# Input Switched Affine Recurrent Networks: An RNN Architecture Designed for Interpretability

Jakob N. Foerster* [1], Justin Gilmer*[1], Jascha Sohl-Dickstein[1],Jan Chorowski[1], David Sussillo [1]

[1]Google Brain

ICML,2017
Presenter: Arshdeep Sekhon

1. Interpreting Neural Networks
2. Crucial in many applications: self driving cars, medical diagnosis, power grid control, etc.

# Related Work

1. Post Hoc Analysis: After training a network, try and analyze it.
   + High Accuracy
   – Hard to interpret
   For example, break down LSTM model errors into classes
2. Design interpretability into the architecture
   + Better understanding
   – accuracy suffers
   For example, decision trees, logistic regression, etc.

## Vanilla RNN

$$h_{t+1} = \sigma(Ux_t + Wh_t + b) \tag{1}$$

$$l_t = \sigma(W_{ro}h_t + b_{ro}) \tag{2}$$

## ISAN

$$h_{t+1} = W_{x_t}h_t + b_{x_t} \tag{3}$$

$$l_t = W_{ro}h_t + b_{ro} \tag{4}$$

# ISAN: Accuracy Comparison

| Parameter count | 8e4 | 3.2e5 | 1.28e6 |
|:---:|:---:|:---:|:---:|
| RNN | 1.88 | 1.69 | 1.59 |
| IRNN | 1.89 | 1.71 | 1.58 |
| GRU | 1.83 | 1.66 | 1.59 |
| LSTM | 1.85 | 1.68 | 1.59 |
| ISAN | 1.92 | 1.71 | 1.58 |

Figure: Accuracy



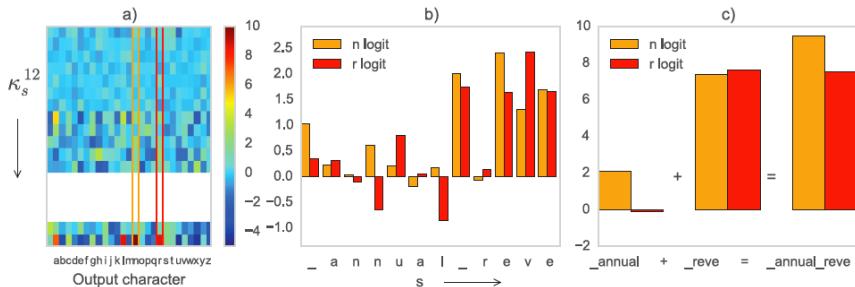Figure: ISAN makes fuller and more uniform use of its hidden state

$$\boldsymbol{h}_{t+1} = \boldsymbol{W}_{x_t} \boldsymbol{h}_t + \boldsymbol{b}_{x_t} \tag{5}$$

$$\mathbf{h}_t = \sum_{s=0}^{t} \left( \prod_{s'=s+1}^{t} \mathbf{W}_{\mathbf{x}_{s'}} \right) \mathbf{b}_{\mathbf{x}_s},$$

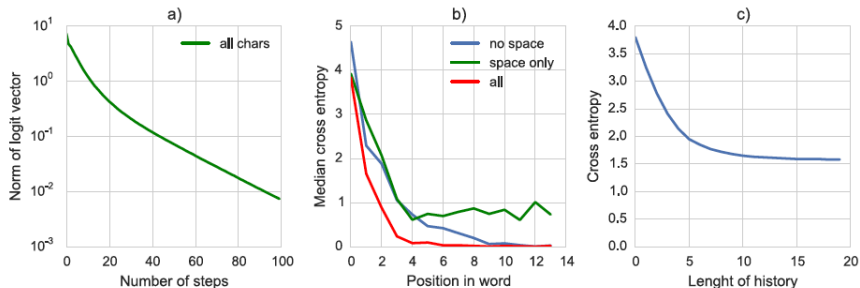$$\mathbf{l}_t = \mathbf{b}_{ro} + \sum_{s=0}^{t} \boldsymbol{\kappa}_s^t$$

$$\boldsymbol{\kappa}_s^t = \mathbf{W}_{ro} \left( \prod_{s'=s+1}^{t} \mathbf{W}_{\mathbf{x}_{s'}} \right) \mathbf{b}_{\mathbf{x}_s},$$
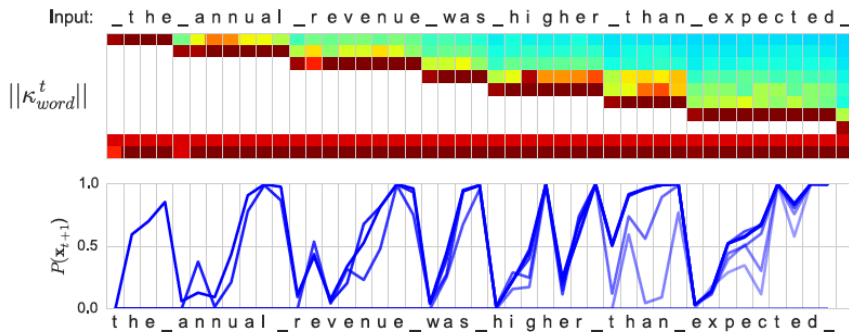
# Characters to Words

Jakob N. Foerster* , Justin Gilmer*, Jascha Input Switched Affine Recurrent Networks:An

1. Divide the hidden space into a subspace $\boldsymbol{P}_{\parallel}^{ro}$ spanned by the rows of the readout matrix $\boldsymbol{W}_{ro}$ and its orthogonal complement $\boldsymbol{P}_{\perp}^{ro}$
2. Thus, 27 dimensions for readout and (216-27) for computational subspace.
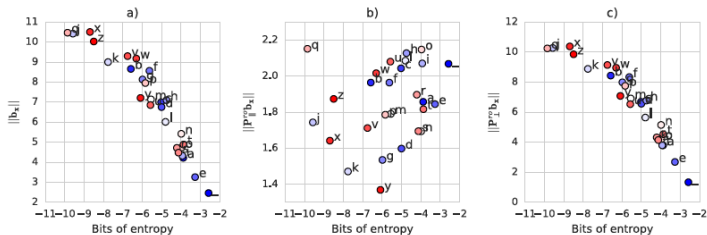
Figure: Information content related to the computation subspace.

# Change of basis



Figure: Correlation in $\boldsymbol{b}_x$. High correlation between vowels and consonants explained by $\boldsymbol{P}_{\parallel}^{ro}$

# Parantheses Counting Task

1. The Task: Count the number of opened parens [, (
2. Input: One hot encoded vector
3. Target Output: nesting level at previous timestep
4. output: two-hot encoded 0-5 count (12 dimensional 2-hot encoded vector)

## Paranthesis Counting

Using an augmented matrix and an augmented vector, it is possible to represent both the translation and the linear map using a single matrix multiplication:

ISAN:

$$\boldsymbol{h}_{t+1} = \boldsymbol{W}\boldsymbol{h}_t + \boldsymbol{b} \tag{6}$$

$$\boldsymbol{W}' = \left[ \begin{array}{cc} \boldsymbol{W} & \boldsymbol{b} \\ \boldsymbol{0}^T & 1 \end{array} \right]$$

$$\boldsymbol{h}_t' = \left[ \begin{array}{c} \boldsymbol{h}_t \\ \boldsymbol{1} \end{array} \right]$$

$$\boldsymbol{h}_{t+1}' = \boldsymbol{W}' \boldsymbol{h}_t' \tag{7}$$

# Paranthesis Counting: Change of Bases

1. Divide the hidden space into a subspace $P_\parallel^{ro}$ and its orthogonal complement $P_\perp^{ro}$
2. Learn bases by linear regression to encourage augmented matrices and hidden states to be sparse

$$\mathbf{W}'_x = \begin{bmatrix} \mathbf{W}^{rr}_x & \mathbf{W}^{rc}_x & \mathbf{b}^r_x \\ \mathbf{W}^{cr}_x & \mathbf{W}^{cc}_x & \mathbf{b}^c_x \\ \mathbf{0}^T & \mathbf{0}^T & 1 \end{bmatrix} \quad \mathbf{h}'_t = \begin{bmatrix} \mathbf{h}^r_t \\ \mathbf{h}^c_t \\ 1 \end{bmatrix}$$

and the update equation can be written as

$$\mathbf{h}'_{t+1} = \mathbf{W}'_x \mathbf{h}'_t = \begin{bmatrix} \mathbf{W}^{rr}_x \mathbf{h}^r_t + \mathbf{W}^{rc}_x \mathbf{h}^c_t + \mathbf{b}^r_x \\ \mathbf{W}^{cr}_x \mathbf{h}^r_t + \mathbf{W}^{cc}_x \mathbf{h}^c_t + \mathbf{b}^c_x \\ 1 \end{bmatrix} .$$

Equations after subspace decomposition
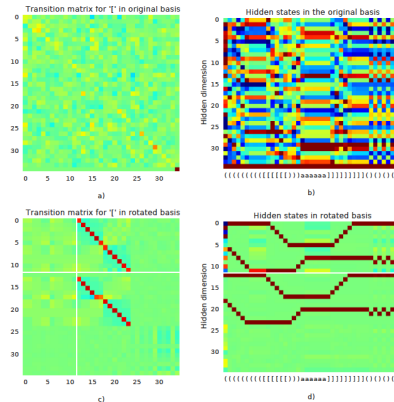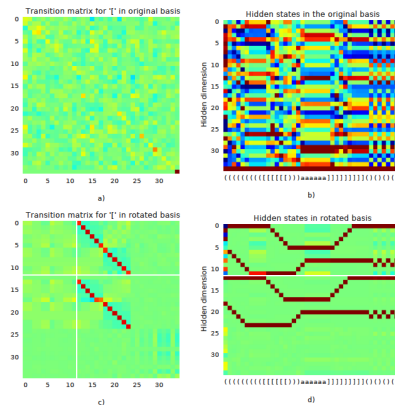
# Paranthesis Counting: Interpretation



Figure: Dynamics of ISAN for '['

1. leftmost 12 columns $\boldsymbol{W}_{[}^{rr}$ $\boldsymbol{W}_{[}^{cr}$ are zero
2. $h_t^r$ has no influence on $\boldsymbol{h}_{t+1}$

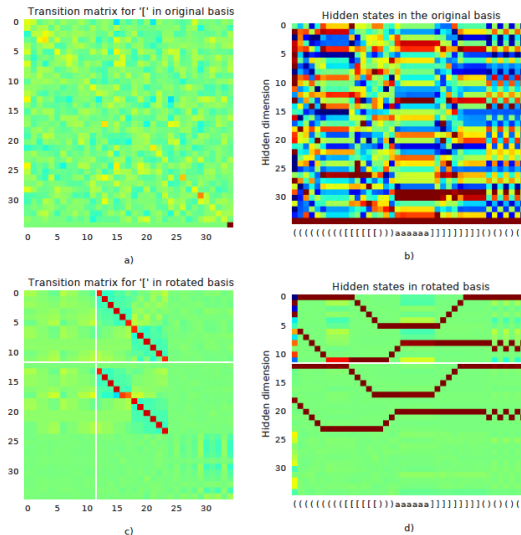Figure: Dynamics of ISAN for '['

1. $W_{[}^{rc}$ is identity; $h_t^r = h_{t-1}^c$

Figure: Dynamics of ISAN for '[': Delay Line Dynamics