**: predict where humans will look.

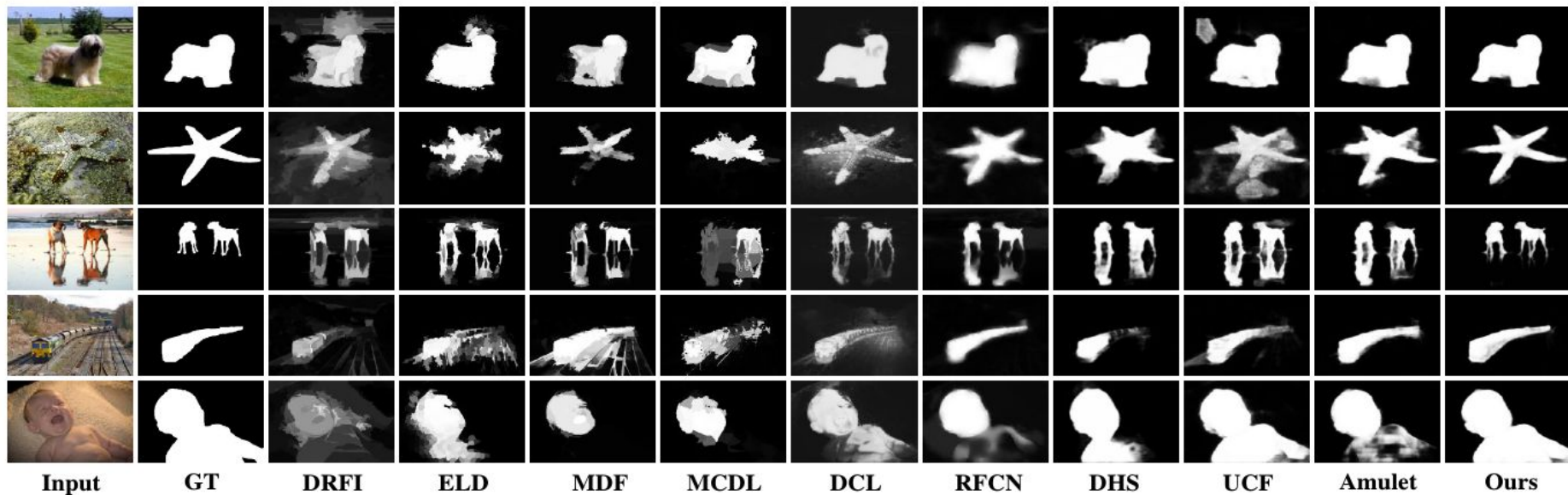
Often want to map an image to a **heatmap** (saliency map).





# Salient object detection?

Often confused with saliency prediction, but a different task.



# Datasets

# MIT 300

300 natural indoor and outdoor scenes.

39 observers. 3 sec free view.

ETL 400 ISCAN eye tracker

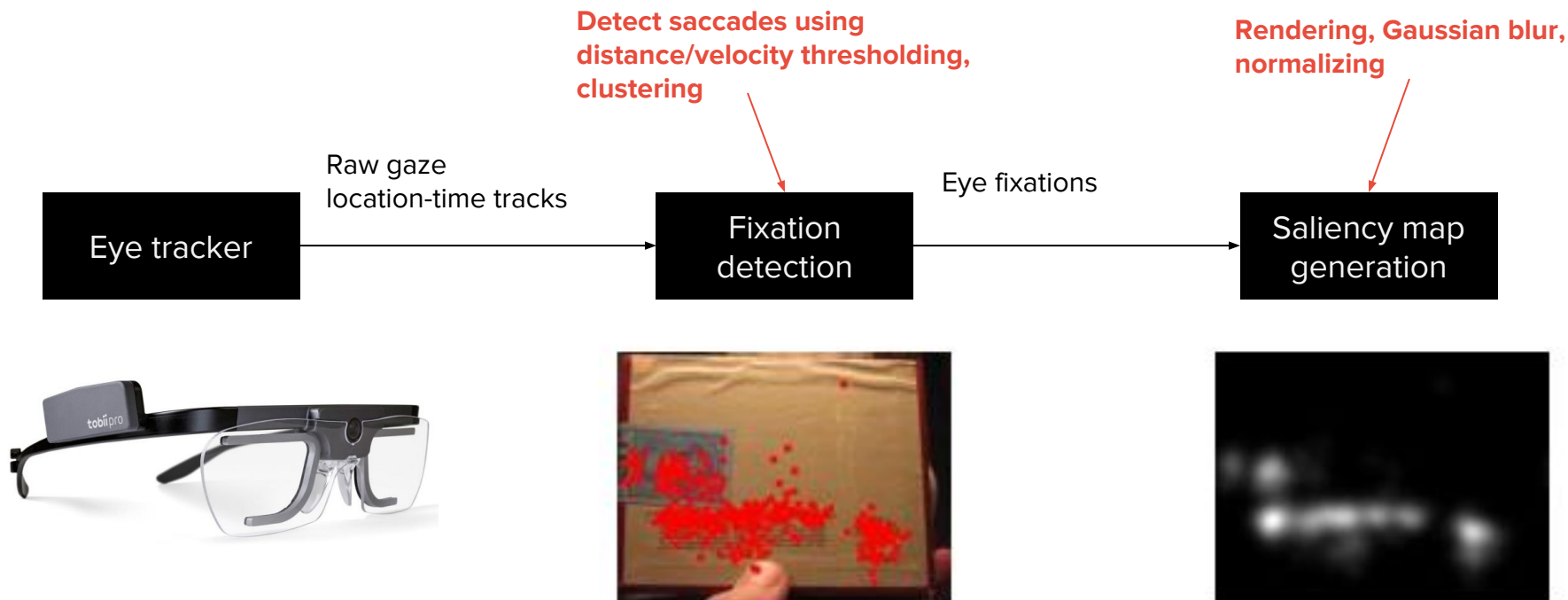
**Test set only:** no training data  
or public ground truth

[http://saliency.mit.edu/results\\_mit300.html](http://saliency.mit.edu/results_mit300.html)



# Fixations and saliency maps

Raw eye tracker data needs to be processed to produce saliency maps





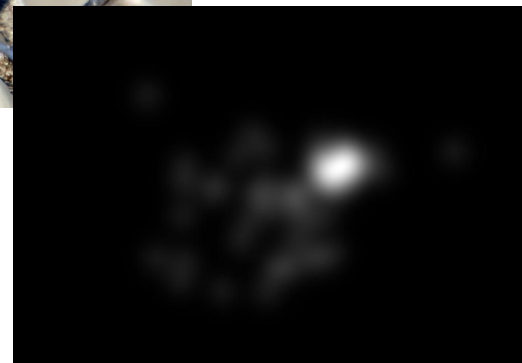
# MIT 1003

1003 natural indoor and outdoor scenes.

15 observers. 3 sec free view.

ETL 400 ISCAN eye tracker

**Training dataset for MIT 300**



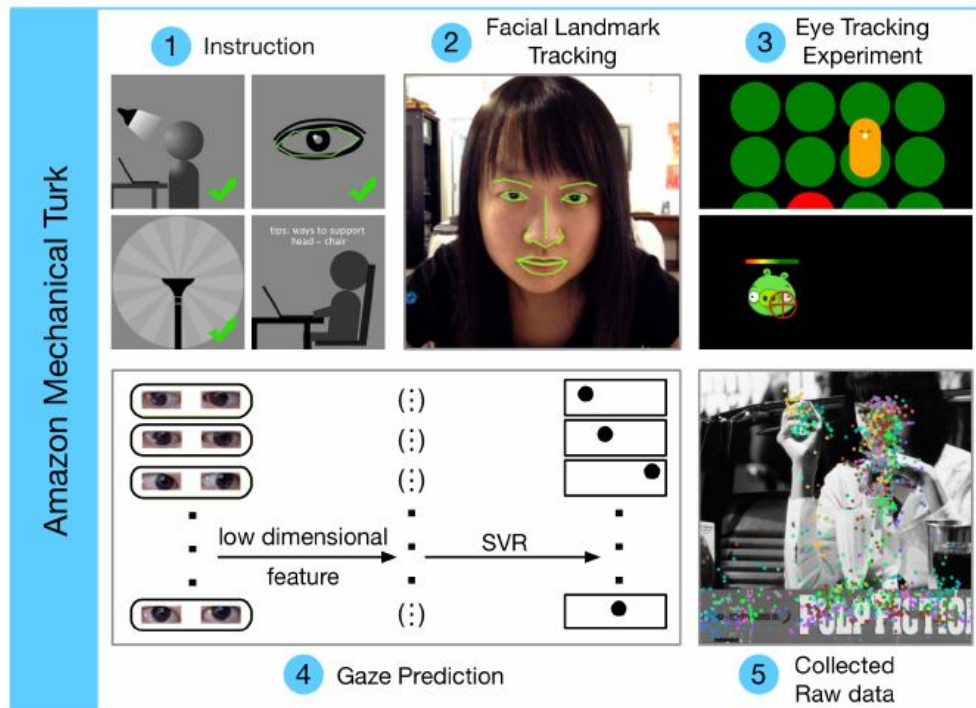
# iSUN

Large scale dataset of natural scenes

20,608 images with avg. 3 observers each

Collected using webcams and Amazon Mechanical Turk

Used in [LSUN challenge](#)  
2015/2016



# SALICON

Another large scaled dataset of images from MS COCO dataset

10K train, 5K val, 5K test

Simulated crowdsourced attention using **mouse** movements and **simulated artificial foveation**.



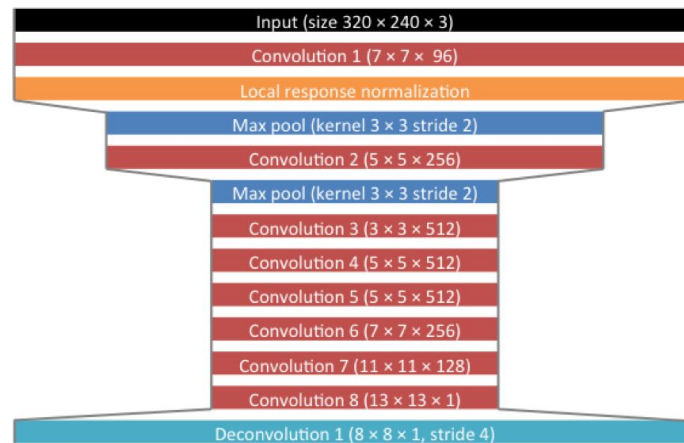
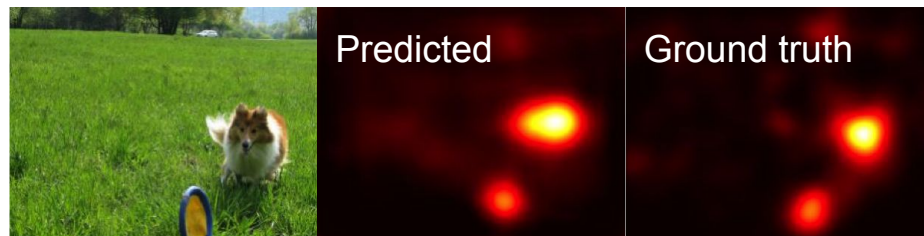
# Models



# SalNet: deep visual saliency model

Predict map of visual attention from image pixels  
(find the parts of the image that stand out)

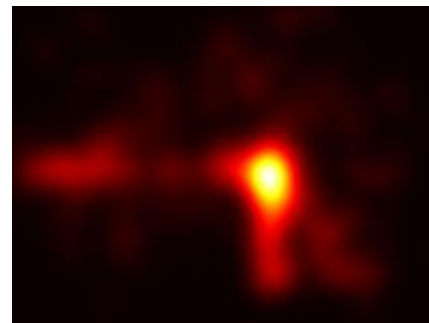
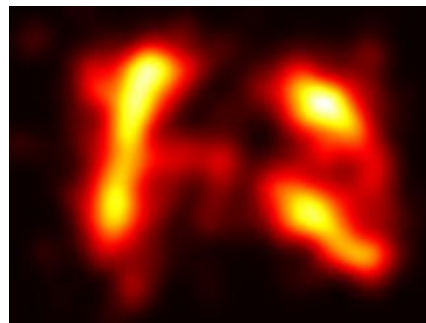
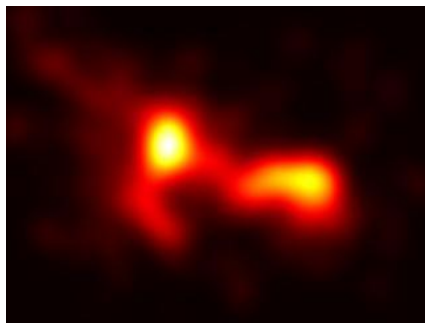
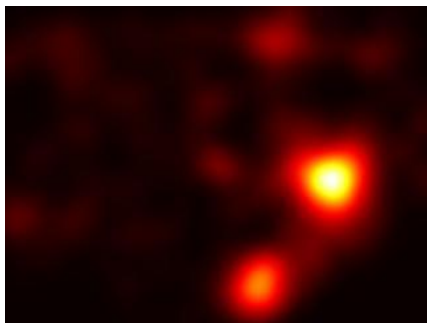
- Feedforward 8 layer “fully convolutional” architecture
- Transfer learning in bottom 3 layers from pretrained VGG-M model on ImageNet
- Trained on SALICON dataset



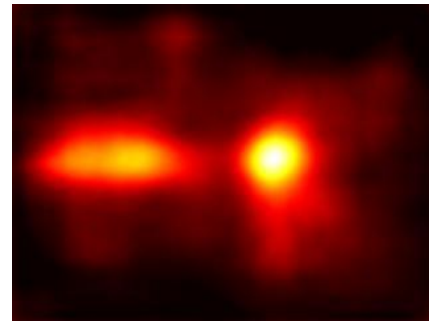
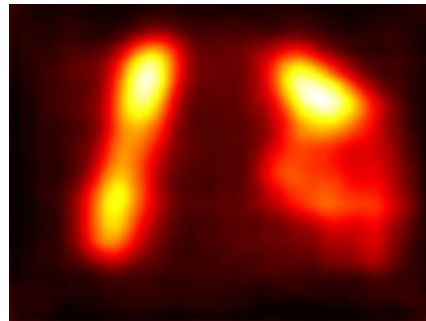
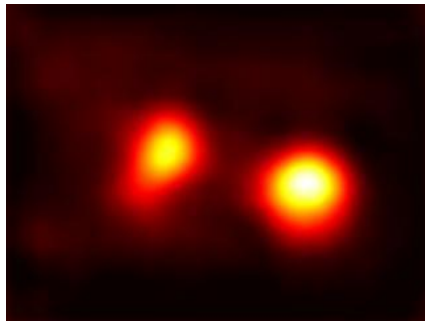
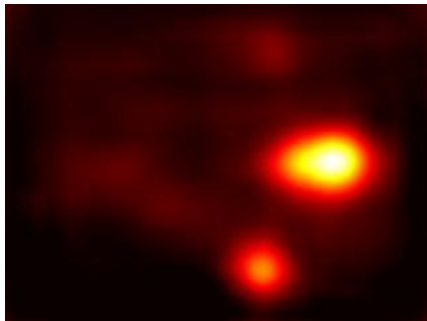
Image



Ground truth



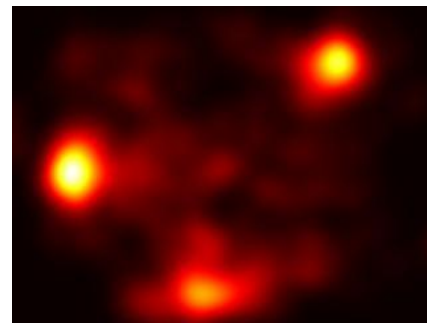
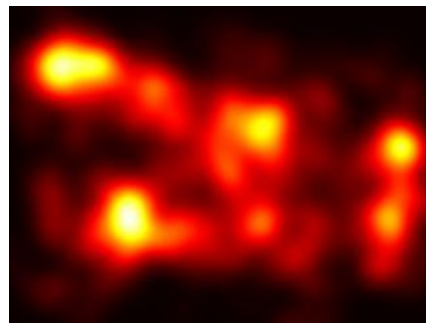
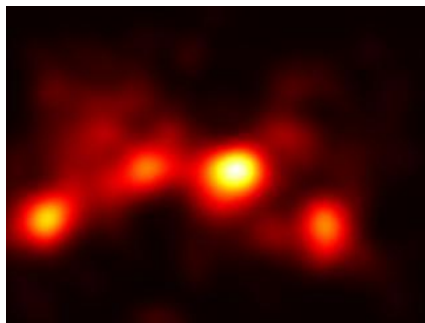
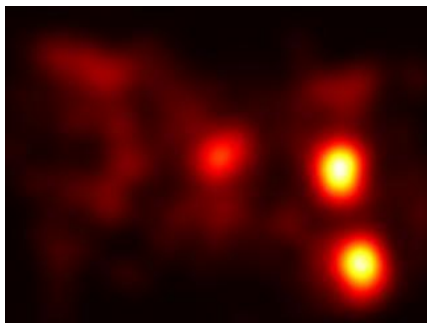
Prediction



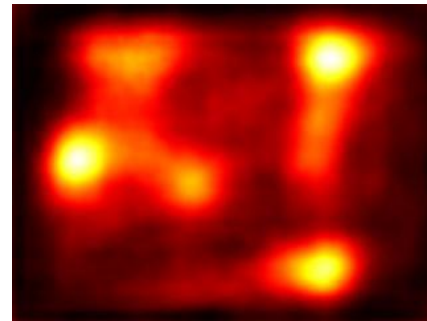
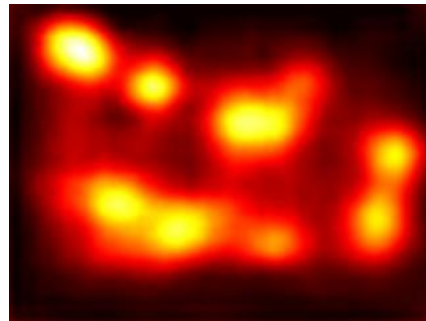
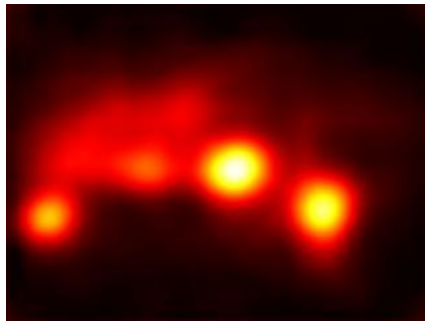
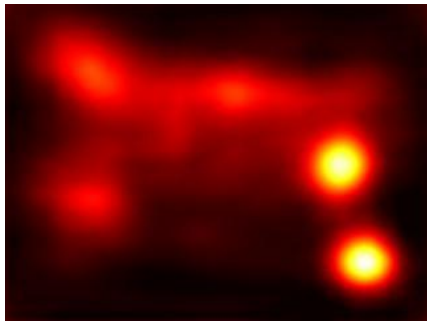
Image



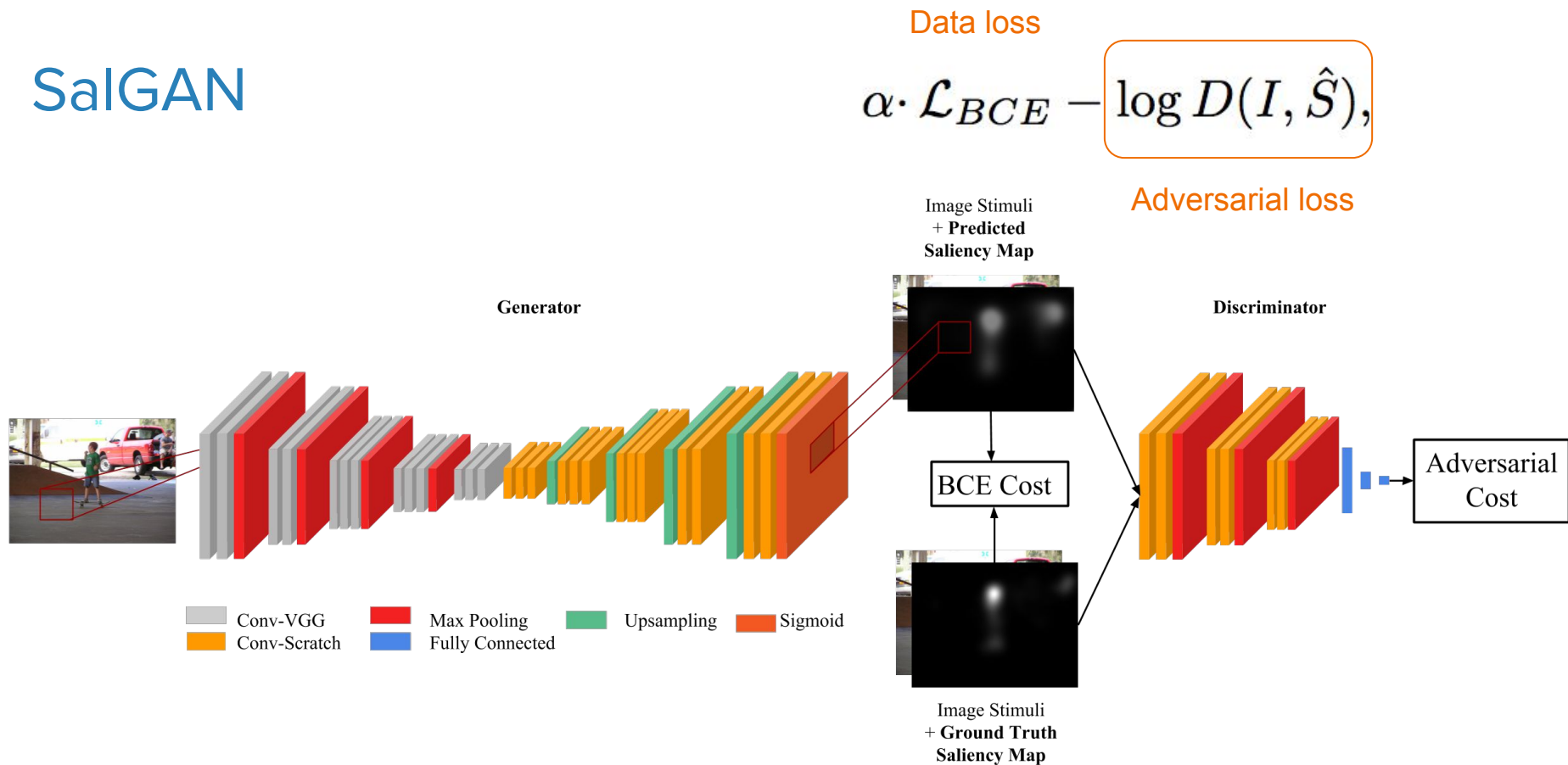
Ground truth



Prediction



# SalGAN





# SalNet and SalGAN benchmarks

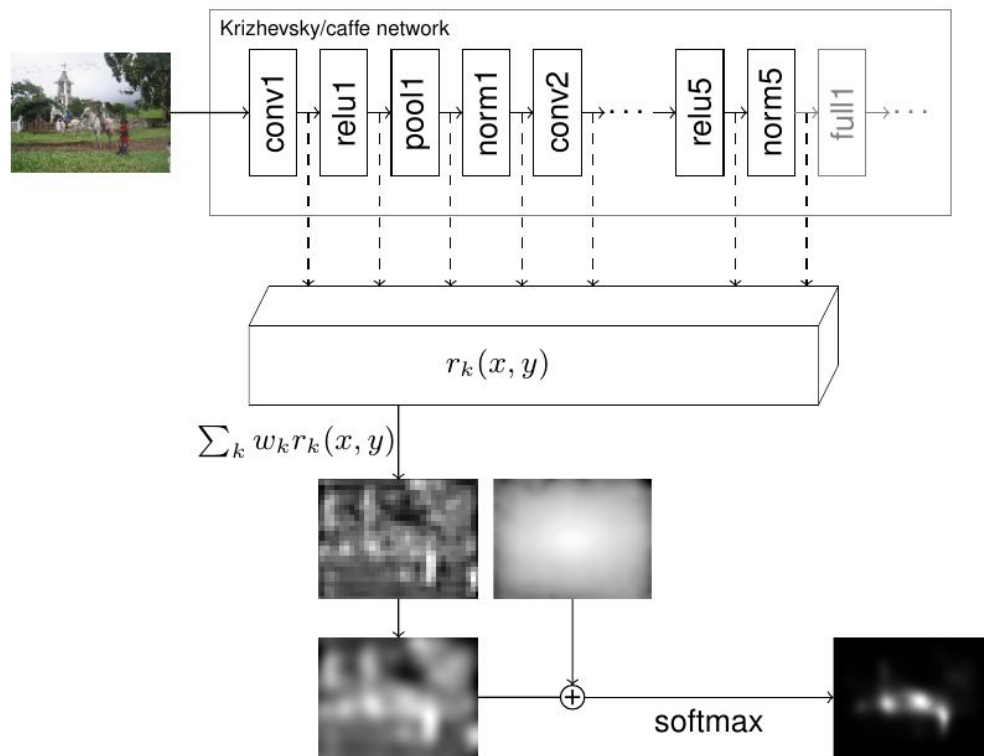
SALICON (test)	AUC-J $\uparrow$	Sim $\uparrow$	EMD $\downarrow$	AUC-B $\uparrow$	sAUC $\uparrow$	CC $\uparrow$	NSS $\uparrow$	KL $\downarrow$
DSCLRCN [24](*)	-	-	-	0.884	0.776	0.831	3.157	-
<b>SalGAN</b>	-	-	-	<b>0.884</b>	<b>0.772</b>	<b>0.781</b>	<b>2.459</b>	-
ML-NET [5]	-	-	-	(0.866)	(0.768)	(0.743)	2.789	-
SalNet [25]	-	-	-	(0.858)	(0.724)	(0.609)	(1.859)	-
MIT300	AUC-J $\uparrow$	Sim $\uparrow$	EMD $\downarrow$	AUC-B $\uparrow$	sAUC $\uparrow$	CC $\uparrow$	NSS $\uparrow$	KL $\downarrow$
Humans	0.92	1.00	0.00	0.88	0.81	1.0	3.29	0.00
Deep Gaze II [21](*)	0.88	(0.46)	(3.98)	0.86	0.72	(0.52)	(1.29)	(0.96)
DSCLRCN [24](*)	0.87	0.68	2.17	(0.79)	0.72	0.80	2.35	0.95
DeepFix [17](*)	0.87	0.67	2.04	(0.80)	(0.71)	0.78	2.26	0.63
SALICON [9]	0.87	(0.60)	(2.62)	0.85	0.74	0.74	2.12	0.54
<b>SalGAN</b>	<b>0.86</b>	<b>0.63</b>	<b>2.29</b>	<b>0.81</b>	<b>0.72</b>	<b>0.73</b>	<b>2.04</b>	<b>1.07</b>
PDP [11]	(0.85)	(0.60)	(2.58)	(0.80)	0.73	(0.70)	2.05	0.92
ML-NET [5]	(0.85)	(0.59)	(2.63)	(0.75)	(0.70)	(0.67)	2.05	(1.10)
Deep Gaze I [19]	(0.84)	(0.39)	(4.97)	0.83	(0.66)	(0.48)	(1.22)	(1.23)
iSEEL [29](*)	(0.84)	(0.57)	(2.72)	0.81	(0.68)	(0.65)	(1.78)	0.65
SalNet [25]	(0.83)	(0.52)	(3.31)	0.82	(0.69)	(0.58)	(1.51)	0.81
BMS [31]	(0.83)	(0.51)	(3.35)	0.82	(0.65)	(0.55)	(1.41)	0.81

# Deep Gaze

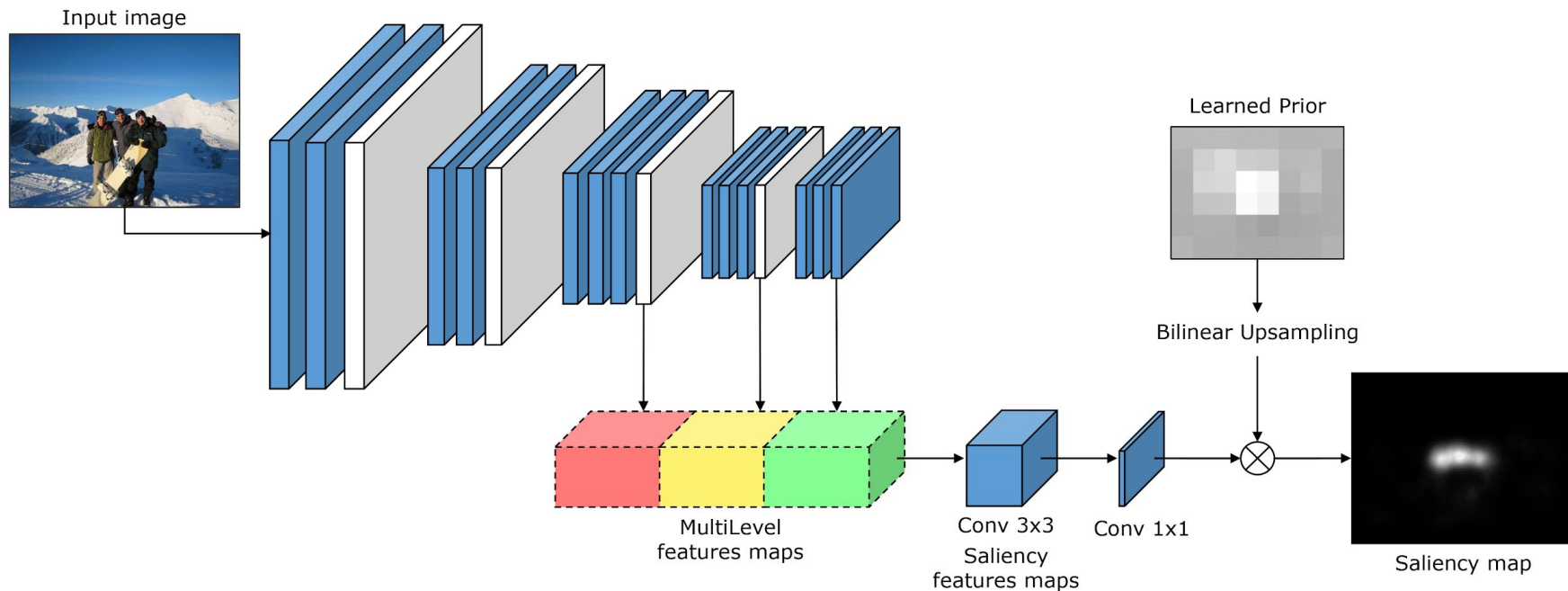
Simple linear model trained on activations of all conv layers (upsampled) from AlexNet

Softmax output over full image, categorical cross entropy.

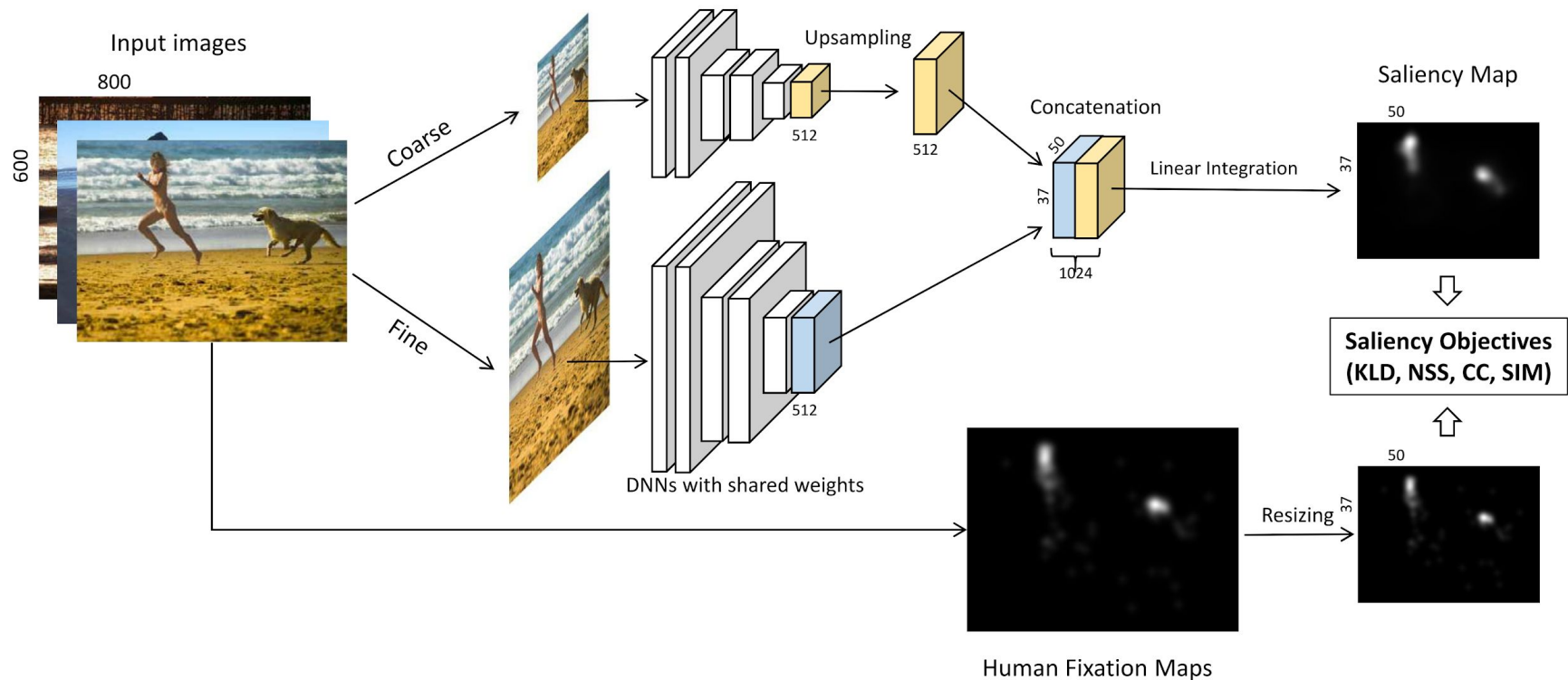
$L_1$  regularization used to encourage sparsity.



# MLNet



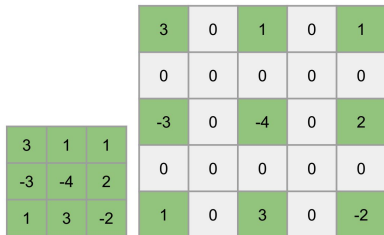
# SALICON



# DeepFix

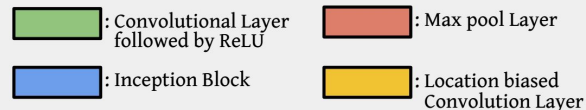
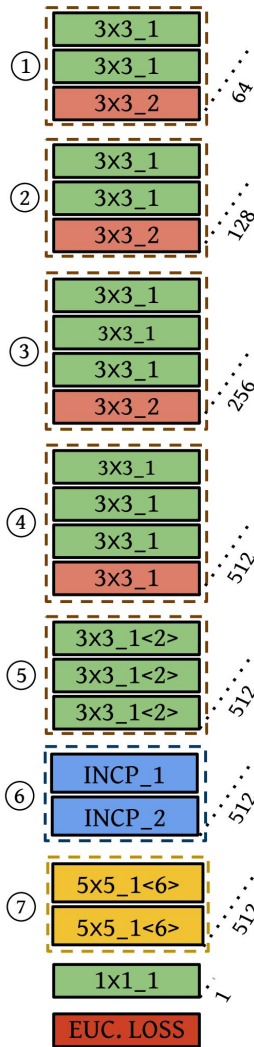
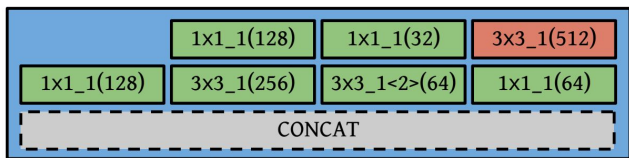
Kruthiventi et al. **DeepFix: A Fully Convolutional Neural Network for predicting Human Eye Fixations**

<https://arxiv.org/abs/1510.02927>



Dilated convolutions

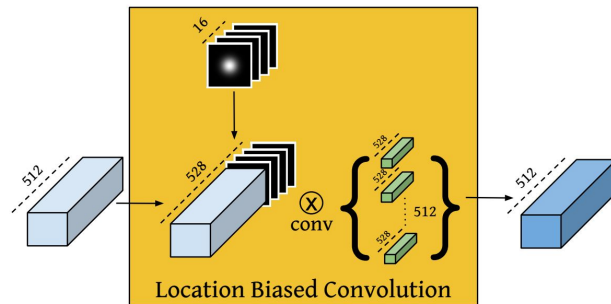
Inception layers



$w \times h_s \times h$ : Layer with kernels of width - w height - h stride - s hole - h

..... : No. of channels in the block's output

Weights initialized from VGG16 trained on ImageNet

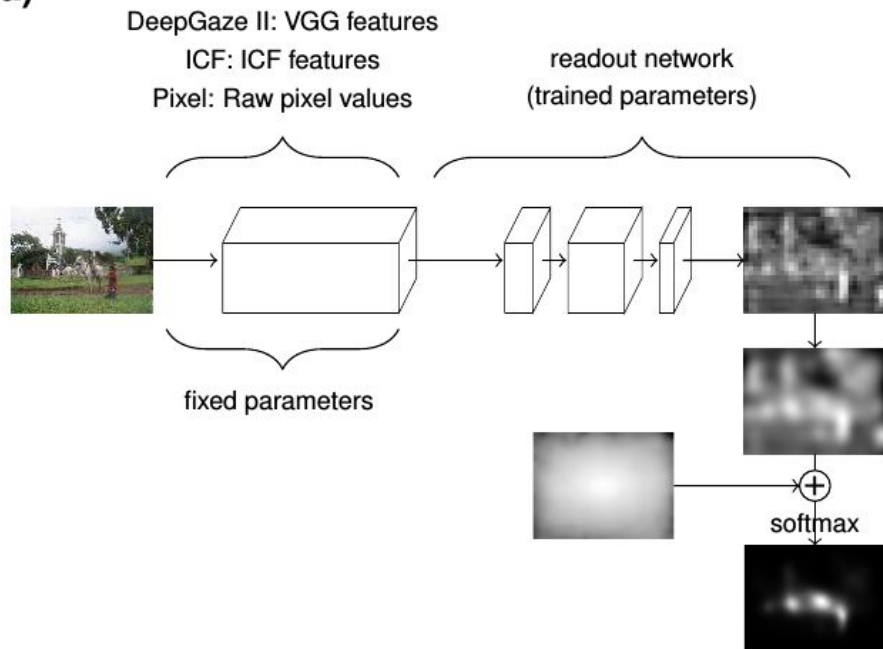


Location biased convolutions

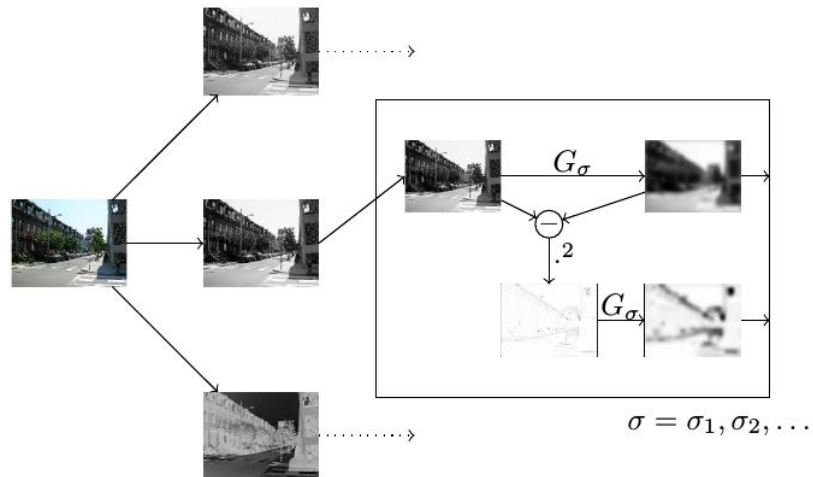


# Deep Gaze II

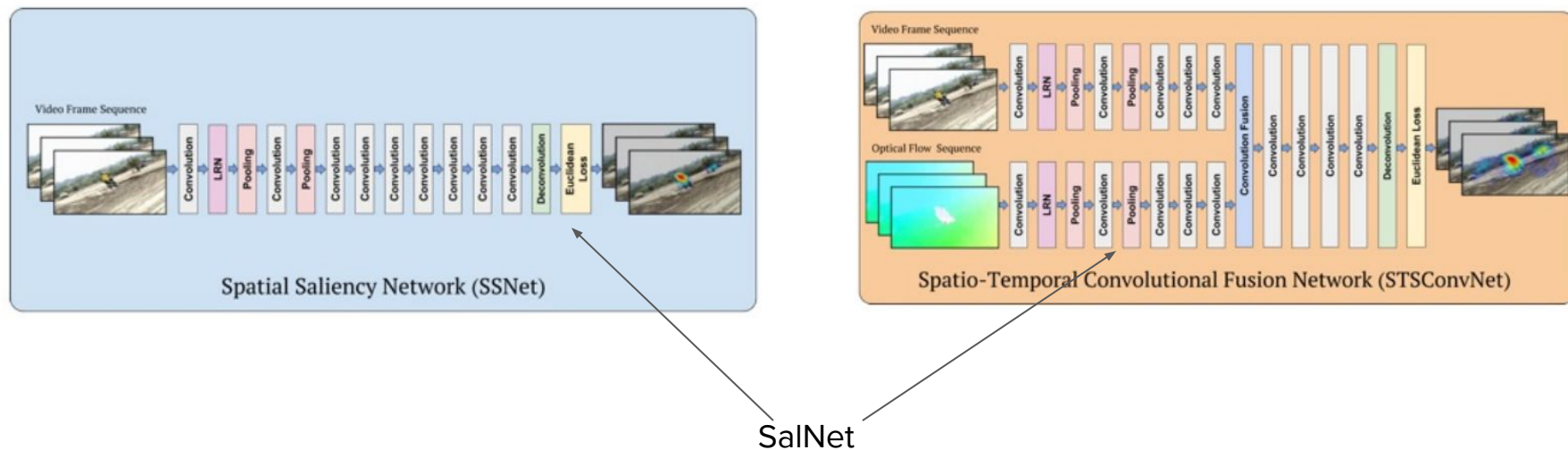
(a)



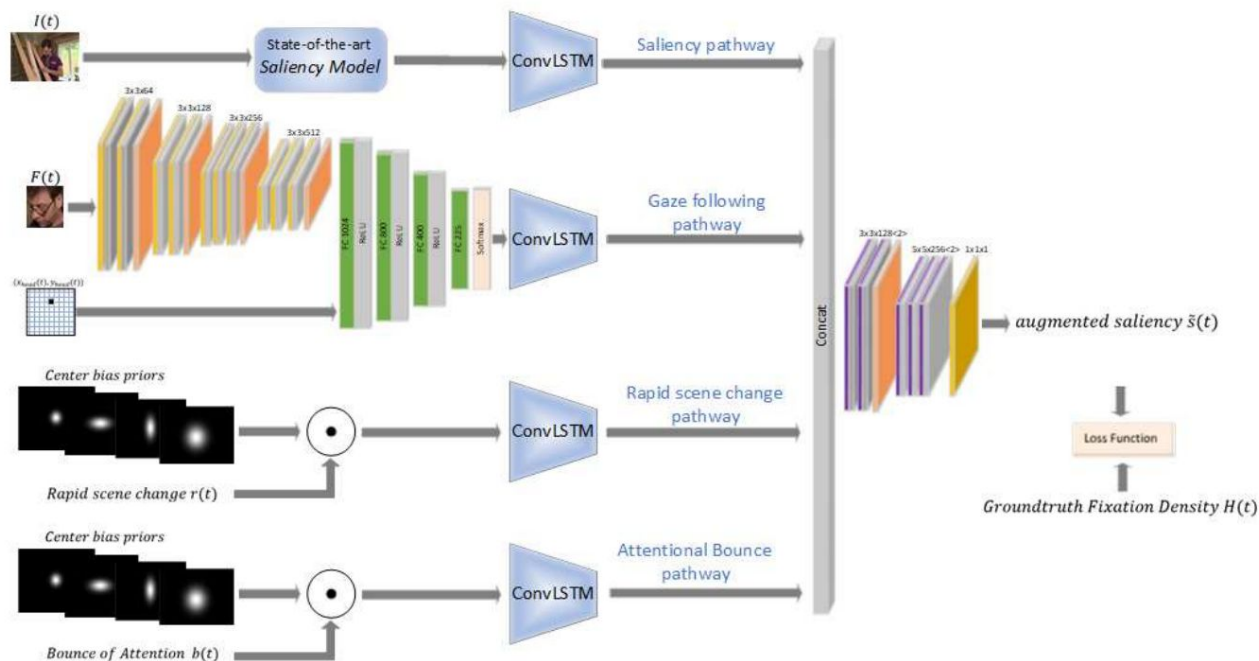
(b)



# From image to video saliency?



# From image to video saliency



Questions?