# Task Specific Visual Saliency Prediction with Memory Augmented Conditional Generative Adversarial Networks

Tharindu Fernando      Simon Denman      Sridha Sridharan      Clinton Fookes

Image and Video Research Laboratory, Queensland University of Technology (QUT), Australia
{t.warnakulasuriya, s.denman, s.sridharan, c.fookes}@qut.edu.au

## Abstract

*Visual saliency patterns are the result of a variety of factors aside from the image being parsed, however existing approaches have ignored these. To address this limitation, we propose a novel saliency estimation model which leverages the semantic modelling power of conditional generative adversarial networks together with memory architectures which capture the subject's behavioural patterns and task dependent factors. We make contributions aiming to bridge the gap between bottom-up feature learning capabilities in modern deep learning architectures and traditional top-down hand-crafted features based methods for task specific saliency modelling. The conditional nature of the proposed framework enables us to learn contextual semantics and relationships among different tasks together, instead of learning them separately for each task. Our studies not only shed light on a novel application area for generative adversarial networks, but also emphasise the importance of task specific saliency modelling and demonstrate the plausibility of fully capturing this context via an augmented memory architecture.*

## 1. Introduction

Visual saliency patterns are the result of a number of factors, including the task being performed, user preferences and domain knowledge. However existing approaches to predict saliency patterns [1–6] ignore these factors, and instead learn a model specific to a single task while disregarding factors such as user preferences.

Based on empirical results, the human visual system is driven by both bottom-up and top-down factors [7]. The first category (bottom-up) is entirely driven by the visual scene where humans deploy their attention towards salient, informative regions such as bright colours, unique textures or sudden movements.

The bottom-up factors are typically exhibited during the free viewing mechanism. In contrast, the top-down attention component, where the observer is performing a task specific search, is modulated by the task at hand [8] and the subject's prior knowledge [9]. For example, when the observer is searching for people in the scene, they can selectively attend to the scene regions which are most likely to contain the targets [10, 11]. Furthermore, a subject's prior knowledge, such as the scene layout, scene categories and statistical regularities will influence the search mechanisms and the fixations [12–14], rendering task specific visual saliency prediction a highly subjective, situational and challenging task [8, 9, 15], which motivates the need for a working memory [9, 16, 17]. Even within groups of subjects completing the same task, due to the differences in a subject's behavioural goals, expectations and domain knowledge, unique saliency patterns are generated. Ideally, this user related context information should be captured via a working memory.

Fig. 1 shows the variability of the saliency maps when observers are performing action recognition and context recognition on the same image. In Fig. 1 (a) the observer is asked to recognise the action performed by the human in the scene; while in Fig. 1 (b) the saliency map is generated when the observer is searching for cars/ trees in the scene. It is evident that there exists variability in the resultant saliency patterns, yet accurate modelling of human fixations in the application areas specified above requires task specific models. For example, semantic video search, content aware image resizing, video surveillance, and video/ scene classification may require a search for pedestrians, for different objects, or recognising human actions depending on the task and video context.

In recent years, motivated by the success of deep learning techniques [18, 19], there have been several attempts to model visual saliency of the human free viewing mechanism with the aid of deep convolutional

(a) Human action recogni-
tion task

(b) Searching for cars/
trees

Figure 1: Variability of the saliency maps when ob-
servers are performing different tasks



(a) Conditional GAN

(b) MC-GAN (proposed
model)

Figure 2: A comparison of conditional GAN architec-
ture with the proposed model

networks [9, 20–22]. Yet, the usual approach when
modelling visual saliency for task specific viewing is to
hand-craft the visual features. For instance, in [8] the
authors utilise the features from person detectors [23]
when estimating the search for humans in the scene;
while in [15] the authors use the features from HoG
descriptors [24] when searching for the objects in the
scene. Therefore, these approaches are application de-
pendent and fail to capture the top-down attentional
mechanism of humans which is driven by factors such
as a subject's prior knowledge and expectation.

In this work we propose a deep learning architecture
for task specific visual saliency estimation. We draw
our inspiration from the recent success of Generative
Adversarial Networks (GAN) [25–29] for pixel to pixel
translation tasks. We exploit the capability of the con-
ditional GAN framework [30] for automatic learning of
task specific features in contrast to hand-crafted fea-
tures that are tailored for specific applications [8, 15].
This results in a unified, simpler architecture enabling
direct application to a variety of tasks. The condi-
tional nature of the proposed architecture enables us to
learn one network for all the tasks of interest, instead of
learning separate networks for each of the tasks. Apart
from the advantage of a simpler learning process, this
enables the capability of learning semantic correspon-
dences among different tasks and propagating these
contextual relationships from one task to another.

Fig. 2 (a) shows the conditional GAN architecture
where the discriminator $D$ learns to classify between
real and synthesised pairs of saliency maps $y$, given the
observed image $x$ and task specific class label $c$. The
generator $G$ tries to fool the discriminator. It also ob-
serves the observed image $x$ and task specific class label
$c$. We compare this model to the proposed model given
in Fig. 2 (b). The differences arise in the utilisation
of memory $M$, where we capture subject specific be-
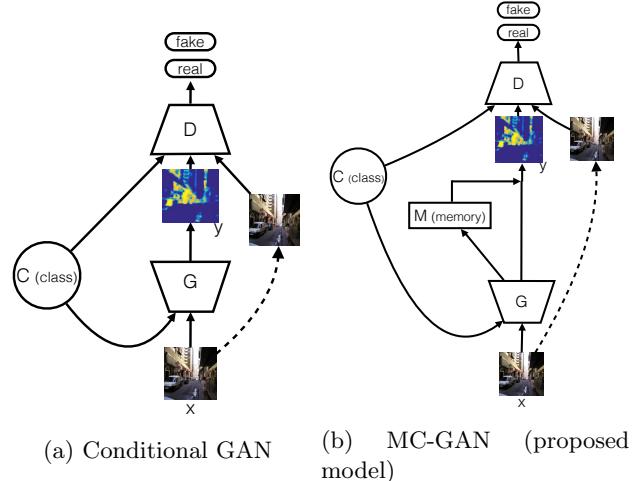havioural patterns. This incorporates a subject's prior

knowledge, behavioural goals and expectations.

## 2. Related Work

Literature related to this work can be broadly cate-
gorised into "Visual Saliency Prediction" and "Gener-
ative Adversarial Networks", and these two areas are
addressed in Sections 2.1 and 2.2 respectively.

### 2.1. Visual Saliency Prediction

Since the first attempts to model human saliency
through feature integration [31], the area of saliency
prediction has been widely explored. Building upon
this bottom-up approach, Koch and Ullman [32] and
Itti et al. [33] proposed approaches based on extract-
ing image features such as colour, intensity and orienta-
tion. These methods generate centre-biased acceptable
saliency predictions for free viewing but are highly in-
accurate in complex real world scenes [9]. Recent stud-
ies such as [34–38] have looked into the development of
more complex features for saliency estimation.

In contrast, motivated by information theory, au-
thors in [8, 15, 39] have taken the top-down approach
where task dependent features comes into play. They
incorporate local information from regions of interest
for the task at hand, such as features from person de-
tectors and HoG descriptors. These models [8, 15, 39]
are completely task specific, rendering adaptation from
one task to another nearly impossible. Furthermore,
they neglect the fact that different subjects may exhibit
different behavioural patterns when achieving the same
goal which generates unique strategies or sub goals that
we term user preferences.

In order to exploit the representative power of deep
architectures, more recent studies have been driven to-

wards the utilisation of convolution neural networks. In contrast to the above approaches, which use hand crafted features, deep learning based approaches offer automatic feature learning. In [22] the authors propose the usage of feature representations from a pre-trained model that has been trained for object classification. This work was followed by [9, 20] where authors train end-to-end saliency prediction models from scratch and their experimental evaluations suggest that deep models trained for saliency prediction itself can outperform off-the-shelf CNN models.

Liu et al. [40] proposed a mulitresolution-CNN for predicting saliency, which has been trained on multiple scales of the observed image. The motivation behind this approach is to capture low and high level features. Yet this design has an inherit deficiency due to the use of isolated image patches which fail to capture the global context, composed of the context of the observed image, the task at hand (i.e free viewing, recognising actions, searching for objects) and user preferences. Even though the context of the observed image is well represented in deep single scale architectures such as [9, 20] they ignore the rest of the global context, the task description and user behavioural patterns, which are often crucial for saliency estimation.

The proposed work bridges the research gap between deep architectures that capture bottom-up saliency features; and top-down methods [8, 15, 39] that are purely driven by the hand crafted features. We investigate the plausibility of the complete automatic learning of global context, which has been ill represented in literature thus far, through a memory augmented conditional generative adversarial model.

### 2.2. Generative Adversarial Networks

Generative adversarial networks (GAN), which belong to the family of generative models, have achieved promising results for pixel-to-pixel synthesis [41]. Several works have looked into numerous architectural augmentations when synthesising natural images. For instance, in [42] the authors utilise a recurrent network approach where as in [43] a de-convolution network approach is used to generate higher quality images. Most recently authors in [44] have utilised the GAN architecture for visual saliency prediction and proposed a novel loss function which is proven to be effective for both initialising the generator, and stabilising adversarial training. Yet their work fails to delineate the ways of achieving task specific saliency estimation and of incorporating task specific dependencies and the subject behavioural patterns for saliency estimation.

The proposed work draws inspiration from conditional GANs [45–52]. This architecture is extensively applied for image based prediction problems such as image prediction from normal maps [52], future frame prediction in videos [46], image style transfer [45], image manipulation guided by user preferences [51], etc. In [30] the authors proposed a novel U-Net [53] based architecture for conditional GANs. Their evaluations suggested that this network is capable of capturing local semantics with applications to a wide range of problems. We investigate the possibility of merging the discriminative learning power of conditional GANs together with a local memory to fully capture the global context, contributing a novel application area and structural argumentation for conditional GANs.

## 3. Visual Saliency Model

### 3.1. Objectives

Generative adversarial networks learn a mapping from a random noise vector $z$ to an output image $y$, $G : z \to y$ [25]; where as conditional GANs learn a mapping from an observed image $x$ and random noise vector $z$, to output $y$, given auxiliary information $c$, where $c$ can be class labels or data from other modalities. $G : \{x, z | c\} \to y$. When we incorporate the notion of time into the system, then the observed image at time instance $t$ will be $x_t$, the respective noise vector $z_t$ and the relevant class label will be $c_t$. Then the objective function of a conditional GAN can be written as,

$$L_{cGAN}(G, D) = E_{x_t, y_t \sim p_{data}(x_t, y_t)}[log(D(x_t, y_t | c_t))] + E_{x_t \sim p_{data}(x_t), z_t \sim p_z(z_t)}[log(1 - D(x_t, G(x_t, z_t | c_t)))]. \tag{1}$$

Let $M \in \mathbb{R}^{k*l}$, shown in Fig. 3, be the working memory with $k$ memory slots and $l$ is the embedding dimension of the generator output,

$$o_t = G(x_t, z_t | c_t). \tag{2}$$

If the representation of memory at time instance $t-1$ is given by $M_{t-1}$ and $f_r^{LSTM}$ is a read function, then we can generate a key vector $a_t$, representing the similarity between the current memory content and the current generator embedding via attending over the memory slots such that,

$$\grave{o}_t = f_r^{LSTM}(o_t), \tag{3}$$

$$a_t = softmax(\grave{o}_t^T, M_{t-1}), \tag{4}$$

and

$$h_t = a_t^T M_{t-1}. \tag{5}$$

Then we retrieve the current memory state by,

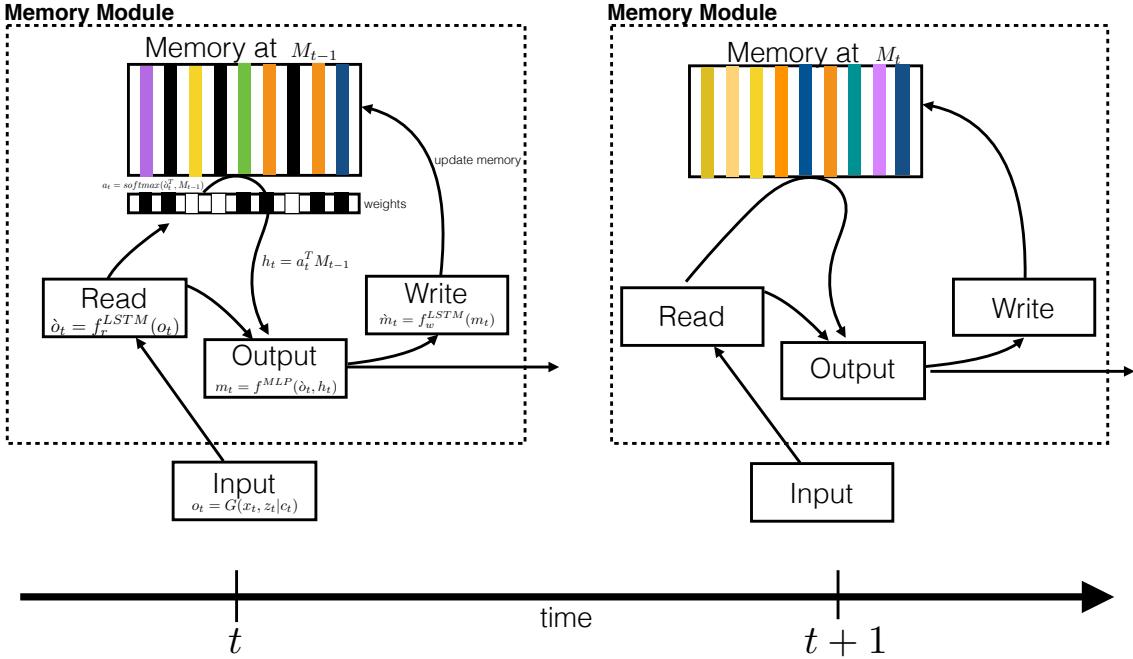$$m_t = f^{MLP}(\grave{o}_t, h_t), \tag{6}$$

Figure 3: Memory architecture: The model receives an input $o_t$ at time instance $t$. A function $f_r^{LSTM}$ is used to embed this input and retrieve the content of the memory $M_{t-1}$. When reading, we use a $softmax$ function to weight the association between each memory slot and the input $o_t$, deriving a weighted retrieval $h_t$. The final output $m_t$ is derived using both $\grave{o}$ and $h_t$. Finally we update the memory using memory write function $f_w^{LSTM}$. This generates the memory representation $M_t$ at time instance $t+1$, shown to the right.

where $f^{MLP}$ is a neural network composed of multi-layer perceptrons (MPL) trained jointly with other networks. Then we generate the vector for the memory update $\dot{m}_t$ via passing it through a write function $f_w^{LSTM}$

$$\dot{m}_t = f_w^{LSTM}(m_t), \qquad (7)$$

and finally we completely update the memory using,

$$M_t = M_{t-1}(1 - (a_t \otimes e_k)^T) + (\dot{m}_t \otimes e_l)(a_t \otimes e_k)^T. \quad (8)$$

where 1 is a matrix of ones, $e_l \in \mathbb{R}^l$ and $e_k \in \mathbb{R}^k$ be vectors of ones and $\otimes$ denotes the outer product which duplicates its left vector $l$ or $k$ times to form a matrix. Now the objective of the proposed memory augmented cGAN can be written as,

$$L_{cGAN}^*(G, D) = E_{x_t, y_t \sim p_{data}(x_t, y_t)}[log(D(x_t, y_t | c_t))] +$$
$$E_{x_t \sim p_{data}(x_t), z_t \sim p_z(z_t)}[log(1 - D(x_t, o_t \otimes tanh(m_t)))]. \quad (9)$$

We would like to emphasise that we are learning a single network for all the tasks at hand, rendering a simpler but informative framework, which can be directly applied to a variety of tasks without any fine tuning.
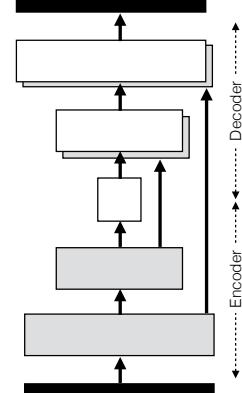


Figure 4: U-Net architecture

### 3.2. Network Architecture

For the generator we adapt the U-Net architecture of [30] (see Fig. 4). Let $Ck$ denote a Convolution-BatchNorm-ReLU layer group with $k$ filters. $CDk$ denotes a Convolution-BatchNorm-Dropout-ReLU layer with a dropout rate of 50%. Then the generator architecture can be written as, Encoder: C64-C128-C256-C512-C512-C512-C512-C512 followed by a U-Net

decoder: CD512-CD1024-CD1024-C1024-C1024-C512-C256-C128 where there are skip connections between each $i^{th}$ layer in the encoder and the $n - i^{th}$ layer of the decoder, and there are $n$ total layers in the generator (see [30] for details). The discriminator architecture is: C64-C128-C256-C512-C512-C512.

All convolutions are 4 x 4 spatial filters applied with stride 2. Convolutions in the encoder and in the discriminator down sample by a factor of 2, whereas in the decoder they up sample by a factor of 2. A description of our memory architecture is as follows. For functions $f_r^{LSTM}$ and $f_w^{LSTM}$ we utilise two, one-layer LSTM networks [54] with 100 hidden units and for $f^{MLP}$ we use a neural network with a single hidden layer and 1024 units with ReLU activations. This memory module is fully differentiable, and we learn it jointly with other networks. We trained the proposed model with the Adam [55] optimiser, with a batch size of 32 and an initial learning rate to 1e-5, for 10 epochs.

## 4. Experimental Evaluations

### 4.1. Datasets

We evaluate our proposed approach on 2 publicly available datasets, VOCA-2012 [15] and MIT person detection (MIT-PD) [8].

The VOCA-2012 dataset consists of 1,085,381 human eye fixation from 12 volunteers (5 male and 7 female) divided into 2 groups based on the given task. It contains 8 subjects performing action recognition in the given image where as the rest of the subjects are performing context dependent visual search. The subjects in this group are searching for furniture, paintings/ wallpapers, bodies of water, buildings, cars/ trucks, mountain/ hills, road/ trees in the given scene. The MIT-PD dataset consist of 12,768 fixations from 14 subjects (between 18-40 years), where the subjects search for people in 912 urban scenes. MIT-PD contains only a single task, and we use this dataset to demonstrate the effectiveness of the memory network.

### 4.2. Evaluation metrics

Let $N$ denote the number of examples in the testing set, $y$ denotes the ground truth saliency map and $\hat{y}$ is the predicted saliency map. Following this notation we define the following metrics:

- **Area Under Curve (AUC):** AUC is a widely used metrics for evaluating saliency models. We use the formulation of this metric defined in [56].

- **Normalised Scan path Saliency (NSS):** NSS [57] is calculated by taking the mean scores assigned by the unit normalised saliency map $\hat{y}^{norm}$

at human eye fixations.

$$NSS = \frac{1}{N} \sum_{i=1}^{N} \hat{y}_i^{norm} \qquad (10)$$

- **Linear Correlation Coefficient (CC):** In order to measure the linear relationship between the ground truth and predicted map we utilise the linear correlation coefficient,

$$CC = \frac{cov(y, \hat{y})}{\sigma_y * \sigma_{\hat{y}}}, \qquad (11)$$

where $cov(y, \hat{y})$ is referred to as the covariance between distributions $y$ and $\hat{y}$, $\sigma_y$ is the standard deviation of distributions $y$ and $\sigma_{\hat{y}}$ is the standard deviation of distributions $\hat{y}$. As the name implies $CC = 1$ denotes a perfect linear relationship between distributions $y$ and $\hat{y}$ where as a value of 0 implies that there is no linear relationship.

- **KL divergence (KL):** To measure the non-symmetric difference between two distributions we utilise the KL divergence measure given by,

$$KL = \sum_{i=1}^{N} \hat{y}_i log(\frac{\hat{y}_i}{y_i}). \qquad (12)$$

As ground truth and predicted saliency maps can be seen as 2D distributions, we can use KL divergence to measure the difference between them.

- **Similarity metric (SM):** Computes the sum of the minimum values at each pixel location between $\hat{y}^{norm}$ and $y^{norm}$ distributions.

$$SM = \sum_{i=1}^{P} min(\hat{y}_i^{norm}, y_i^{norm}), \qquad (13)$$

where $\sum_{i=1}^{P} \hat{y}_i^{norm} = 1$ and $\sum_{i=1}^{P} y_i^{norm} = 1$ are the normalised probability distributions and $P$ denotes all the pixel location in the 2D maps.

### 4.3. Results

Quantitative evaluations on the VOCA-2012 dataset is presented in Tab. 1. In the proposed model, in order to retain the user dependent factors such as user preference in memory, we feed the examples in order such that examples from each specific user go through in sequence. We compare our model with 8 state-of-the-art methods. The row 'human' stands for the human saliency predictor, which computes the saliency map derived from the fixations made by half of the human subjects performing the same task. This predictor is

evaluated with respect to the rest of the subjects, as opposed to the entire group [15].

The evaluations suggest that the bottom-up model of Itti et. al [58] generates poor results as it does not incorporate task specific information. Even with high level object detectors, the models of Judd et al. [59] and HOG detector [15] fail to render accurate predictions.

Deep learning models PDP [60] and ML-net [20] are able to out perform the techniques stated above but they lack the ability to learn task dependent information. We note the accuracy gain of cGAN model over PDP, ML-net and SalGAN, where the model incorporates a conditional variable to discriminate between the 'action recognition' and 'context recognition' tasks instead of learning two separate networks or fine-tuning on them individually. Our proposed approach builds upon this by incorporating an augmented memory architecture with conditional learning. We learn different user patterns and retain the dependencies among different tasks, and outperform all baselines considered (MC-GAN (proposed), Tab. 1).

As further study, in Tab. 1, row M-GAN (separate), we show the evaluations for training two separate memory augmented GAN networks for the tasks in the VOCA-2012 test set without using the conditional learning process. The results emphasise the importance of learning a single network for all the tasks, leveraging semantic relationships between different tasks. The accuracy of the networks learned for separate tasks are lower than the combined MC-GAN and cGAN approaches (rows MC-GAN (proposed) and cGAN, Tab. 1), highlighting the importance of learning the different tasks together and allowing the model to discriminate between the tasks and learn the complimentary information, rather than keeping the model completely blind regarding the existence of another task category.

To provide qualitative insight, some predicted maps along with ground truth and baseline ML-net [20] predictions are given in Fig. 5. In the first column we show the input image, and columns "Action rec GT" and "Context rec GT" depict the ground truth saliency maps for the respective tasks. In columns "Our action rec" and "Our context rec" we show the respective predictions from our model, and finally the column 'ml-Net' contains the prediction from the ML-net [20] baseline. Observing columns "Action rec GT " and "Context rec GT" one can clearly see how the tasks differ based on the different saliency patterns. Yet, the proposed model is able to capture these different semantics within a single network which is trained together for all the tasks. As shown in Fig. 5, it has efficiently identified the image saliency from low level features as

| Saliency Models | Task | | | |
| | Action Rec | | Context Rec | |
| | AUC | KL | AUC | KL |
|---|---|---|---|---|
| HOG detector [15] | 0.736 | 8.54 | 0.646 | 8.10 |
| Judd et al. [59] | 0.715 | 11.00 | 0.636 | 9.66 |
| Itti et. al [58] | 0.533 | 16.53 | 0.512 | 15.04 |
| central bias [15] | 0.780 | 9.59 | 0.685 | 8.82 |
| PDP [60] | 0.875 | 8.23 | 0.690 | 7.98 |
| ML-net [20] | 0.847 | 8.51 | 0.684 | 8.02 |
| SalGAN [44] | 0.848 | 8.47 | 0.679 | 8.00 |
| cGAN [30] | 0.852 | 8.24 | 0.701 | 7.95 |
| M-GAN (separate) | 0.848 | 8.54 | 0.704 | 8.00 |
| **MC-GAN (proposed)** | **0.901** | **8.07** | **0.734** | **7.65** |
| Human [15] | 0.922 | 6.14 | 0.813 | 5.90 |

Table 1: Experimental evaluation on VOCA-2012 test set. We augment the current state-of-the-art GAN method (SalGAN [44]) by adding a conditional variable (cGAN [30]) to mimic the joint learning process instead of learning two separate networks. To capture user and task specific behavioural patterns we add a memory module to cGAN, MC-GAN (proposed), and outperform all the baselines. We also compare training 2 separate memory augmented GAN networks, M-GAN (separate) without the conditional learning process.

| Saliency Models | Task | | | | | |
| | Action Rec | | | Context Rec | | |
| | NSS | CC | SM | NSS | CC | SM |
|---|---|---|---|---|---|---|
| ML-net [20] | 2.05 | 0.71 | 0.51 | 2.03 | 0.64 | 0.42 |
| SalGAN [44] | 2.10 | 0.73 | 0.51 | 2.10 | 0.68 | 0.44 |
| cGAN [30] | 2.23 | 0.76 | 0.55 | 2.14 | 0.71 | 0.57 |
| **MC-GAN (proposed)** | **2.23** | **0.79** | **0.60** | **2.20** | **0.77** | **0.69** |

Table 2: Comparison between ML-Net, SalGAN, cGAN and MC-GAN (proposed) on VOCA-2012

well as task dependent saliency factors from high level cues such as trees, furniture and roads. Furthermore, the single learning process and the incorporation of a memory architecture renders the plausibility of retaining the semantical relationships among different tasks and how users adapt to those.

Tab. 3 shows the performance of the proposed model with 5 baselines for the MIT-PD test set. The first baseline "Scene Context" [8] utilises colour and orientation features where as "Combined" [8] incorporates both scene context features and higher level features from a person detector [23]. Even with such explicit modelling of the task, this baseline fails to generate accurate predictions suggesting the subjective nature of the task specific viewing. With the aid of the associative memory of the proposed model we successfully capture those underlying factors.

In Tab. 2 and Tab. 4 we present the evaluations of NSS, CC and SM metrics. In order to evaluate ML-
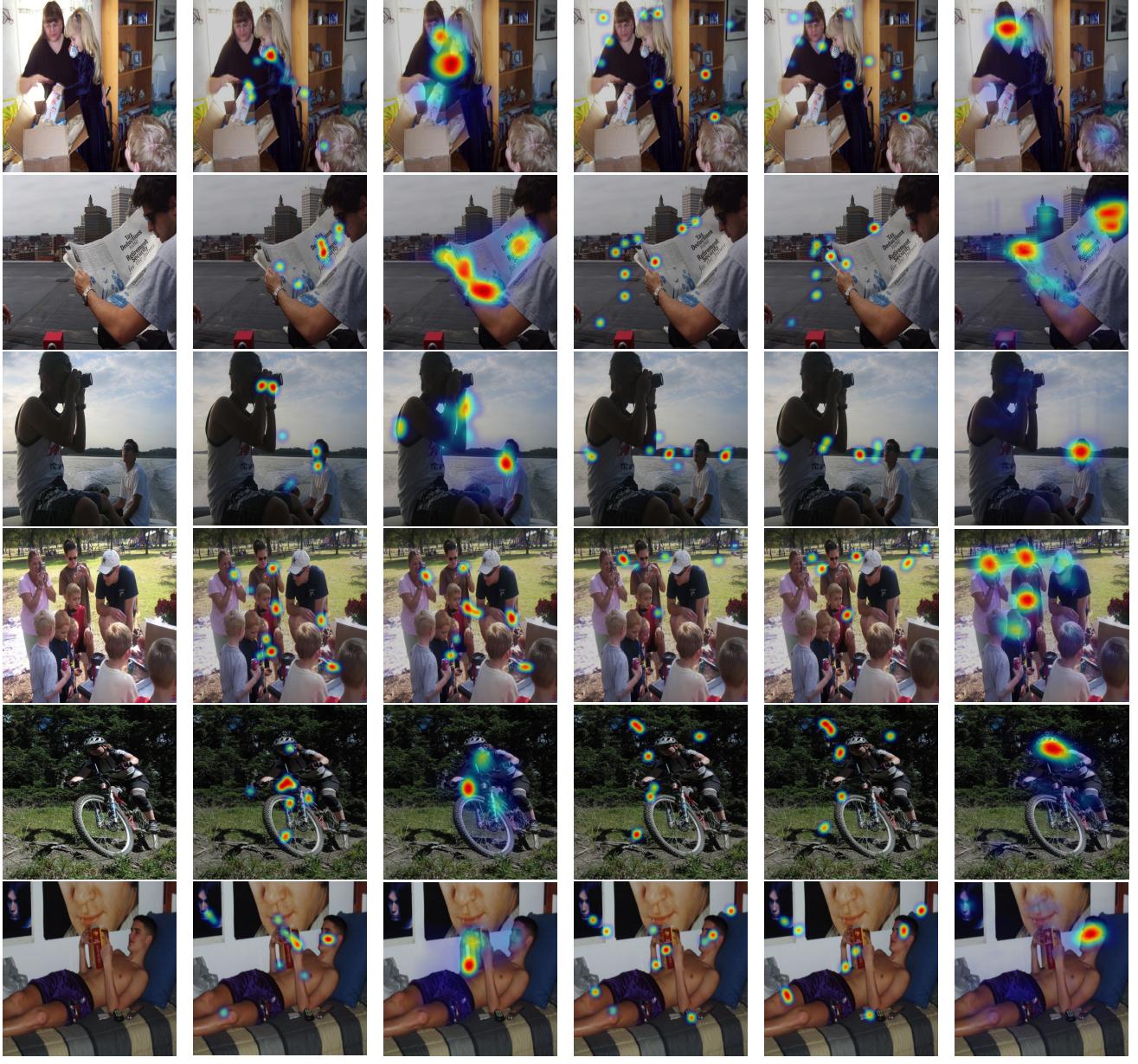
|            |                |                   |                  |                     |            |
| :--------: | :------------: | :---------------: | :--------------: | :-----------------: | :--------: |
| (a) Image  | (b) Action rec GT | (c) Our action rec | (d) Context rec GT | (e) Our context rec | (f) ML-net |

Figure 5: Qualitative results for VOCA-2012 dataset and comparisons to the state-of-the-art.

net, SalGAN and cGAN models we utilise the implementation of the algorithm released by the authors. When comparing the results between the ML-net [20], cGAN [30] and SalGAN [44] models and proposed MC-GAN model a considerable gain in performance is observed in all the metrics considered, which emphasises a greater agreement between predicted and ground truth saliency maps. We were unable to compare other baselines using these metrics due to the unavailability of public implementations.

The qualitative results obtained from the proposed

model along with the ML-net [20] network on a few examples from the MIT-PD dataset are shown in Fig. 6. We would like to emphasise the usage of a subject's prior knowledge in the task of searching for people in the urban scene. The subjects selectively attend the areas such as high rise buildings (see rows 2, 6) and pedestrian walkways (see rows 1, 3-6), which are more likely to contain humans, which our model has effectively captured. With the lack of capacity to model such user preferences, the baseline ML-Net model generates centre biased saliency without effectively under-

| Saliency Models | Scene Type | |
| --- | --- | --- |
| | Target Present | Target absent |
| | AUC | AUC |
| Scene Context [8] | 0.844 | 0.845 |
| Combined [8] | 0.896 | 0.877 |
| ML-net [20] | 0.901 | 0.881 |
| SalGAN [44] | 0.910 | 0.887 |
| cGAN [30] | 0.923 | 0.899 |
| **MC-GAN (proposed)** | **0.942** | **0.903** |
| Human [8] | 0.955 | 0.930 |

Table 3: Experimental evaluation on MIT-PD test set

| Saliency Models | Task | | | | | |
| --- | --- | --- | --- | --- | --- | --- |
| | Target Present | | | Target absent | | |
| | NSS | CC | SM | NSS | CC | SM |
| ML-Net [20] | 1.41 | 0.55 | 0.41 | 1.22 | 0.43 | 0.38 |
| SalGAN [44] | 1.41 | 0.53 | 0.44 | 1.20 | 0.42 | 0.35 |
| cGAN [30] | 1.67 | 0.51 | 0.59 | 2.02 | 0.41 | 0.52 |
| **MC-GAN (proposed)** | **2.17** | **0.76** | **0.71** | **2.34** | **0.75** | **0.78** |

Table 4: Comparison between ML-Net, SalGAN, cGAN and MC-GAN (proposed) on MIT-PD test set

standing the subject's strategy.

### 4.4. Task Specific Generator Activations

In Fig. 7 we visualise the activations from the 2nd (conv-l-2) and 2nd last (conv-l-8) convolution layers of the generator. The task specific learning of the proposed conditional GAN architecture is clearly evident in the activations. For instance, when the task at hand is to recognise actions the generator activations are highly concentrated around the foreground of the image (see (b), (g)), while for context recognition the model has learned that the areas of interest are in the background of the image (see (c), (h)). These task specific salient features are combined and compressed hierarchically and in latter layers (i.e conv-l-8), the networks has learned the most specific areas to focus when generating the output saliency map.

## 5. Conclusion

This work introduces a novel human saliency estimation architecture which combines task and user specific information together in a generative adversarial pipeline. We show the importance of fully capturing the context information which incorporates the task information, subject behavioural goals and image context. The resultant frame work offers several advantages compared to task specific handcrafted features, enabling direct transferability among different tasks. Qualitative and quantitative experimental evaluations on two public datasets demonstrates superior performance with respect to the current state-of-the-art.

(a) Image    (b) GT    (c) Our    (d) ML-net

Figure 6: Qualitative results for MIT-PD dataset and comparisons to the state-of-the-art.



(a) Input (b) l-2 AR (c) l-2 CR (d) l-8 AR (e) l-8 CR

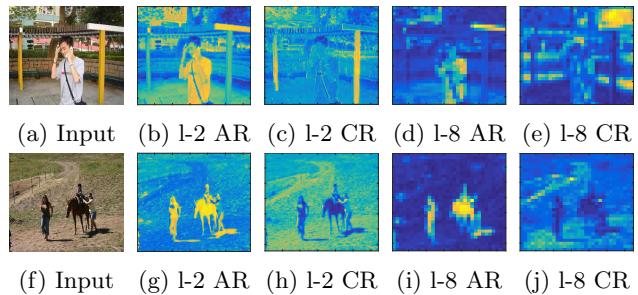(f) Input (g) l-2 AR (h) l-2 CR (i) l-8 AR (j) l-8 CR

Figure 7: Visualisation of the generator activations from 2nd (l-2) and 2nd last (l-8) convolution layers for action recognition (AR) and context recognition (CR) tasks. The importance varies from blue to yellow where blue represents the areas of least importance and yellow represents areas of more importance.

# References

[1] D. Walther, L. Itti, M. Riesenhuber, T. Poggio, and C. Koch, "Attentional selection for object recognitiona gentle way," in *International Workshop on Biologically Motivated Computer Vision*. Springer, 2002, pp. 472–479.

[2] H. Hadizadeh and I. V. Bajic, "Saliency-aware video compression," *IEEE Transactions on Image Processing*, vol. 23, no. 1, pp. 19–33, 2014.

[3] W. Wang, J. Shen, and F. Porikli, "Saliency-aware geodesic video object segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 3395–3402.

[4] R. Achanta and S. Süsstrunk, "Saliency detection for content-aware image resizing," in *Image Processing (ICIP), 2009 16th IEEE International Conference on*. IEEE, 2009, pp. 1005–1008.

[5] G. Sharma, F. Jurie, and C. Schmid, "Discriminative spatial saliency for image classification," in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on*. IEEE, 2012, pp. 3506–3513.

[6] T. Yubing, F. A. Cheikh, F. F. E. Guraya, H. Konik, and A. Trémeau, "A spatiotemporal saliency model for video surveillance," *Cognitive Computation*, vol. 3, no. 1, pp. 241–263, 2011.

[7] C. E. Connor, H. E. Egeth, and S. Yantis, "Visual attention: bottom-up versus top-down," *Current Biology*, vol. 14, no. 19, pp. R850–R852, 2004.

[8] K. A. Ehinger, B. Hidalgo-Sotelo, A. Torralba, and A. Oliva, "Modelling search for people in 900 scenes: A combined source model of eye guidance," *Visual cognition*, vol. 17, no. 6-7, pp. 945–978, 2009.

[9] S. S. Kruthiventi, K. Ayush, and R. V. Babu, "Deepfix: A fully convolutional neural network for predicting human eye fixations," *arXiv preprint arXiv:1510.02927*, 2015.

[10] W. Einhã, U. Rutishauser, C. Koch *et al.*, "Task-demands can immediately reverse the effects of sensory-driven saliency in complex visual stimuli," *Journal of vision*, vol. 8, no. 2, pp. 2–2, 2008.

[11] G. J. Zelinsky, "A theory of eye movements during target acquisition." *Psychological review*, vol. 115, no. 4, p. 787, 2008.

[12] M. S. Castelhano and J. M. Henderson, "Initial scene representations facilitate eye movement guidance in visual search." *Journal of Experimental Psychology: Human Perception and Performance*, vol. 33, no. 4, p. 753, 2007.

[13] M. Chaumon, V. Drouet, and C. Tallon-Baudry, "Unconscious associative memory affects visual processing before 100 ms," *Journal of vision*, vol. 8, no. 3, pp. 10–10, 2008.

[14] M. B. Neider and G. J. Zelinsky, "Scene context guides eye movements during visual search," *Vision research*, vol. 46, no. 5, pp. 614–621, 2006.

[15] S. Mathe and C. Sminchisescu, "Action from still image dataset and inverse optimal control to learn task specific visual scanpaths," in *Advances in neural information processing systems*, 2013, pp. 1923–1931.

[16] T. Fernando, S. Denman, A. McFadyen, S. Sridharan, and C. Fookes, "Tree memory networks for modelling long-term temporal dependencies," *arXiv preprint arXiv:1703.04706*, 2017.

[17] T. Fernando, S. Denman, S. Sridharan, and C. Fookes, "Going deeper: Autonomous steering with neural memory networks," in *The IEEE International Conference on Computer Vision (ICCV)*, Oct 2017.

[18] H. Gammulle, S. Denman, S. Sridharan, and C. Fookes, "Two stream lstm: A deep fusion framework for human action recognition," in *Applications of Computer Vision (WACV), 2017 IEEE Winter Conference on*. IEEE, 2017, pp. 177–186.

[19] T. Fernando, S. Denman, S. Sridharan, and C. Fookes, "Soft+ hardwired attention: An lstm framework for human trajectory prediction and abnormal event detection," *arXiv preprint arXiv:1702.05552*, 2017.

[20] M. Cornia, L. Baraldi, G. Serra, and R. Cucchiara, "A deep multi-level network for saliency prediction," *arXiv preprint arXiv:1609.01064*, 2016.

[21] E. Vig, M. Dorr, and D. Cox, "Large-scale optimization of hierarchical features for saliency prediction in natural images," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 2798–2805.

[22] M. Kümmerer, L. Theis, and M. Bethge, "Deep gaze i: Boosting saliency prediction with feature maps trained on imagenet," *arXiv preprint arXiv:1411.1045*, 2014.

[23] N. Dalal, B. Triggs, and C. Schmid, "Human detection using oriented histograms of flow and appearance," in *European conference on computer vision*. Springer, 2006, pp. 428–441.

[24] N. Dalal and B. Triggs, "Histograms of oriented gradients for human detection," in *Computer Vision and Pattern Recognition, 2005. CVPR 2005. IEEE Computer Society Conference on*, vol. 1. IEEE, 2005, pp. 886–893.

[25] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in neural information processing systems*, 2014, pp. 2672–2680.

[26] E. L. Denton, S. Chintala, R. Fergus *et al.*, "Deep generative image models using a laplacian pyramid of adversarial networks," in *Advances in neural information processing systems*, 2015, pp. 1486–1494.

[27] A. Radford, L. Metz, and S. Chintala, "Unsupervised representation learning with deep convolutional generative adversarial networks," *arXiv preprint arXiv:1511.06434*, 2015.

[28] T. Salimans, I. Goodfellow, W. Zaremba, V. Cheung, A. Radford, and X. Chen, "Improved techniques for training gans," in *Advances in Neural Information Processing Systems*, 2016, pp. 2226–2234.

[29] J. Zhao, M. Mathieu, and Y. LeCun, "Energy-based generative adversarial network," *arXiv preprint arXiv:1609.03126*, 2016.

[30] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros, "Image-to-image translation with conditional adversarial networks," *arXiv preprint arXiv:1611.07004*, 2016.

[31] A. M. Treisman and G. Gelade, "A feature-integration theory of attention," *Cognitive psychology*, vol. 12, no. 1, pp. 97–136, 1980.

[32] C. Koch and S. Ullman, "Shifts in selective visual attention: towards the underlying neural circuitry," in *Matters of intelligence*. Springer, 1987, pp. 115–141.

[33] L. Itti, C. Koch, and E. Niebur, "A model of saliency-based visual attention for rapid scene analysis," *IEEE Transactions on pattern analysis and machine intelligence*, vol. 20, no. 11, pp. 1254–1259, 1998.

[34] R. Valenti, N. Sebe, and T. Gevers, "Image saliency by isocentric curvedness and color," in *Computer Vision, 2009 IEEE 12th International Conference on.* IEEE, 2009, pp. 2185–2192.

[35] Z. Liu, O. Le Meur, S. Luo, and L. Shen, "Saliency detection using regional histograms," *Optics letters*, vol. 38, no. 5, pp. 700–702, 2013.

[36] C. Lang, T. Nguyen, H. Katti, K. Yadati, M. Kankanhalli, and S. Yan, "Depth matters: Influence of depth cues on visual saliency," *Computer Vision–ECCV 2012*, pp. 101–115, 2012.

[37] E. Erdem and A. Erdem, "Visual saliency estimation by nonlinearly integrating features using region covariances," *Journal of vision*, vol. 13, no. 4, pp. 11–11, 2013.

[38] F. Porikli, O. Tuzel, and P. Meer, "Covariance tracking using model update based on lie algebra," in *Computer Vision and Pattern Recognition, 2006 IEEE Computer Society Conference on*, vol. 1. IEEE, 2006, pp. 728–735.

[39] N. Bruce and J. Tsotsos, "Saliency based on information maximization," *Advances in neural information processing systems*, vol. 18, p. 155, 2006.

[40] N. Liu, J. Han, D. Zhang, S. Wen, and T. Liu, "Predicting eye fixations using convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 362–370.

[41] T. Arici and A. Celikyilmaz, "Associative adversarial networks," *arXiv preprint arXiv:1611.06953*, 2016.

[42] K. Gregor, I. Danihelka, A. Graves, D. J. Rezende, and D. Wierstra, "Draw: A recurrent neural network for image generation," *arXiv preprint arXiv:1502.04623*, 2015.

[43] A. Dosovitskiy, J. Tobias Springenberg, and T. Brox, "Learning to generate chairs with convolutional neural networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1538–1546.

[44] J. Pan, C. Canton, K. McGuinness, N. E. O'Connor, J. Torres, E. Sayrol, and X. Giro-i Nieto, "Salgan: Visual saliency prediction with generative adversarial networks," *arXiv preprint arXiv:1701.01081*, 2017.

[45] C. Li and M. Wand, "Precomputed real-time texture synthesis with markovian generative adversarial networks," in *European Conference on Computer Vision.* Springer, 2016, pp. 702–716.

[46] M. Mathieu, C. Couprie, and Y. LeCun, "Deep multi-scale video prediction beyond mean square error," *arXiv preprint arXiv:1511.05440*, 2015.

[47] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.

[48] D. Pathak, P. Krahenbuhl, J. Donahue, T. Darrell, and A. A. Efros, "Context encoders: Feature learning by inpainting," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2536–2544.

[49] S. Reed, Z. Akata, X. Yan, L. Logeswaran, B. Schiele, and H. Lee, "Generative adversarial text to image synthesis," in *Proceedings of The 33rd International Conference on Machine Learning*, vol. 3, 2016.

[50] D. Yoo, N. Kim, S. Park, A. S. Paek, and I. S. Kweon, "Pixel-level domain transfer," in *European Conference on Computer Vision.* Springer, 2016, pp. 517–532.

[51] Y. Zhou and T. L. Berg, "Learning temporal transformations from time-lapse videos," in *European Conference on Computer Vision.* Springer, 2016, pp. 262–277.

[52] X. Wang and A. Gupta, "Generative image modeling using style and structure adversarial networks," in *European Conference on Computer Vision.* Springer, 2016, pp. 318–335.

[53] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical Image Computing and Computer-Assisted Intervention.* Springer, 2015, pp. 234–241.

[54] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural computation*, vol. 9, no. 8, pp. 1735–1780, 1997.

[55] D. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[56] A. Borji and L. Itti, "Exploiting local and global patch rarities for saliency detection," in *Computer Vision and Pattern Recognition (CVPR), 2012 IEEE Conference on.* IEEE, 2012, pp. 478–485.

[57] R. J. Peters, A. Iyer, L. Itti, and C. Koch, "Components of bottom-up gaze allocation in natural images," *Vision research*, vol. 45, no. 18, pp. 2397–2416, 2005.

[58] L. Itti and C. Koch, "A saliency-based search mechanism for overt and covert shifts of visual attention," *Vision research*, vol. 40, no. 10, pp. 1489–1506, 2000.

[59] T. Judd, K. Ehinger, F. Durand, and A. Torralba, "Learning to predict where humans look," in *Computer Vision, 2009 IEEE 12th international conference on.* IEEE, 2009, pp. 2106–2113.

[60] S. Jetley, N. Murray, and E. Vig, "End-to-end saliency mapping via probability distribution prediction," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 5753–5761.