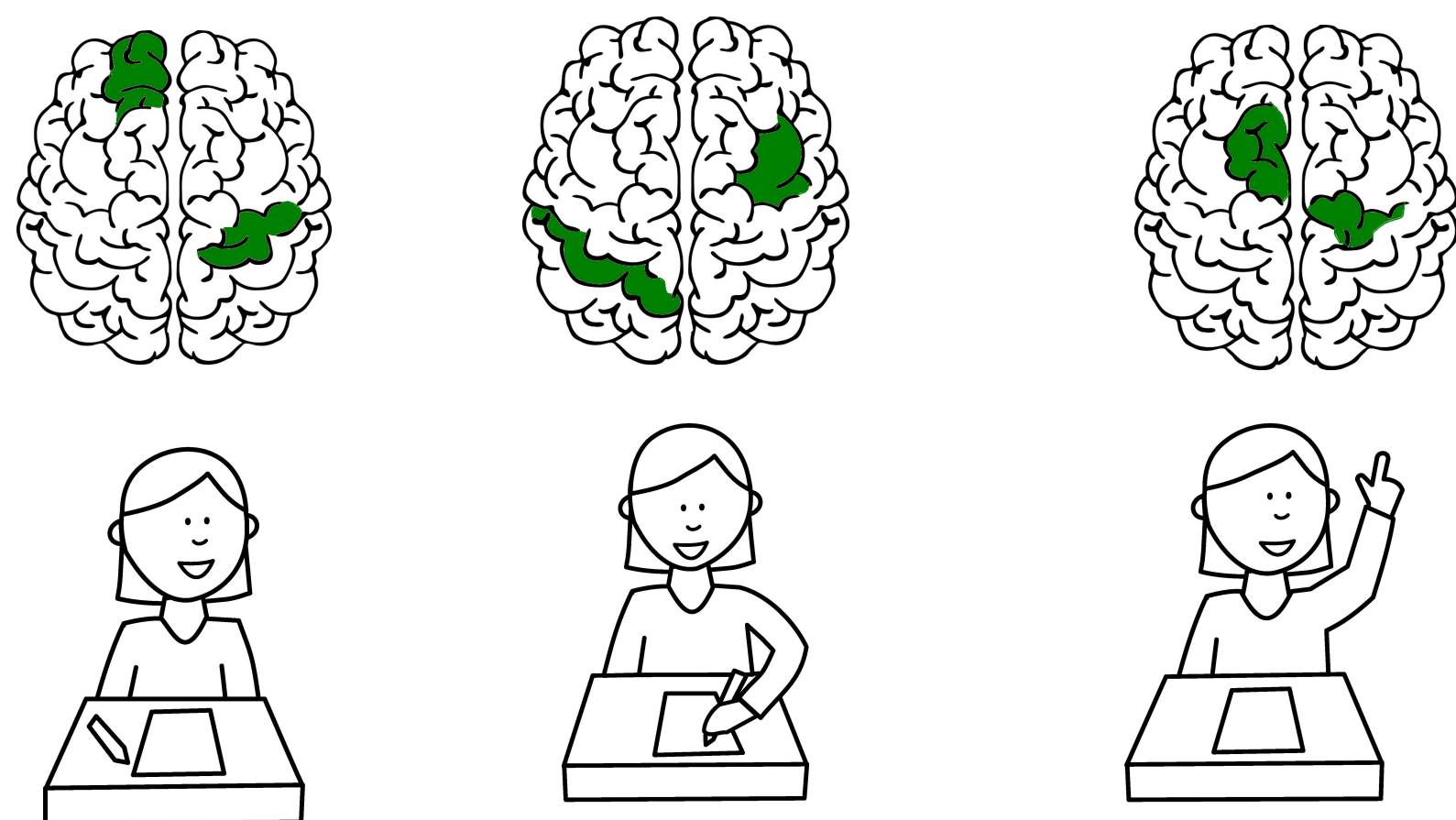


# INPUT-CELL ATTENTION REDUCES VANISHING SALIENCY OF RECURRENT NEURAL NETWORKS

Aya Abdelsalam Ismail<sup>1</sup>, Mohamed Gunady<sup>1</sup>, Luiz Pessoa<sup>2</sup>, Héctor Corrada Bravo<sup>\*1</sup> and Soheil Feizi<sup>\*1</sup>

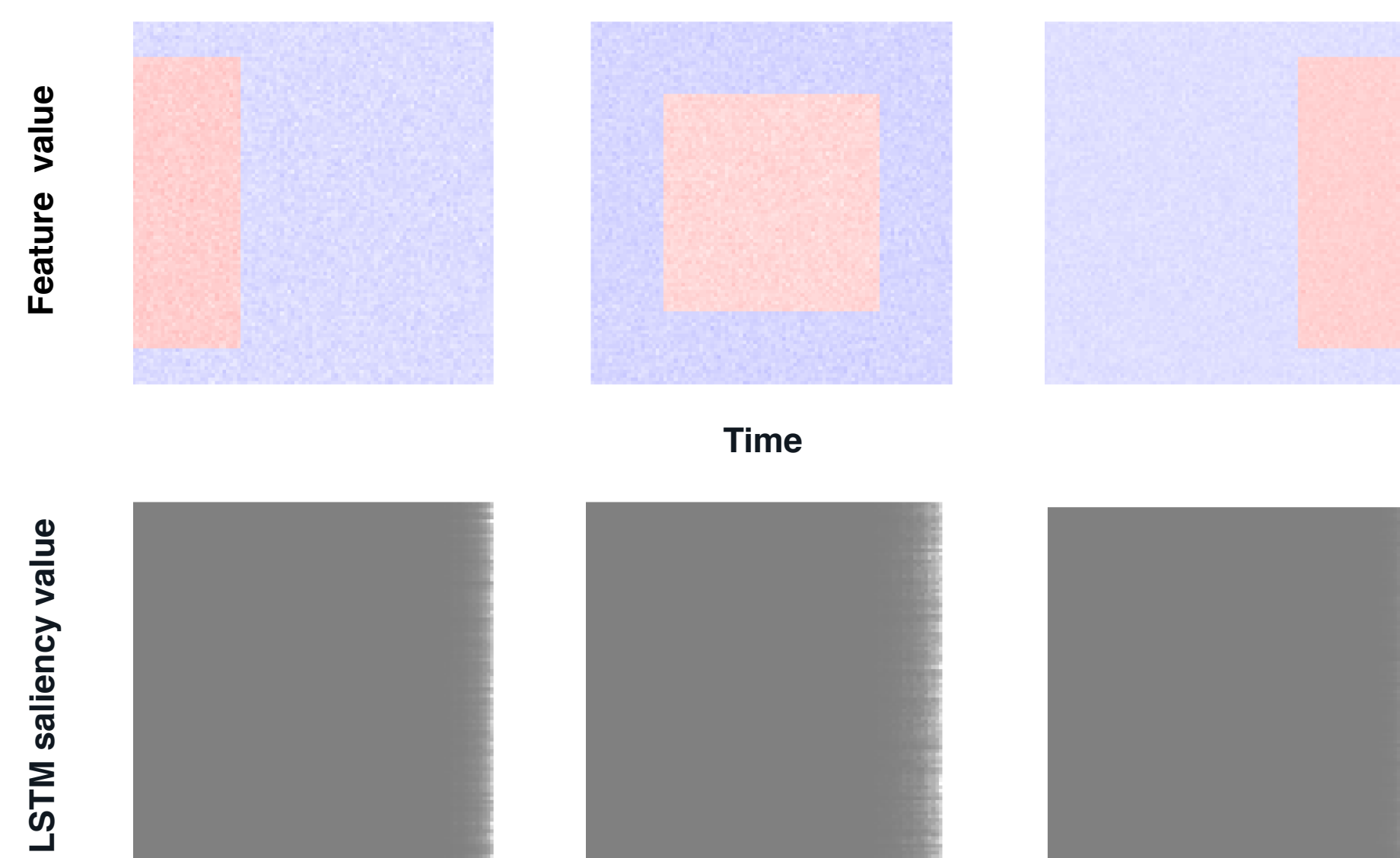
## Saliency for Time Series



To characterize brain region activity pertinent to the task being performed by the subject, a saliency method should capture changes in feature importance (corresponding to brain regions) over time.

## Saliency in RNNs Vanishes over Time

In RNNs, specifically LSTMs, saliency vanishes over time, biasing detection of salient features only to last few time steps.



**Top row:** samples from synthetic datasets, where red represents important features and blue is Gaussian noise. **Bottom row:** saliency maps produced by LSTM for each case, importance is reported only in the last few time steps.



For more details please check the original paper.

### LSTM

$$\begin{aligned} i_t &= \sigma(W_{xi}x_t + W_{hi}h_{t-1} + b_i) \\ f_t &= \sigma(W_{xf}x_t + W_{hf}h_{t-1} + b_f) \\ o_t &= \sigma(W_{xo}x_t + W_{ho}h_{t-1} + b_o) \\ \tilde{c}_t &= \tanh(W_{xc}x_t + W_{hc}h_{t-1} + b_c) \\ c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \\ h_t &= o_t \odot \tanh(c_t) \end{aligned}$$

### Input-cell At.

$$\begin{aligned} i_t &= \sigma(W_{Mi}M_t + W_{hi}h_{t-1} + b_i) \\ f_t &= \sigma(W_{Mf}M_t + W_{hf}h_{t-1} + b_f) \\ o_t &= \sigma(W_{Mo}M_t + W_{ho}h_{t-1} + b_o) \\ \tilde{c}_t &= \tanh(W_{Mc}M_t + W_{hc}h_{t-1} + b_c) \\ c_t &= f_t \odot c_{t-1} + i_t \odot \tilde{c}_t \\ h_t &= o_t \odot \tanh(c_t) \end{aligned}$$

## Why does LSTM Saliency Vanish?

- The relevance ( $R$ ) of each input feature at a given time step ( $x_{t_i}$ ) to the output of network ( $S$ ) is given by:

$$\begin{aligned} R_T(x_t) &= \left| \frac{\partial S(x_T)}{\partial x_t} \right| \\ &= \left| \frac{\partial S}{\partial h_T} \left( \prod_{i=T}^{t+1} \frac{\partial h_i}{\partial h_{i-1}} \right) \frac{\partial h_t}{\partial x_t} \right| \end{aligned}$$

- $\frac{\partial h_t}{\partial h_{t-1}}$  is the only term affected by the number of time steps.
- Solving the partial derivative:

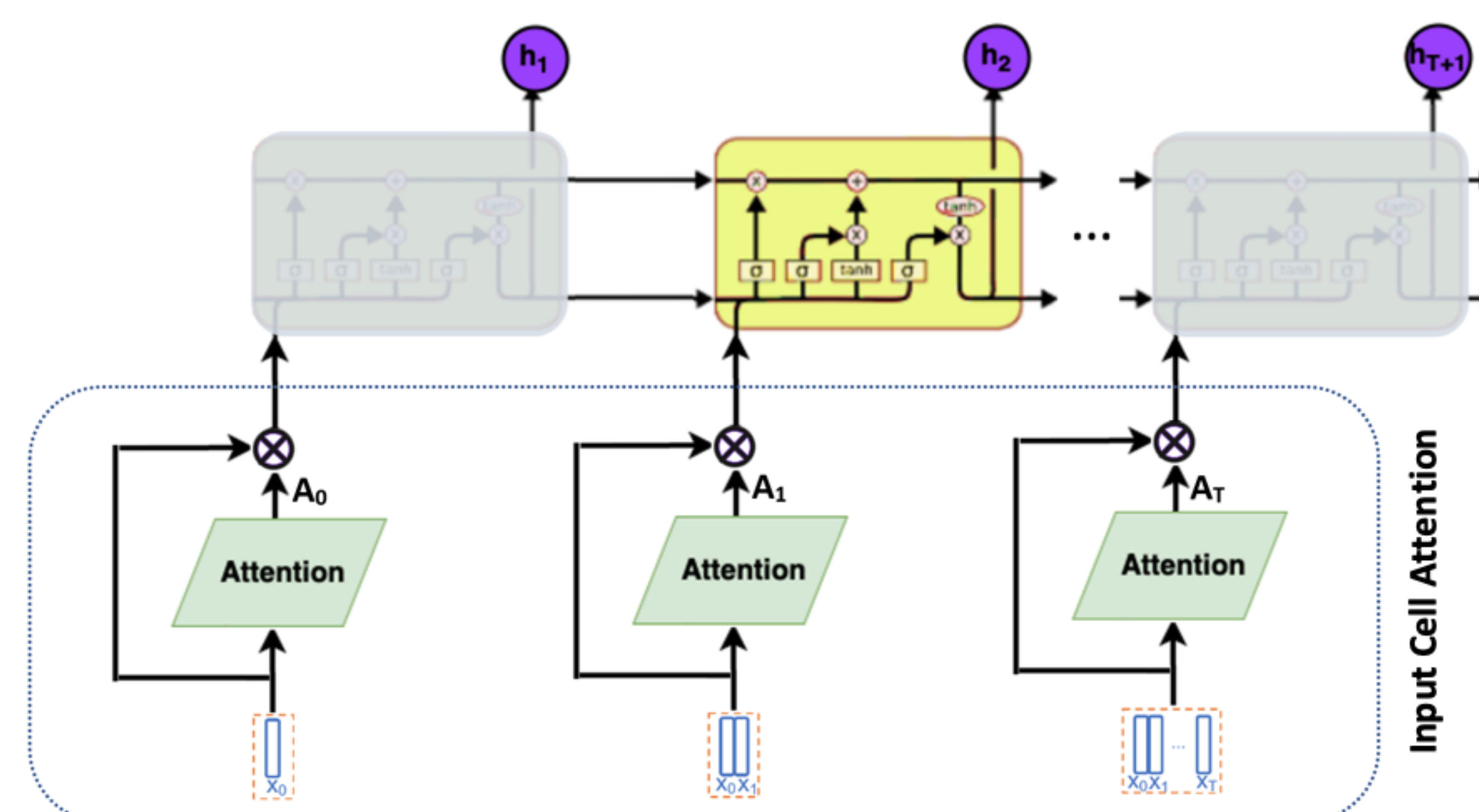
$$\begin{aligned} \frac{\partial h_t}{\partial h_{t-1}} &= \tanh(c_t) \left( [W_{ho}] (o_t \odot (1 - o_t)) \right) + o_t \odot (1 - \tanh^2(c_t)) \left[ c_{t-1} \left( [W_{hf}] (f_t \odot (1 - f_t)) \right) + \right. \\ &\quad \left. \tilde{c}_t \left( [W_{hi}] (i_t \odot (1 - i_t)) \right) + \right. \\ &\quad \left. i_t \left( [W_{hc}] (1 - \tilde{c}_t \odot \tilde{c}_t) \right) + f_t \right] \end{aligned}$$

- As  $t$  decreases, those terms multiplied by the weight matrix will eventually vanish,  $\frac{\partial h_t}{\partial h_{t-1}}$  will be reduced to :

$$\frac{\partial h_t}{\partial h_{t-1}} \approx o_t \odot (1 - \tanh^2(c_t)) \left[ f_t \right]$$

- The amount of information preserved depends on "forget gate" ( $f_t$ ).
- The relevance of input at time  $t$  eventually disappears as  $t$  decreases.

## Input-Cell Attention Reduces Vanishing Saliency

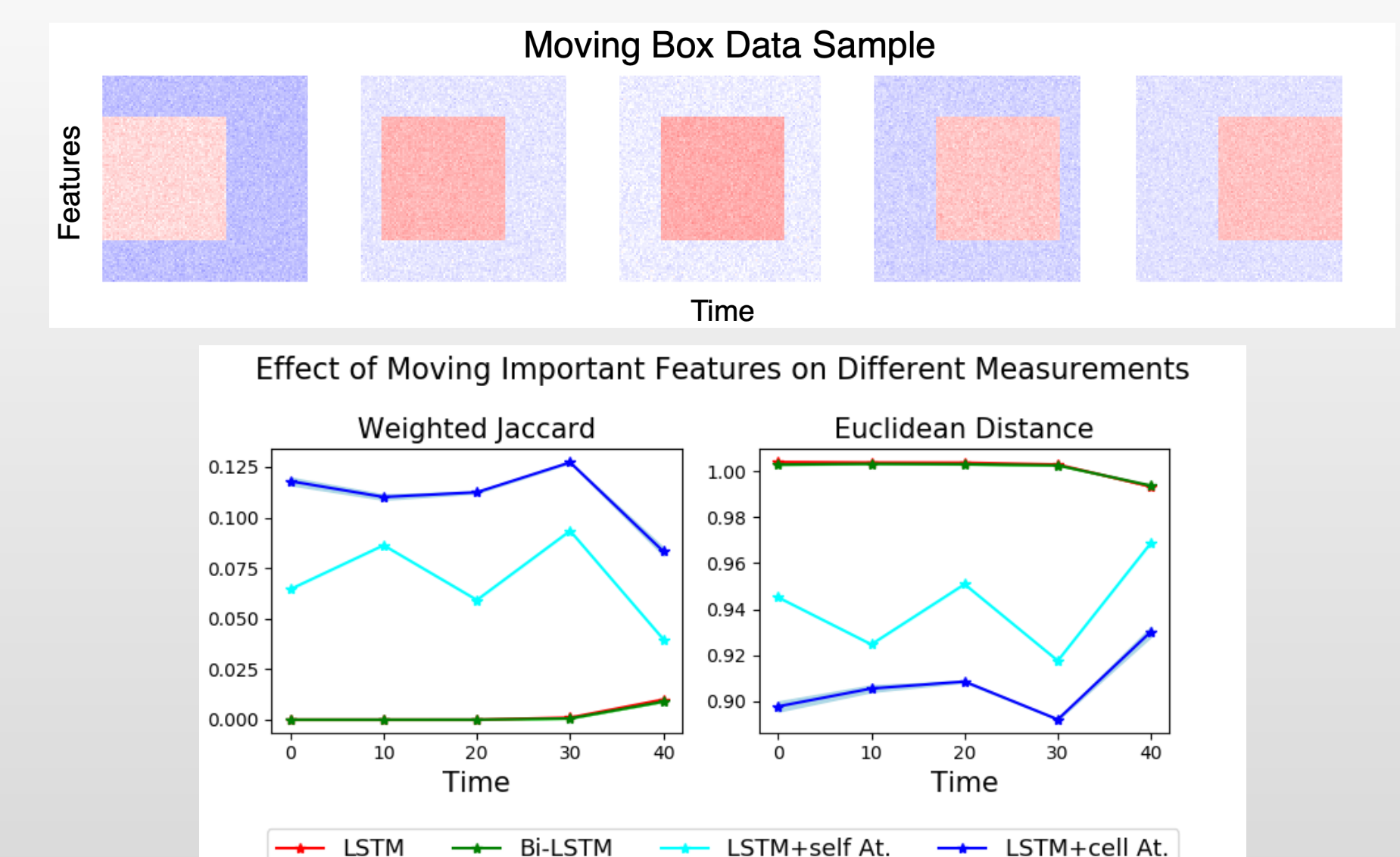


Input-cell attention uses a fixed-size matrix embedding, each row of the matrix attending to different inputs from current or previous time steps.

## Experiments

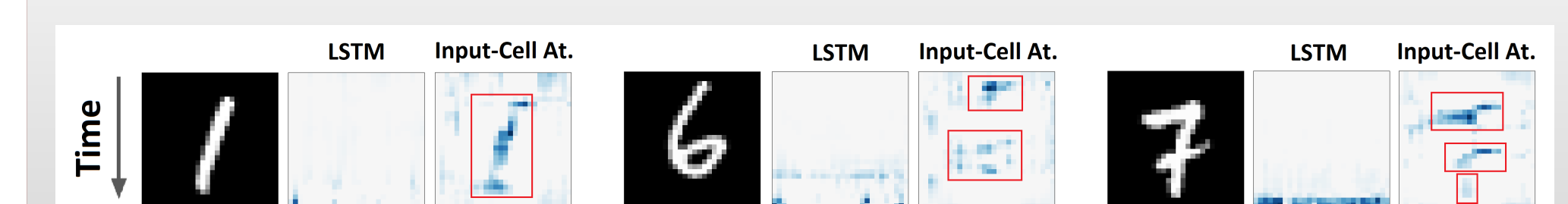
### Synthetic Datasets

The effect of changing the location of importance features in time on weighted Jaccard similarity and Euclidean distance for different architectures.



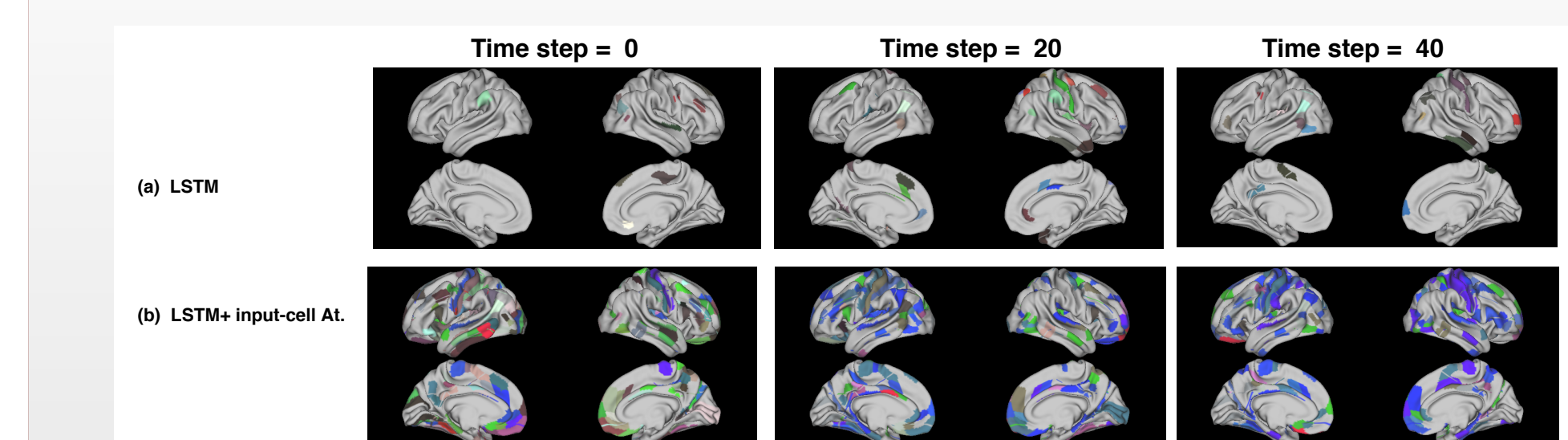
### MNIST as a Time Series

To validate the resulting saliency maps in cases where important features have more structured distributions of different shapes, we treat MNIST dataset as a time series.

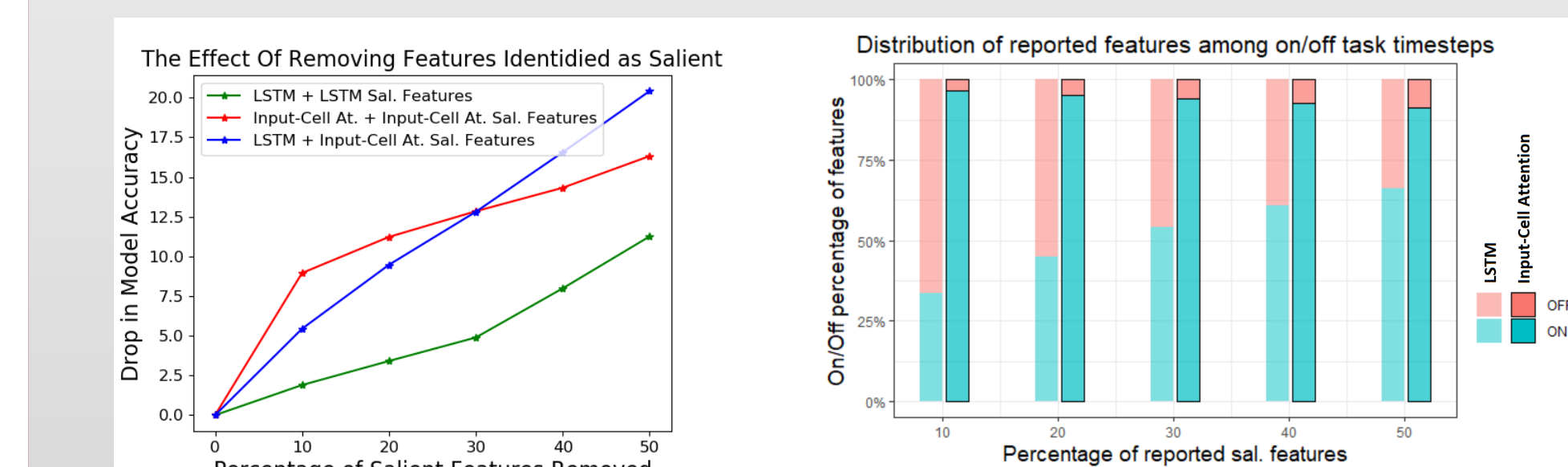


### Human Connectome Project fMRI Data

RNNs are used to classify time series based on the task performed by the subject while scanned by an fMRI machine. The saliency map produced by each model is shown below:



Studying features identified as salient by different models:



<sup>1</sup> Department of Computer Science, University of Maryland

<sup>2</sup> Department of Psychology, University of Maryland