

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

Jakob N. Foerster*¹, Justin Gilmer*¹, Jascha Sohl-Dickstein¹, Jan Chorowski¹, David Sussillo¹

¹Google Brain

ICML, 2017

Presenter: Arshdeep Sekhon

Motivation

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

- 1 Post Hoc Analysis: After training a network, try and analyze it.

- 1 Post Hoc Analysis: After training a network, try and analyze it.
- 2 Design interpretability into the architecture

- 1 Post Hoc Analysis: After training a network, try and analyze it.
For example, break down LSTM model errors into classes
 - + High Accuracy
 - Hard to interpret
- 2 Design interpretability into the architecture

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

Input Switched Affine Networks: ISAN

Vanilla RNN

$$\mathbf{h}_{t+1} = \sigma(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{h}_t + \mathbf{b}) \quad (1)$$

$$\mathbf{l}_t = \sigma(\mathbf{W}_{ro}\mathbf{h}_t + \mathbf{b}_{ro}) \quad (2)$$

Input Switched Affine Networks: ISAN

Vanilla RNN

$$\mathbf{h}_{t+1} = \sigma(\mathbf{U}\mathbf{x}_t + \mathbf{W}\mathbf{h}_t + \mathbf{b}) \quad (1)$$

$$\mathbf{l}_t = \sigma(\mathbf{W}_{ro}\mathbf{h}_t + \mathbf{b}_{ro}) \quad (2)$$

ISAN

$$\mathbf{h}_t = \mathbf{W}_{x_t}\mathbf{h}_{t-1} + \mathbf{b}_{x_t} \quad (3)$$

$$\mathbf{l}_t = \mathbf{W}_{ro}\mathbf{h}_t + \mathbf{b}_{ro} \quad (4)$$

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: *

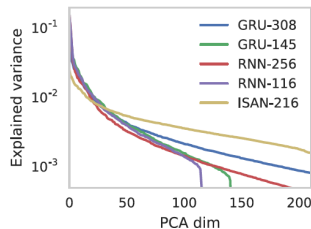
ISAN performs as well as other recurrent architectures

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: *

ISAN performs as well as other recurrent architectures



ISAN

$$\mathbf{h}_t = \mathbf{W}_{x_t} \mathbf{h}_{t-1} + \mathbf{b}_{x_t} \quad (5)$$

$$\mathbf{l}_t = \mathbf{W}_{r_o} \mathbf{h}_t + \mathbf{b}_{r_o} \quad (6)$$

ISAN

$$\mathbf{h}_t = \sum_{s=0}^t \left(\prod_{s'=s+1}^t \mathbf{W}_{x'_{s'}} \right) \mathbf{b}_{x_s} \quad (7)$$

ISAN

$$\kappa_s^t = \mathbf{w}_{ro} \left(\prod_{s'=s+1}^t \mathbf{w}_{x_{s'}} \right) \mathbf{b}_{x_s} \quad (8)$$

$$\mathbf{l}_t = \mathbf{b}_{ro} + \sum_{s=0}^t \kappa_s^t \quad (9)$$

Linearity of κ

Consider string: "_annual_revenue"

How does "_annual" affect output after "_rev"?

$$I_t = \mathbf{b}_{ro} + \sum_{s=0}^{t'} \kappa_s^t + \sum_{s=t'}^t \kappa_s^t \quad (10)$$

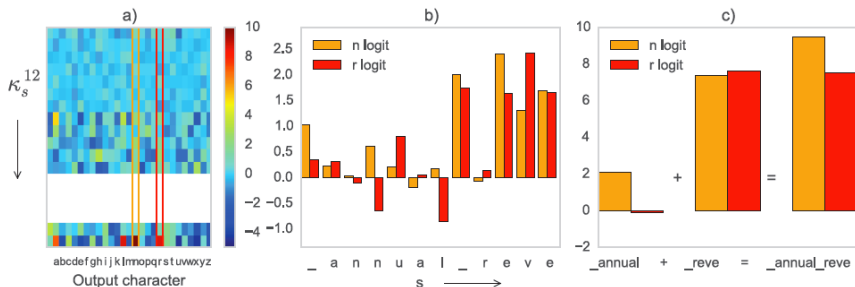


Figure: *

ISAN: information timescales of network

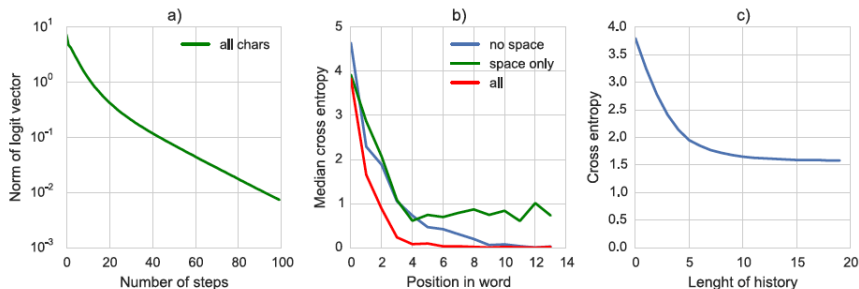


Figure: *

A κ_s^t averaged for all characters as a function of t-s

B Importance of " _ " character in decoding

C Cross entropy as a function of number of characters considered for prediction

Characters to Words

we can aggregate all of the κ_s^t belonging to a given word and visualize them as a single contribution to the prediction of the letters in the next word

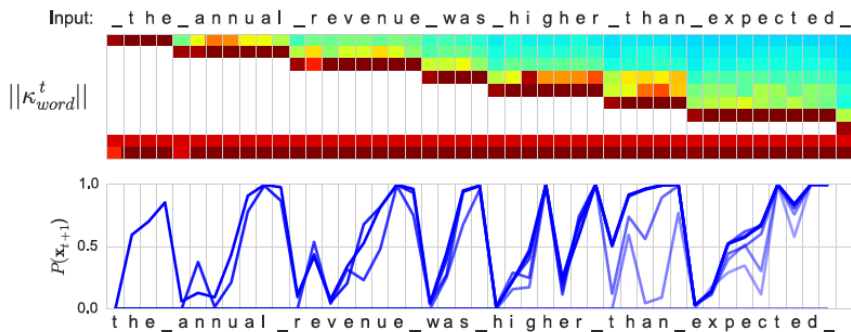


Figure: *

Change of Basis

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

Change of basis

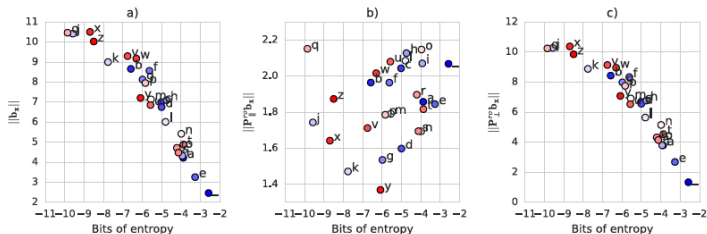


Figure: *

Information content related to the computation subspace.

- A the norm of the learnt b_x is strongly correlated to the log-probability of the unigram x in the training data.

Change of basis

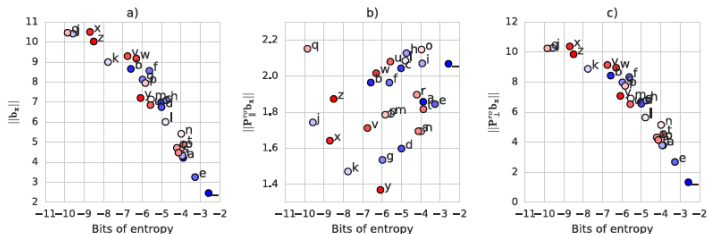


Figure: *

Information content related to the computation subspace.

- A the norm of the learnt b_x is strongly correlated to the log-probability of the unigram x in the training data.
- B this correlation is not related to reading out the next-step prediction

Change of basis

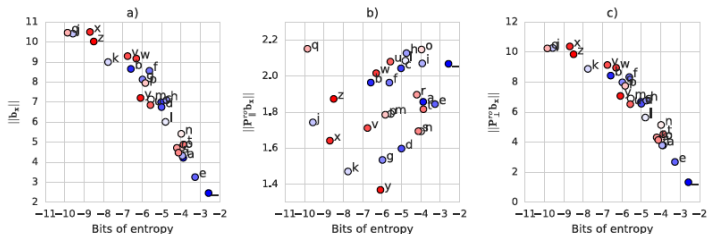


Figure: *

Information content related to the computation subspace.

- A the norm of the learnt b_x is strongly correlated to the log-probability of the unigram x in the training data.
- B this correlation is not related to reading out the next-step prediction
- C This implies a connection between information or surprise and distance in the computational subspace of state space.

Change of basis



- A Cosine distance/ correlation in original space
- B Cosine distance/ correlation in readout space or $P_{||}^{ro}$. two blocks of high correlations between the vowels and consonants respectively, while b_{\perp} is uncorrelated to either
- C Cosine distance/ correlation in readout space or P_{\perp}

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:

$$\mathbf{h}_{t+1} = \mathbf{W}\mathbf{h}_t + \mathbf{b} \quad (11)$$

$$\mathbf{h}'_{t+1} = \mathbf{W}'\mathbf{h}'_t \quad (12)$$

Paranthesis Counting: Change of Bases

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

Paranthesis Counting: Change of Bases

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

and the update equation can be written as

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

Figure: Equations after subspace decomposition

Paranthesis Counting: Interpretation

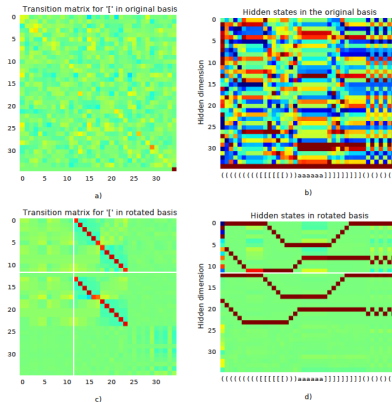


Figure: Dynamics of ISAN for '['

- 1 leftmost 12 columns $W_{[}^{rr}$ $W_{[}^{cr}$ are zero
- 2 h_t^r has no influence on h_{t+1}

Paranthesis Counting: Interpretation

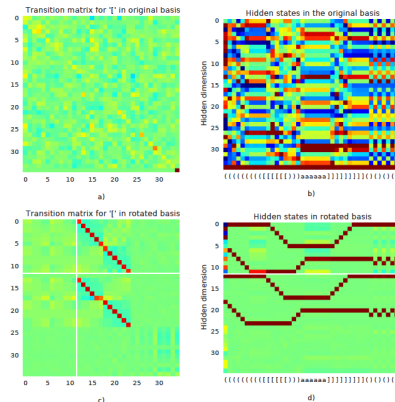


Figure: Dynamics of ISAN for '['

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

Paranthesis Counting: Interpretation

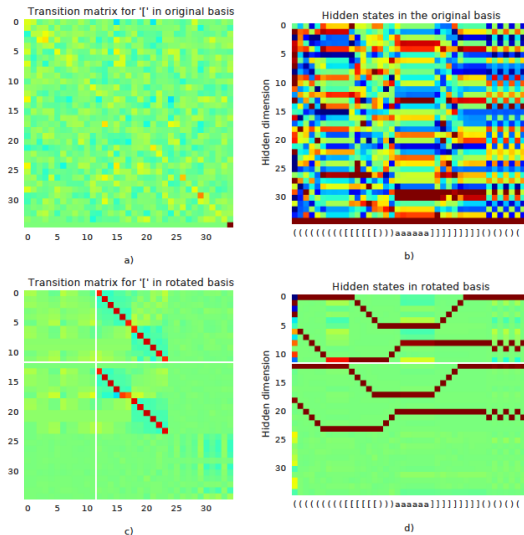


Figure: Dynamics of ISAN for 'l': Delay Line Dynamics