(Srivastava et al., 2015).

$$\mathbf{c}_t = \mathrm{s}\left(\mathbf{W^{ch}}\,\mathbf{h}_{t-1} + \mathbf{W^{cx}}\,\mathbf{x}_t + \mathbf{b^c}\right) \tag{2}$$

$$\mathbf{g}_t = \sigma\left(\mathbf{W^{gh}}\,\mathbf{h}_{t-1} + \mathbf{W^{gx}}\,\mathbf{x}_t + \mathbf{b^g} + b^{fg}\right) \tag{3}$$

$$\mathbf{h}_t = \mathbf{g}_t \cdot \mathbf{h}_{t-1} + (\mathbf{1} - \mathbf{g}_t) \cdot \mathbf{c}_t \tag{4}$$

GRU - GATED RECURRENT UNIT (CHO ET AL., 2014)

$$\mathbf{r}_t = \sigma\left(\mathbf{W^{rh}}\,\mathbf{h}_{t-1} + \mathbf{W^{rx}}\,\mathbf{x}_t + \mathbf{b^r}\right) \tag{5}$$

$$\mathbf{u}_t = \sigma\left(\mathbf{W^{uh}}\,\mathbf{h}_{t-1} + \mathbf{W^{ux}}\,\mathbf{x}_t + \mathbf{b^u} + b^{fg}\right) \tag{6}$$

$$\mathbf{c}_t = \mathrm{s}\left(\mathbf{W^{ch}}\,(\mathbf{r}_t \cdot \mathbf{h}_{t-1}) + \mathbf{W^{cx}}\,\mathbf{x}_t + \mathbf{b^c}\right) \tag{7}$$

$$\mathbf{h}_t = \mathbf{u}_t \cdot \mathbf{h}_{t-1} + (\mathbf{1} - \mathbf{u}_t) \cdot \mathbf{c}_t \tag{8}$$

LSTM - LONG SHORT TERM MEMORY(HOCHREITER & SCHMIDHUBER, 1997)

$$\mathbf{i}_t = \sigma\left(\mathbf{W^{ih}}\,\mathbf{h}_{t-1} + \mathbf{W^{ix}}\,\mathbf{x}_t + \mathbf{b^i}\right) \tag{9}$$

$$\mathbf{f}_t = \sigma\left(\mathbf{W^{fh}}\,\mathbf{h}_{t-1} + \mathbf{W^{fx}}\,\mathbf{x}_t + \mathbf{b^f} + b^{fg}\right) \tag{10}$$

$$\mathbf{c}_t^{in} = \mathrm{s}\left(\mathbf{W^{ch}}\,\mathbf{h}_{t-1} + \mathbf{W^{cx}}\,\mathbf{x}_t + \mathbf{b^c}\right) \tag{11}$$

$$\mathbf{c}_t = \mathbf{f}_t \cdot \mathbf{c}_{t-1} + \mathbf{i}_t \cdot \mathbf{c}_t^{in} \tag{12}$$

$$\mathbf{o}_t = \sigma\left(\mathbf{W^{oh}}\,\mathbf{h}_{t-1} + \mathbf{W^{ox}}\,\mathbf{x}_t + \mathbf{b^o}\right) \tag{13}$$

$$\mathbf{h}_t = \mathbf{o}_t \cdot \tanh(\mathbf{c}_t) \tag{14}$$

+RNN - INTERSECTION RNN

Due to the success of the UGRNN for shallower architectures in this study (see later figures on trainability), as well as some of the observed trainability problems for both the LSTM and GRU for deeper architectures (e.g. Figure 4h) we developed the Intersection RNN (denoted with a '+') architecture with a coupled depth gate in addition to a coupled recurrent gate. Additional influences for this architecture were the recurrent gating of the LSTM and GRU, and the depth gating from the highway network (Srivastava et al., 2015). This architecture has recurrent input, $\mathbf{h}_{t-1}$, and depth input, $\mathbf{x}_t$. It also has recurrent output, $\mathbf{h}_t$, and depth output, $\mathbf{y}_t$. Note that this architecture only applies between layers where $\mathbf{x}_t$ and $\mathbf{y}_t$ have the same dimension, and is not appropriate for networks with a depth of 1 (we exclude depth one +RNNs in our experiments).

$$\mathbf{y}_t^{in} = \mathrm{s1}\left(\mathbf{W^{yh}}\,\mathbf{h}_{t-1} + \mathbf{W^{yx}}\,\mathbf{x}_t + \mathbf{b^y}\right) \tag{15}$$

$$\mathbf{h}_t^{in} = \mathrm{s2}\left(\mathbf{W^{hh}}\,\mathbf{h}_{t-1} + \mathbf{W^{hx}}\,\mathbf{x}_t + \mathbf{b^h}\right) \tag{16}$$

$$\mathbf{g}_t^y = \sigma\left(\mathbf{W^{g^yh}}\,\mathbf{h}_{t-1} + \mathbf{W^{g^yx}}\,\mathbf{x}_t + \mathbf{b^{gy}} + b^{fg,y}\right) \tag{17}$$

$$\mathbf{g}_t^h = \sigma\left(\mathbf{W^{g^hh}}\,\mathbf{h}_{t-1} + \mathbf{W^{g^hx}}\,\mathbf{x}_t + \mathbf{b^{gh}} + b^{fg,h}\right) \tag{18}$$

$$\mathbf{y}_t = \mathbf{g}_t^y \cdot \mathbf{x}_t + (\mathbf{1} - \mathbf{g}_t^y) \cdot \mathbf{y}_t^{in} \tag{19}$$

$$\mathbf{h}_t = \mathbf{g}_t^h \cdot \mathbf{h}_{t-1} + (\mathbf{1} - \mathbf{g}_t^h) \cdot \mathbf{h}_t^{in} \tag{20}$$

In practice we used ReLU for s1 and $\tanh$ for s2.

## 2 CAPACITY EXPERIMENTS

### 2.1 PER-PARAMETER CAPACITY

A foundational result in machine learning is that a single-layer perceptron with $N^2$ parameters can store at least 2 bits of information per parameter (Cover, 1965; Gardner, 1988; Baldi & Venkatesh, 1987). More precisely, a perceptron can implement a mapping from $2N$, $N$-dimensional, input vectors to arbitrary $N$-dimensional binary output vectors, subject only to the extremely weak restriction that the input vectors be in general position. RNNs provide a far more complex input-output mapping, with hidden units, recurrent dynamics, and a diversity of nonlinearities. Nonetheless, we wondered if there were analogous capacity results for RNNs that we might be able to observe empirically.