# CEDL Homework01 - Introduce a New NN with Memory
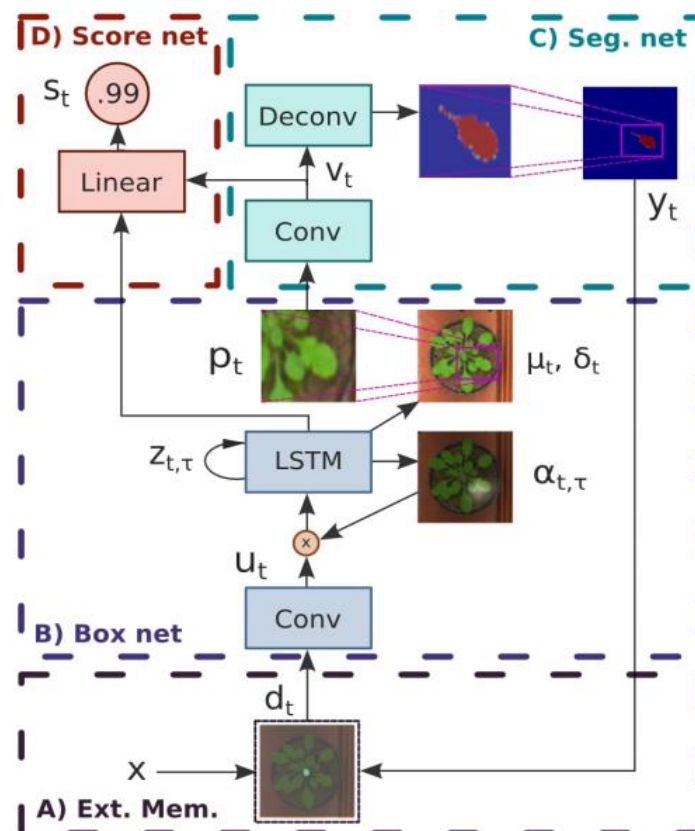
## End-to-End Instance Segmentation and Counting with Recurrent

105062536 劉祐欣  105065514 柯子逸

This paper addresses one of the main challenges of instance segmentation, that is, the object occlusion. Because the non-maximal suppression (NMS) may suppress the detection results for a heavily occluded object, too much overlap with foreground objects in cluttered scenes, the proposed method performs dynamic NMS to reason about occlusion. Counting the instances of an object class was considered jointly with instance segmentation.

The proposed model based on a recurrent neural network (RNN) that utilizes visual attention to perform instance segmentation and also performs dynamic NMS, using an object that is already segmented to aid in the discovery of an occluded object later in the sequence.



The model of this paper has four major components, illustrated in the above figure: (A) External memory (B) Box proposal network (C) Segmentation network (D) Scoring network, we will briefly illustrate the model below.
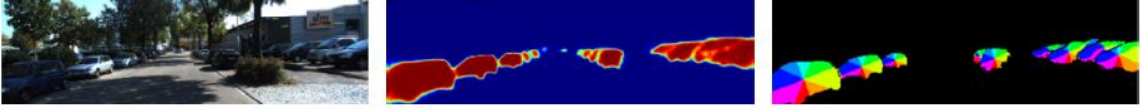
## Part A: External memory

External memory is used to track the state of the already segmented objects, decide where to look next. It provides object boundary details from all previous steps.

The formulation is the cumulative canvas: it stores the full history of segmentation outputs. Canvas has 10 channels in total: the first channel of the canvas keeps adding new pixels from the output of the previous time step, and the other channels store the preprocessed input image.

$$c_t = \begin{cases} \mathbf{0}, if \ t = 0 \\ \max(c_{t-1}, y_{t-1}), otherwise \end{cases}$$

$$d_t = [c_t, x]$$

$c_t$ is the channel of the canvas from the current time step, $y_{t-1}$ is the output segmentation sequences from the previous time step. The figure below illustrates the preprocessed input image, which has one channel for foreground segmentation and 8 channels for quantized angle maps.



## Part B: Box network

Part B localizes objects of interest. The CNN outputs a $H' \times W' \times L$ feature map $u_t$. A single glimpse is indexed by $\tau$. The proposed method allows the glimpse LSTM to look at different locations by feeding a dimension $L$ vector each time. In this part, dynamic pooling is employed to extract useful information along spatial dimensions, weighted by $\alpha_t^{h,w}$, initialized to be uniform over all locations.

$$u_t = \text{CNN}(d_t)$$

$$z_{t,\tau} = \text{LSTM}(z_{t,\tau-1}, \sum_{h,w} \alpha_{t,\tau}^{h,w} u_t^{h,w,l})$$

$$\alpha_{t,\tau} = \begin{cases} 1/(H' \times W'), if \ \tau = 0 \\ \text{MLP}(z_{t,\tau}), otherwise \end{cases}$$

$\text{LSTM}()$ denotes unrolling the long short-term memory by one time-step with the previous hidden state $z_{t-1}$ and current input $u_t$, and returning the current hidden state $z_{t,\tau}$; $\text{MLP}()$ denotes passing an input $z_{t,\tau}$ through a multi-layer perceptron and returning the hidden state $\alpha$.

Then pass LSTM()'s hidden state through a linear layer to obtain predicted box coordinates. The box is parameterized by its normalized center $(\hat{g}X, \hat{g}Y)$, and size $(\log\hat{\delta}X, \log\hat{\delta}Y)$. $\gamma$ is used when re-projecting the patch to the original image size.

$$\left[\hat{g}X, \hat{g}Y, \log\hat{\delta}X, \log\hat{\delta}Y, \log\hat{\sigma}X, \log\hat{\sigma}Y, \gamma\right] = \boldsymbol{w}_b^{\mathsf{T}}\boldsymbol{z}_{t,end} + \boldsymbol{w}_{b0}$$

$$gX = (\hat{g}X + 1)W/2$$

$$gY = (\hat{g}Y + 1)H/2$$

$$\delta X = \hat{\delta}XW$$

$$\delta Y = \hat{\delta}YH$$

Gaussian kernel is used to extract an $\widehat{H} \times \widehat{W}$ patch from the $\hat{x}$, a concatenation of the original image with $d_t$.

i, j index the location in the patch of dimension $\widehat{H} \times \widehat{W}$, and a, b index the location in the original image. $F_X$ and $F_Y$ indicate the contribution of the location $(a, b)$ in the original image towards the location $(i, j)$ in the extracted patch.

$$\mu_X^i = gX + (\delta X + 1) \cdot (i - \frac{\widehat{W}}{2} + 0.5)/\widehat{W}$$

$$\mu_Y^j = gY + (\delta Y + 1) \cdot (j - \frac{\widehat{H}}{2} + 0.5)/\widehat{H}$$

$$F_X^{a,i} = \frac{1}{\sqrt{2\pi}\sigma X}\exp(-\frac{(a-\mu^i X)^2}{2\sigma^2 X})$$

$$F_Y^{b,j} = \frac{1}{\sqrt{2\pi}\sigma Y}\exp(-\frac{(b-\mu^j Y)^2}{2\sigma^2 Y})$$

$$\widehat{\boldsymbol{x}_t} = [\boldsymbol{x_0}, \boldsymbol{d}_t]$$

$$\boldsymbol{p}_t = \text{Extract}(\widehat{\boldsymbol{x}_t}, F_Y, F_X) \equiv F_Y^{\mathsf{T}}\widehat{\boldsymbol{x}_t}F_X$$

**Part C: Segmentation network**

The segmentation network utilizes a convolution network to produce a feature map $v_t$, then input into a deconvolution layers to upsample the low-resolution feature map. Finally produce a patch-level segmentation prediction heat map $\hat{y}_t$ and re-project this patch prediction to the original image using the transpose of $F_Y$ and $F_X$.

$$\boldsymbol{v}_t = \text{CNN}(\boldsymbol{p}_t)$$

$$\widehat{\boldsymbol{y}_t} = \text{D} - \text{CNN}(\boldsymbol{v}_t)$$

$$y_t = \text{sigmoid}(\gamma \cdot \text{Extract}(\widehat{y_t}, F_Y^\mathsf{T}, F_X^\mathsf{T}) - \beta)$$

**Part D: Scoring network**

Part D takes information from the box network ($z_t$) and the segmentation network ($v_t$) to produce a score between 0 and 1 to estimate the number of objects in the image and terminate the sequential process.

$$s_t = \text{sigmoid}(w_{zs}^\mathsf{T} z_{t,end} + w_{vs}^\mathsf{T} v_t + w_{s0})$$

**Loss functions**

The total loss function is a sum of three losses: the segmentation matching IoU loss $\mathcal{L}_y$; the box IoU loss $\mathcal{L}_b$; and the score cross-entropy loss $\mathcal{L}_s$:

$$\mathcal{L}(y, b, s) = \mathcal{L}_y(y, y^*) + \mathcal{L}_b(b, b^*) + \mathcal{L}_s(s, s^*)$$

## Unique properties of the new NN:

--The next fine-grained part is decided to segment based on the segmentation output from the previous time step and the decision of interest region from multiple glimpses in the box network.

--The external memory addresses the segmentation in order, that is, the counting is addressed from the smaller segmentation to the bigger one.

## Applications of taking advantage of the properties:

--Detection of cancer cells, observation of cells growth, and other kinds of applications in the biomedical field.

--The warehousing management, which automatically check the quantity of the products in storage.

--Traffic flow detection and management.