# Wider and Deeper, Cheaper and Faster: Tensorized LSTMs for Sequence Learning

Zhen He[1,2], Shaoging Gao[3], Liang Xiao[1], Daxue Liu[1], Hangen He[1], David Barber[2,4]

[1]National University of Defense Technology  [2]University College London  [3]Sichuan University  [4]Alan Turing Institute
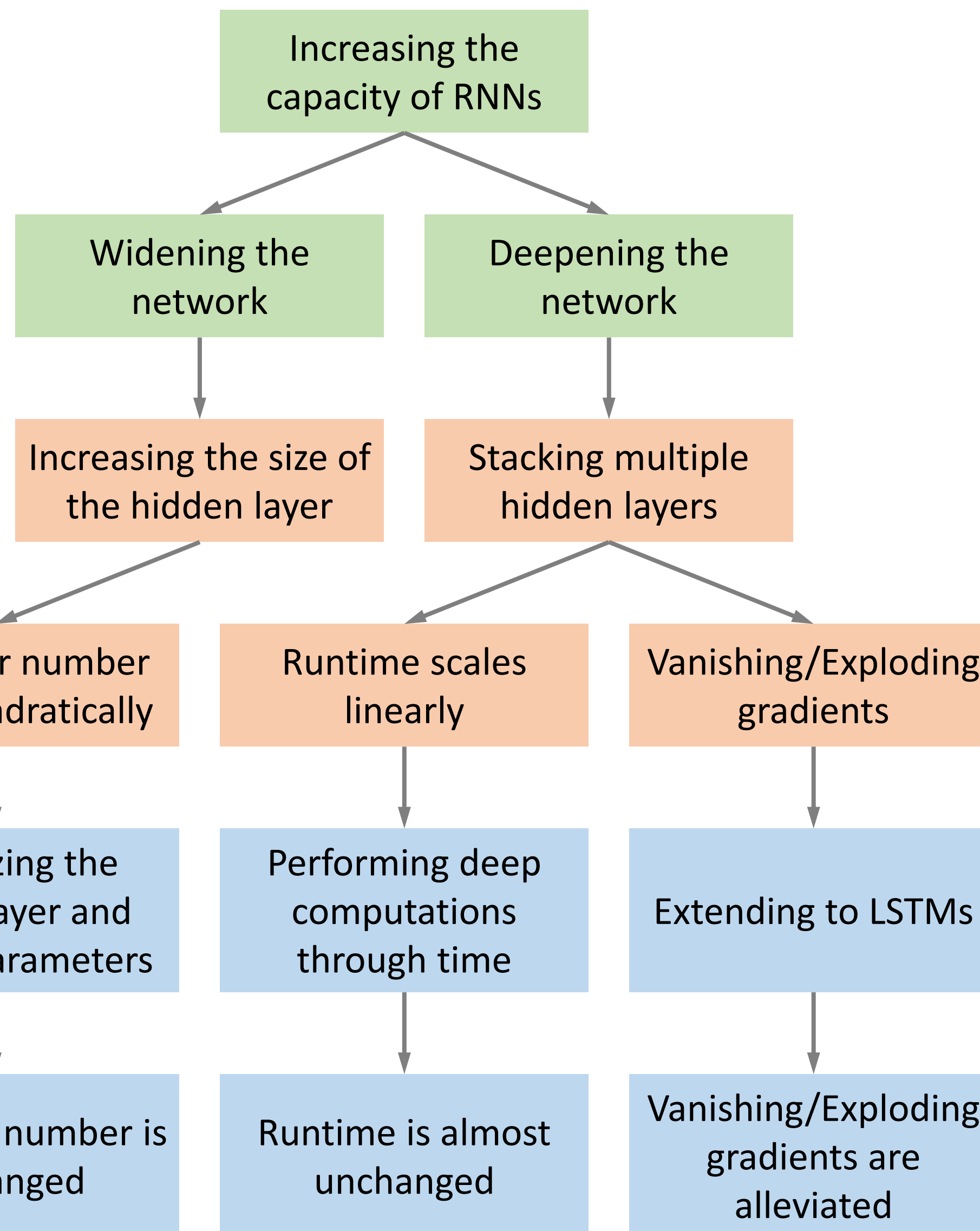
## 1. Introduction



*Motivation:* Increasing the capacity of RNNs

*How to increase their capacity?* Widening the network / Deepening the network

*Common solutions:* Increasing the size of the hidden layer / Stacking multiple hidden layers

*Drawbacks:* Parameter number scales quadratically / Runtime scales linearly / Vanishing/Exploding gradients

*Our solutions:* Tensorizing the hidden layer and sharing parameters / Performing deep computations through time / Extending to LSTMs

*Advantages:* Parameter number is unchanged / Runtime is almost unchanged / Vanishing/Exploding gradients are alleviated
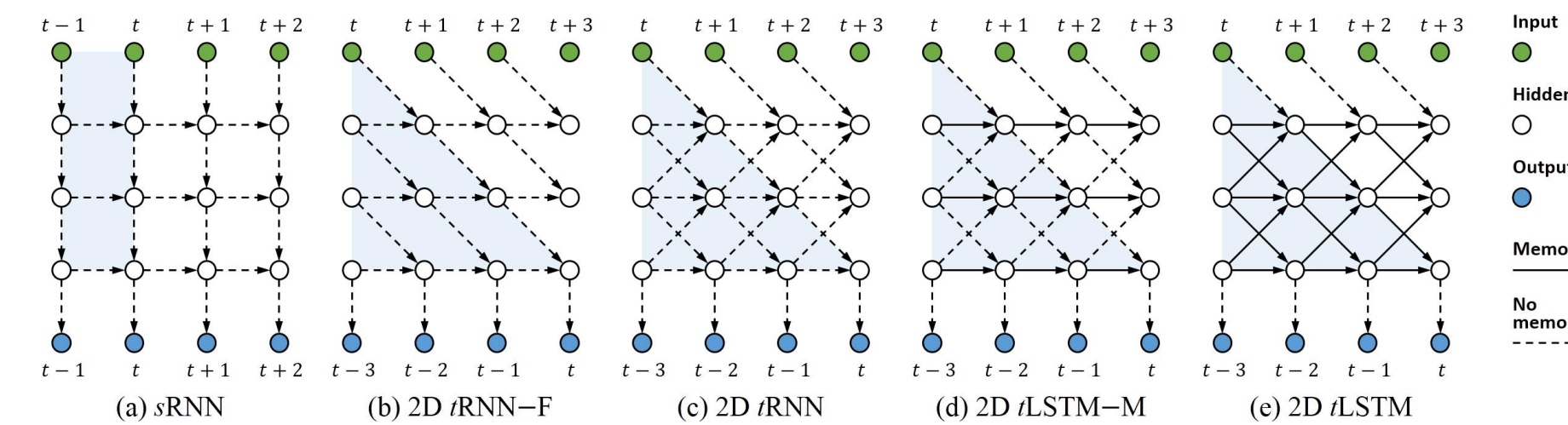
## 2. Method



Figure 1: Examples of *s*RNN, *t*RNNs and *t*LSTMs. (a) A 3-layer *s*RNN. (b) A 2D *t*RNN without (–) feedback (F) connections, which can be thought as a *skewed* version of (a). (c) A 2D *t*RNN. (d) A 2D *t*LSTM without (–) memory (M) cell convolutions. (e) A 2D *t*LSTM. In each model, the blank circles in column 1 to 4 denote the hidden state at timestep $t-1$ to $t+2$, respectively, and the blue region denotes the receptive field of the current output $y_t$. In (b)-(e), the outputs are delayed by $L-1=2$ timesteps, where $L=3$ is the depth.
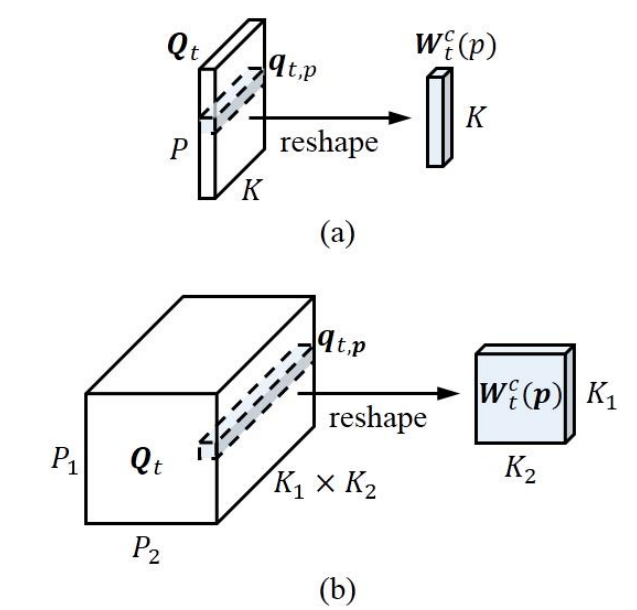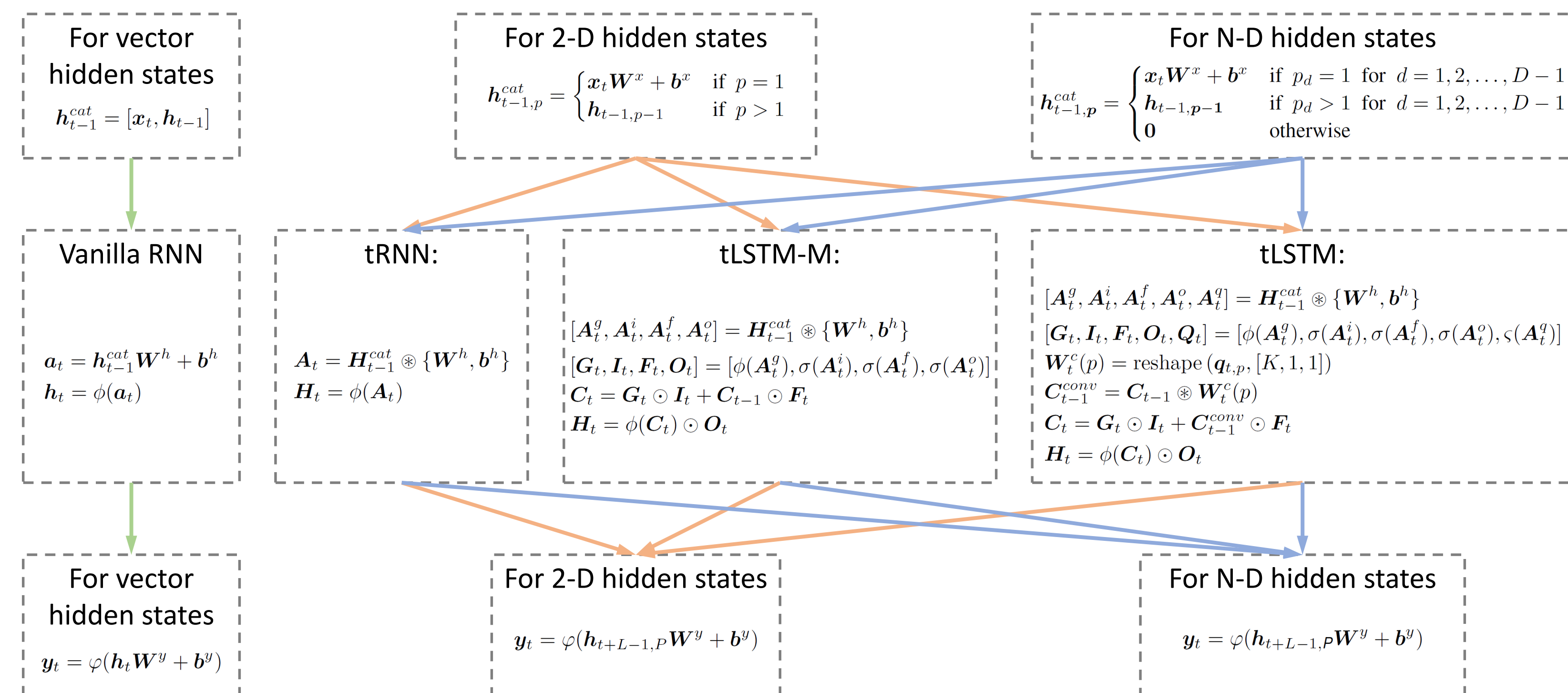
Figure 2: Illustration of generating the memory cell convolution kernel, where (a) is for 2D tensors and (b) for 3D tensors.

**Concatenating the input:**

For vector hidden states
$$h_{t-1,p}^{cat} = [x_t, h_{t-1}]$$

For 2-D hidden states
$$h_{t-1,p}^{cat} = \begin{cases} x_t W^x + b^x & \text{if } p = 1 \\ h_{t-1,p-1} & \text{if } p > 1 \end{cases}$$

For N-D hidden states
$$h_{t-1,p}^{cat} = \begin{cases} x_t W^x + b^x & \text{if } p_d = 1 \text{ for } d = 1,2,\ldots,D-1 \\ h_{t-1,p-1} & \text{if } p_d > 1 \text{ for } d = 1,2,\ldots,D-1 \\ 0 & \text{otherwise} \end{cases}$$

**Updating the hidden state:**

Vanilla RNN
$$a_t = h_{t-1}^{cat} W^h + b^h$$
$$h_t = \phi(a_t)$$

tRNN:
$$A_t = H_{t-1}^{cat} \circledast \{W^h, b^h\}$$
$$H_t = \phi(A_t)$$

tLSTM-M:
$$[A_t^g, A_t^i, A_t^f, A_t^o] = H_{t-1}^{cat} \circledast \{W^h, b^h\}$$
$$[G_t, I_t, F_t, O_t] = [\phi(A_t^g), \sigma(A_t^i), \sigma(A_t^f), \sigma(A_t^o)]$$
$$C_t = G_t \odot I_t + C_{t-1} \odot F_t$$
$$H_t = \phi(C_t) \odot O_t$$

tLSTM:
$$[A_t^g, A_t^i, A_t^f, A_t^o, A_t^q] = H_{t-1}^{cat} \circledast \{W^h, b^h\}$$
$$[G_t, I_t, F_t, O_t, Q_t] = [\phi(A_t^g), \sigma(A_t^i), \sigma(A_t^f), \sigma(A_t^o), \varsigma(A_t^q)]$$
$$W_t^c(p) = \text{reshape}(q_{t,p}, [K,1,1])$$
$$C_{t-1}^{conv} = C_{t-1} \circledast W_t^c(p)$$
$$C_t = G_t \odot I_t + C_{t-1}^{conv} \odot F_t$$
$$H_t = \phi(C_t) \odot O_t$$

**Generating the output:**

For vector hidden states
$$y_t = \varphi(h_t W^y + b^y)$$

For 2-D hidden states
$$y_t = \varphi(h_{t+L-1,P} W^y + b^y)$$

For N-D hidden states
$$y_t = \varphi(h_{t+L-1,P} W^y + b^y)$$

## 3. Experiments

### Comparison of different configurations



Figure 3: Performance and runtime of different configurations on Wikipedia.

Figure 4: Performance and runtime of different configurations on the addition (left) and memorization (right) tasks.

Figure 5: Performance and runtime of different configurations on sequential MNIST (left) and sequential *p*MNIST (right).

- Wider and deeper networks perform better.
- Parameter number and runtime are invariant.
- Memory cell convolutions are crucial to maintain improvement.
- Feedback/tensorization/CN is useful.

### Comparison to the state-of-the-art methods

Table 1: Test BPC on Wikipedia.

| | BPC | # Param. |
|---|---|---|
| MI-LSTM [51] | 1.44 | ≈17M |
| mLSTM [33] | 1.42 | ≈20M |
| HyperLSTM+LN [23] | 1.34 | 26.5M |
| HM-LSTM+LN [11] | 1.32 | ≈35M |
| Large RHN [54] | 1.27 | ≈46M |
| Large FS-LSTM-4 [38] | 1.245 | ≈47M |
| 2 × Large FS-LSTM-4 [38] | **1.198** | ≈94M |
| 3D *t*LSTM+CN ($L=6$, $M=1200$) | 1.264 | 50.1M |

Table 2: Test accuracies on two algorithmic tasks.

| | Addition | | Memorization | |
|---|---|---|---|---|
| | Acc. | # Samp. | Acc. | # Samp. |
| Stacked LSTM [21] | 51% | 5M | >50% | 900K |
| Grid LSTM [30] | >99% | 550K | >99% | 150K |
| 3D *t*LSTM+CN ($L=7$) | >99% | **298K** | >99% | 115K |
| 3D *t*LSTM+CN ($L=10$) | >99% | 317K | >99% | **54K** |

Table 3: Test accuracies (%) on sequential MNIST/*p*MNIST.

| | MNIST | *p*MNIST |
|---|---|---|
| *i*RNN [33] | 97.0 | 82.0 |
| LSTM [2] | 98.2 | 88.0 |
| *u*RNN [2] | 95.1 | 91.4 |
| Full-capacity *u*RNN [49] | 96.9 | 94.1 |
| *s*TANH [53] | 98.1 | 94.0 |
| BN-LSTM [13] | 99.0 | 95.4 |
| Dilated GRU [8] | **99.2** | 94.6 |
| Dilated CNN [40] in [8] | 98.3 | **96.7** |
| 3D *t*LSTM+CN ($L=3$) | **99.2** | 94.9 |
| 3D *t*LSTM+CN ($L=5$) | **99.2** | 95.7 |

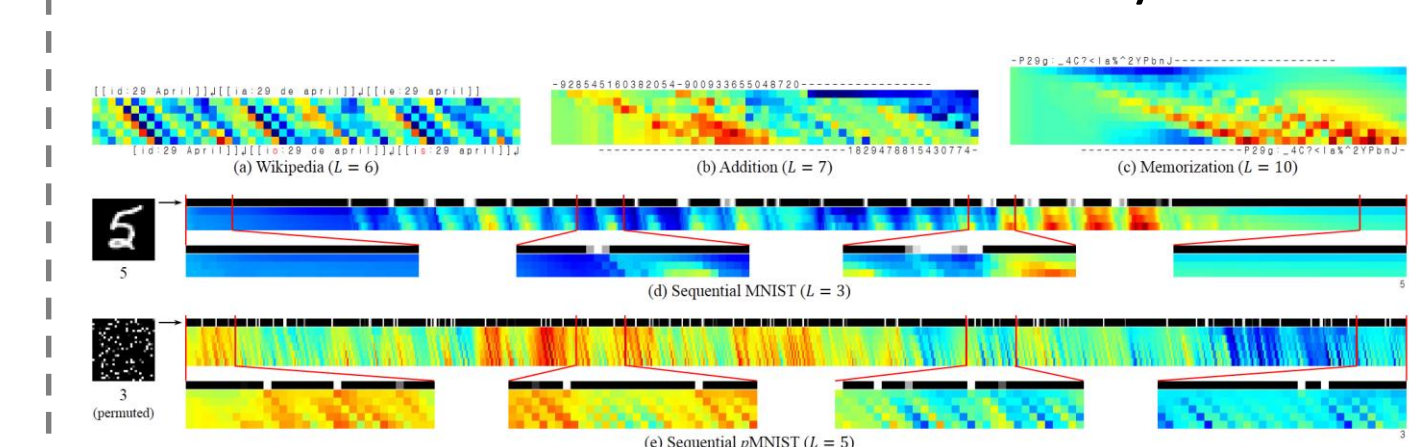### Visualization of the *t*LSTM memory cells



Figure 6: Visualization of the diagonal channel means of the *t*LSTM memory cells for each task. In each horizontal bar, the rows from top to bottom correspond to the diagonal locations from $p^{in}$ to $p^{out}$, the columns from left to right correspond to different timesteps (from 1 to $T+L-1$ for the full sequence, where $L-1$ is the time delay), and the values are normalized to be in range $[0, 1]$ for better visualization. Both full sequences in (d) and (e) are zoomed out horizontally.

- Wider (larger) tensors can encode more information, with less effort to compress it.
- Deep computations are indeed performed together with temporal computations, with long-range dependencies carried by memory cells.
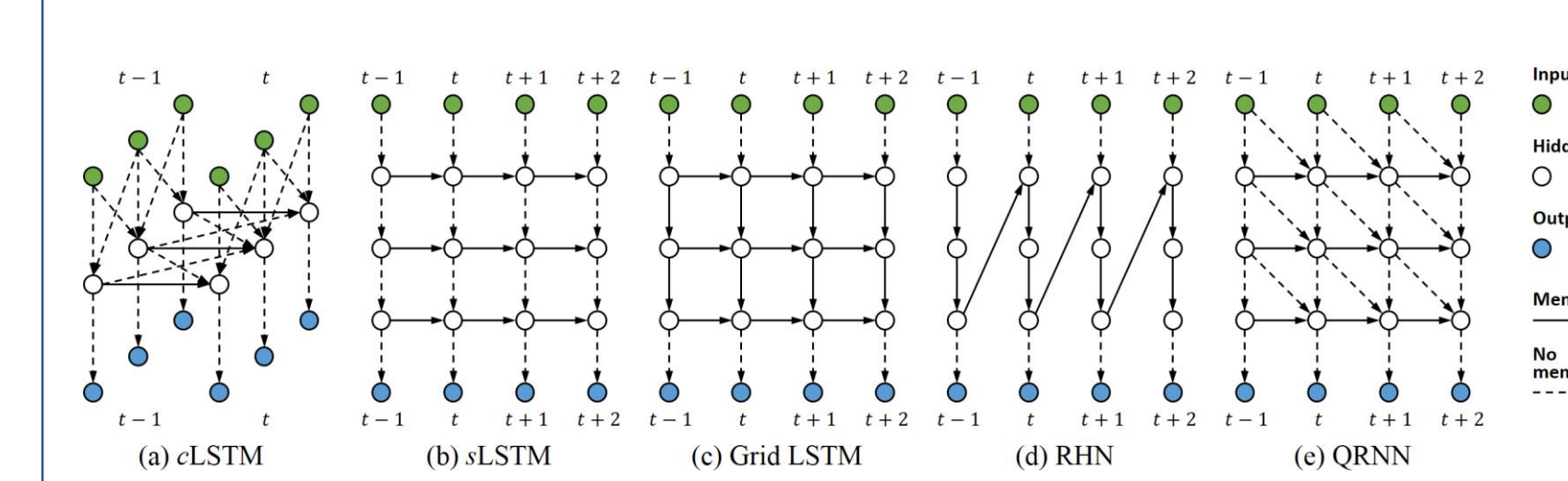
## 4. Related Work



Figure 7: Examples of models related to *t*LSTMs. (a) A single layer *c*LSTM [48] with vector array input. (b) A 3-layer *s*LSTM [21]. (c) A 3-layer Grid LSTM [30]. (d) A 3-layer RHN [54]. (e) A 3-layer QRNN [7] with kernel size 2, where costly computations are done by temporal convolution.

- Convolutional LSTMs (a) are for structured input.
- Stacked/Deep LSTMs (b, c, and d) typically multiply the runtime.
- Temporal parallelization (e) is potentially unsuitable for real-time online inference.