

MC-LSTM: Mass-Conserving LSTM



Pieter-Jan Hoedt & Frederik Kratzert @ DLCT

ELLIS Unit Linz and LIT AI Lab
Institute for Machine Learning
Johannes Kepler University, Linz, Austria



E-mail: {hoedt,kratzert}@ml.jku.at

Contributors



Johannes Kepler
University



Google



Pieter-Jan
Hoedt



Frederik
Kratzert



Daniel
Klotz



Christina
Halmich



Markus
Holzleitner



Grey
Nearing



Sepp
Hochreiter



Günter
Klambauer



Whereabouts

 [@ml_hoedt](https://twitter.com/ml_hoedt)

 [@fkratzert](https://twitter.com/fkratzert)

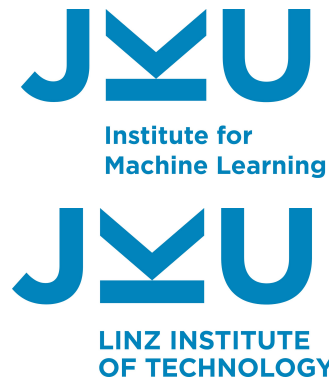
 [@ido87](https://twitter.com/ido87)

 [@tinahalmich](https://twitter.com/tinahalmich)

 [@GreyNearing](https://twitter.com/GreyNearing)



Linz, AUSTRIA



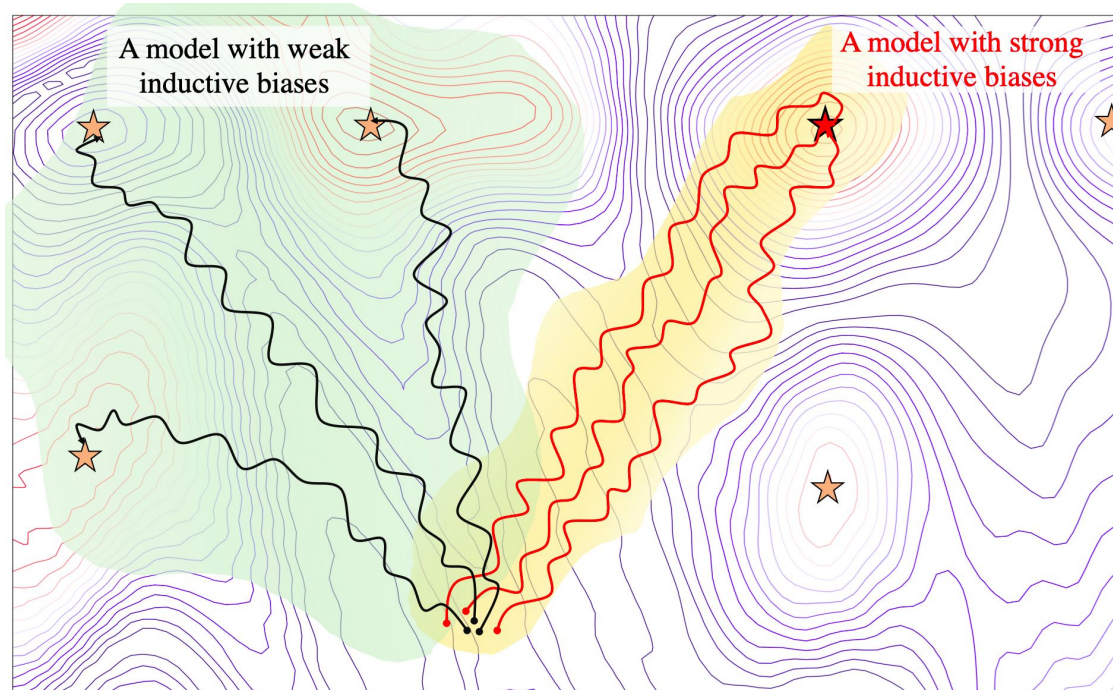
OUTLINE

- Motivation
- Model
- Experiments

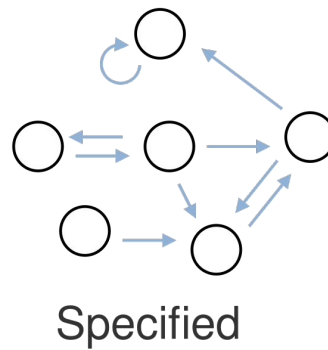
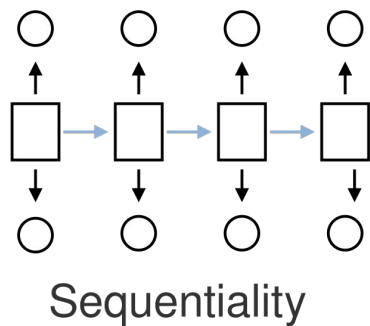
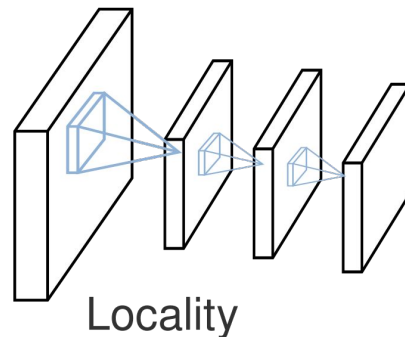
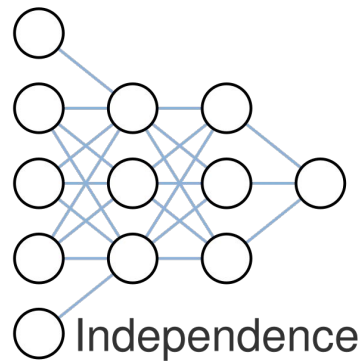
OUTLINE

- Motivation
- Model
- Experiments

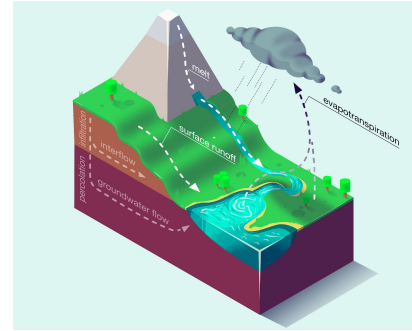
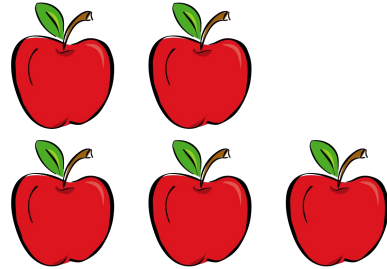
INDUCTIVE BIAS



INDUCTIVE BIAS



CONSERVATION LAWS



MASS CONSERVATION

Theorem 1 (Conservation property). *Let $m_c^\tau = \sum_{k=1}^K c_k^\tau$ be the mass contained in the system and $m_h^\tau = \sum_{k=1}^K h_k^\tau$ be the mass efflux, or, respectively, the accumulated mass in the MC-LSTM storage and the outputs at time τ . At any timestep τ , we have:*

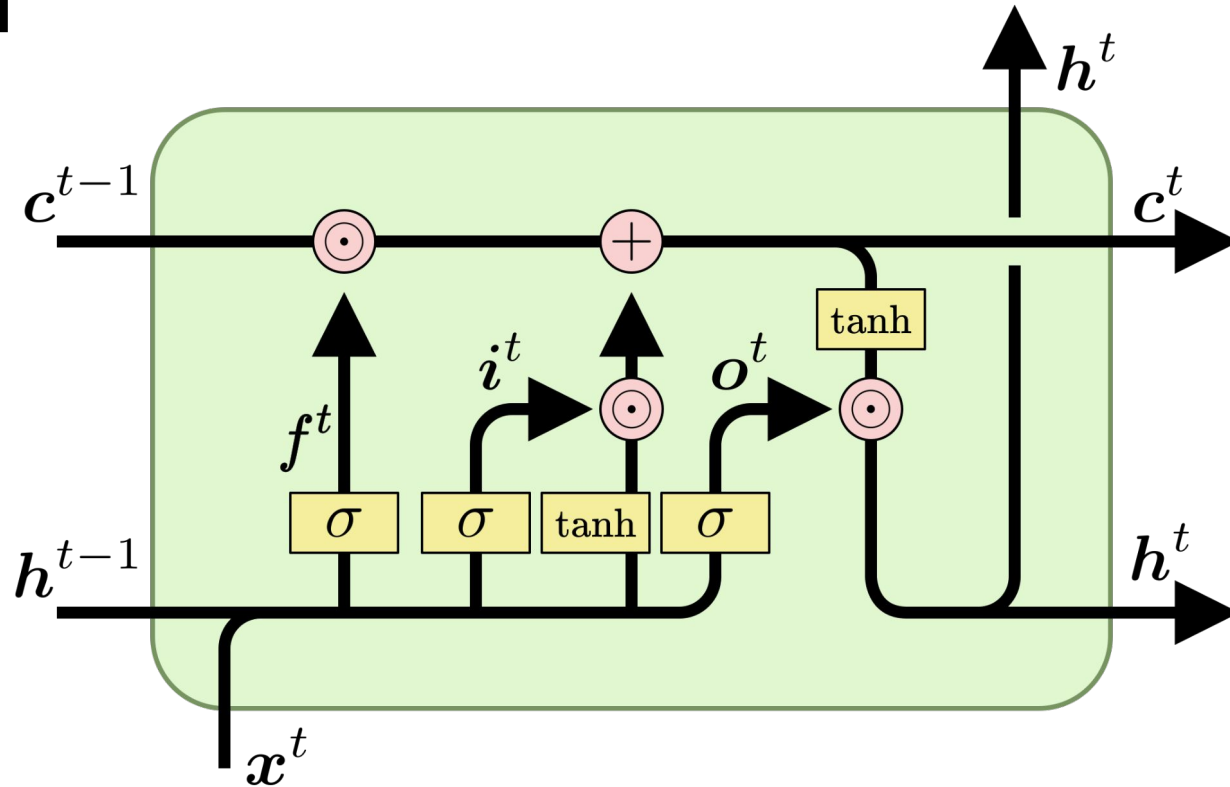
$$m_c^\tau = m_c^0 + \sum_{t=1}^{\tau} x^t - \sum_{t=1}^{\tau} m_h^t. \quad (9)$$

That is, the change of mass in the memory cells is the difference between the input and output mass, accumulated over time.

OUTLINE

- Motivation
- Model
- Experiments

LSTM



MC-LSTM

Total mass

$$m_{\text{tot}}^t = R^t \cdot c^{t-1} + i^t \cdot x^t$$

State mass

$$c^t = (\mathbf{1} - o^t) \odot m_{\text{tot}}^t$$

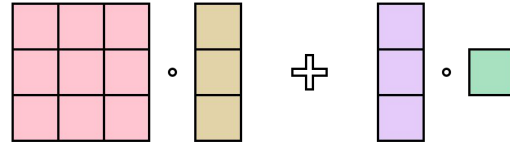
Output mass

$$h^t = o^t \odot m_{\text{tot}}^t.$$

- Cell State
- Mass Input
- Auxiliary Input
- Parameter

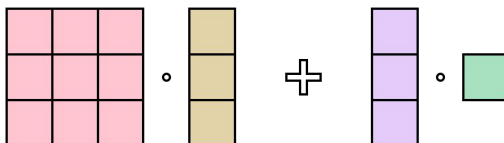
MC-LSTM

Total mass

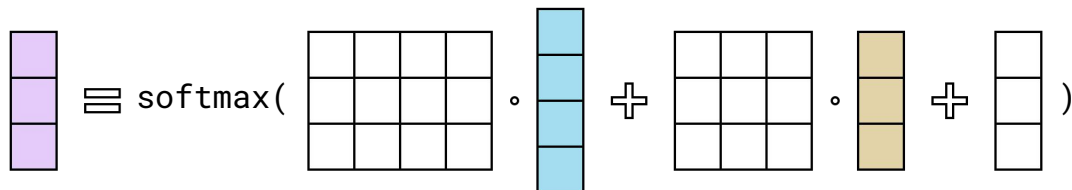


MC-LSTM

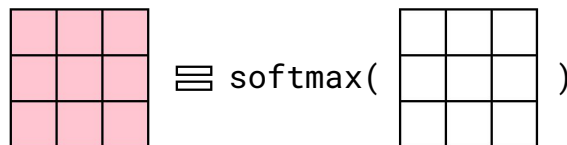
Total mass



Input gate



Redistribution
(static)




MC-LSTM

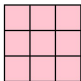
Total mass $m_{\text{tot}}^t = R^t \cdot c^{t-1} + i^t \cdot x^t$

State mass $c^t = (\mathbf{1} - o^t) \odot m_{\text{tot}}^t$

Output mass $h^t = o^t \odot m_{\text{tot}}^t.$

Input gate  $i^t = \text{softmax}(W_i \cdot a^t + U_i \cdot \frac{c^{t-1}}{\|c^{t-1}\|_1} + b_i)$

Output gate $o^t = \sigma(W_o \cdot a^t + U_o \cdot \frac{c^{t-1}}{\|c^{t-1}\|_1} + b_o)$


Redistribution
(static)  $R^t = \text{softmax}(B_r),$

MC-LSTM

Total mass $m_{\text{tot}}^t = R^t \cdot c^{t-1} + i^t \cdot x^t$

State mass $c^t = (1 - o^t) \odot m_{\text{tot}}^t$

Output mass $h^t = o^t \odot m_{\text{tot}}^t.$

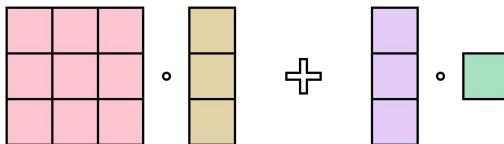
Input gate  $i^t = \text{softmax}(W_i \cdot a^t + U_i \cdot \frac{c^{t-1}}{\|c^{t-1}\|_1} + b_i)$

Output gate $o^t = \sigma(W_o \cdot a^t + U_o \cdot \frac{c^{t-1}}{\|c^{t-1}\|_1} + b_o)$

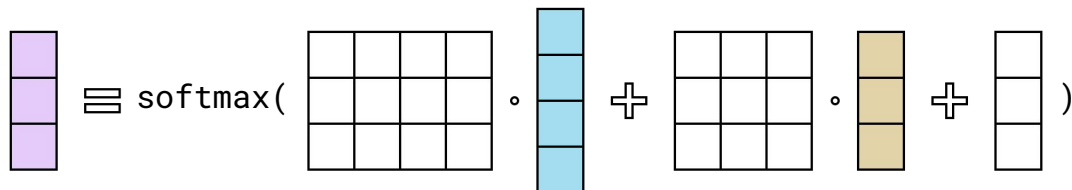
Redistribution (dynamic)  $R^t = \text{softmax}\left(W_r \cdot a^t + U_r \cdot \frac{c^{t-1}}{\|c^{t-1}\|_1} + B_r\right)$

MC-LSTM

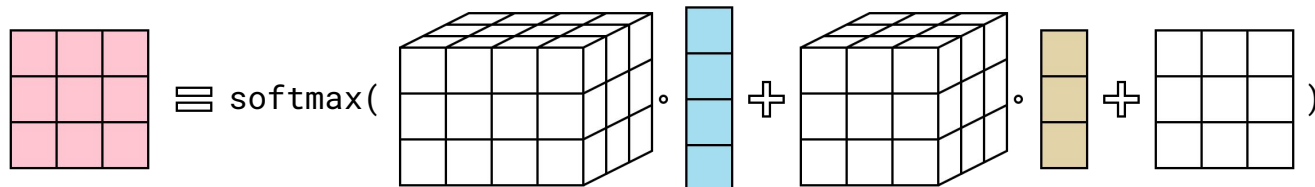
Total mass



Input gate



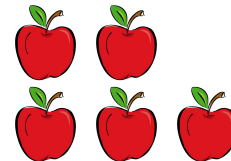
Redistribution
(dynamic)



OUTLINE

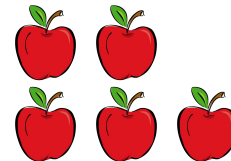
- Motivation
- Model
- Experiments

MC-LSTM ON ARITHMETIC



reference: 0.3, 0.4, 0.2, 0.0, 0.1, 0.4, 0.3, 0.1, 0.5, 0.2

MC-LSTM ON ARITHMETIC



reference: 0.3, 0.4, 0.2, 0.0, 0.1, 0.4, 0.3, 0.1, 0.5, 0.2

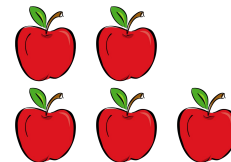
length: 0.3, 0.4, 0.2, 0.0, 0.1, 0.4, 0.3, 0.1, 0.5, 0.2, 0.3, ...

range: 3.4, 4.4, 2.2, 0.1, 1.2, 4.4, 3.0, 1.5, 5.1, 2.8

count: 0.3, 0.4, 0.2, 0.0, 0.1, 0.4, 0.3, 0.1, 0.5, 0.2

combo: 3.4, 4.4, 2.2, 0.1, 1.2, 4.4, 3.0, 1.5, 5.1, 2.8, 3.7, ...

MC-LSTM ON ARITHMETIC



	reference ^a	seq length ^b	input range ^c	count ^d	combo ^e	NaN ^f
MC-LSTM	0.004 \pm 0.003	0.009 \pm 0.004	0.8 \pm 0.5	0.6 \pm 0.4	4.0 \pm 2.5	0
LSTM	0.008 \pm 0.003	0.727 \pm 0.169	21.4 \pm 0.6	9.5 \pm 0.6	54.6 \pm 1.0	0
NALU	0.060 \pm 0.008	0.059 \pm 0.009	25.3 \pm 0.2	7.4 \pm 0.1	63.7 \pm 0.6	93
NAU	0.248 \pm 0.019	0.252 \pm 0.020	28.3 \pm 0.5	9.1 \pm 0.2	68.5 \pm 0.8	24

^a training regime:

summing 2 out of 100 numbers between 0 and 0.5.

^b longer sequence lengths:

summing 2 out of 1 000 numbers between 0 and 0.5.

^c more *mass* in the input:

summing 2 out of 100 numbers between 0 and 5.0.

^d higher number of summands:

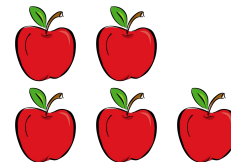
summing 20 out of 100 numbers between 0 and 0.5.

^e combination of previous scenarios:

summing 10 out of 500 numbers between 0 and 2.5.

^f Number of runs that did not converge.

MC-LSTM ON ARITHMETIC

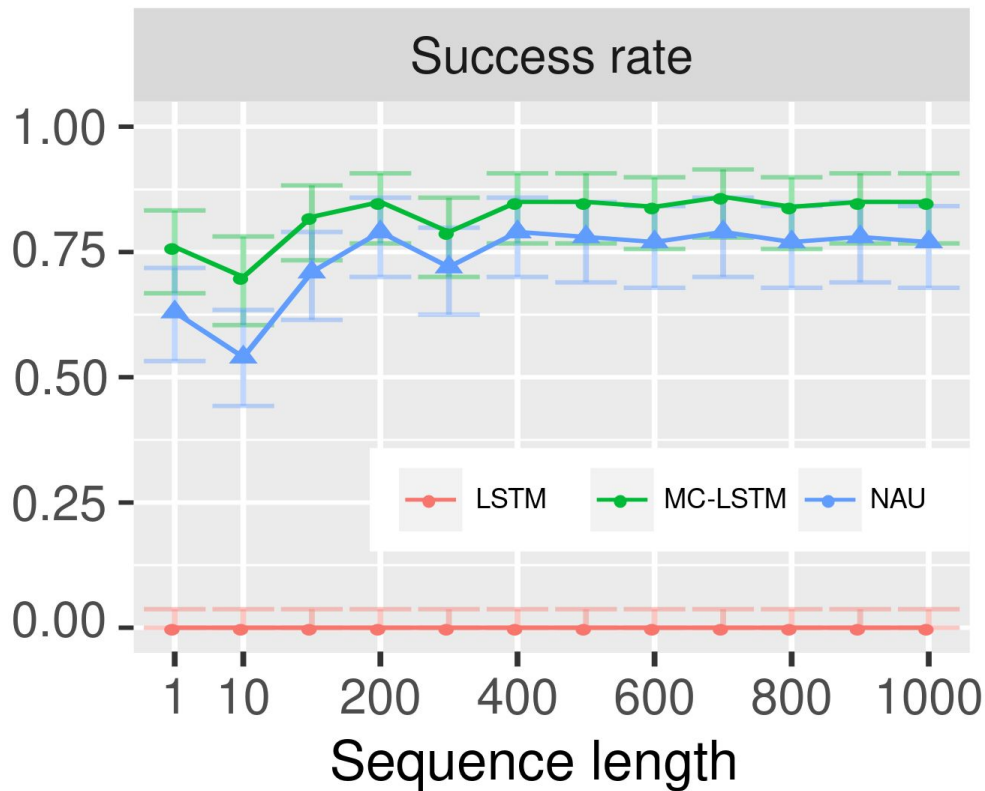
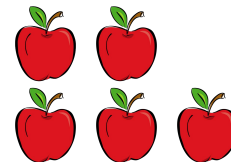


Input:

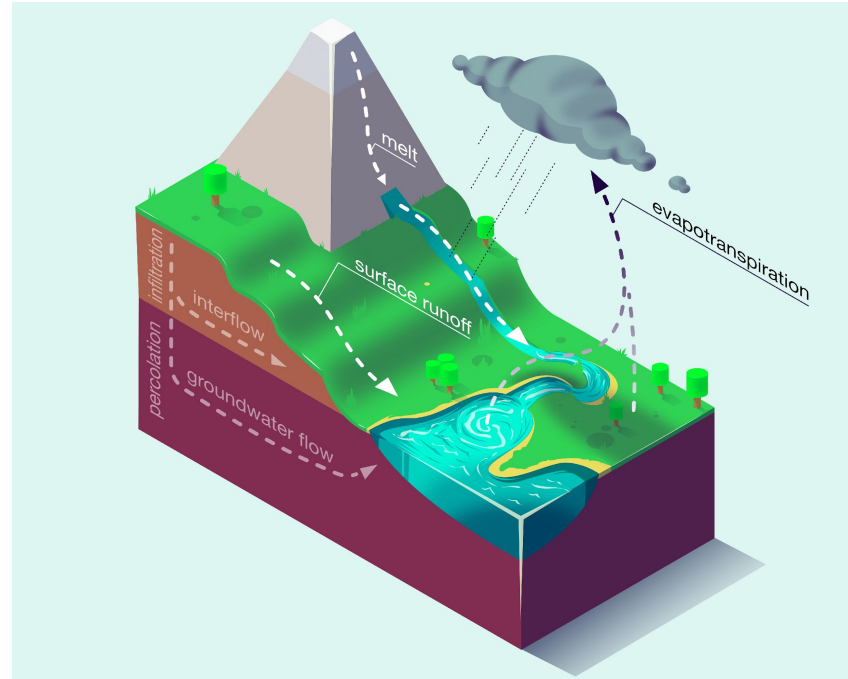
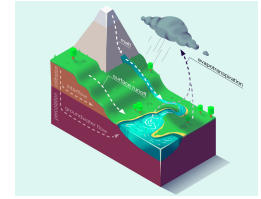


Output: $26 = 4 + 8 + 1 + 5 + 1 + 7$

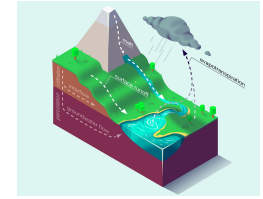
MC-LSTM ON ARITHMETIC



MC-LSTM ON HYDROLOGY



MC-LSTM ON HYDROLOGY



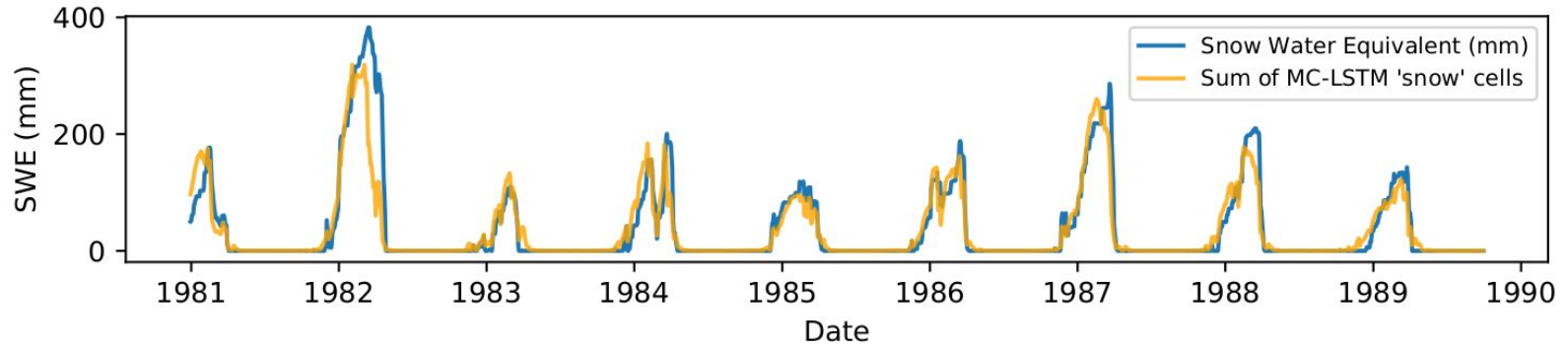
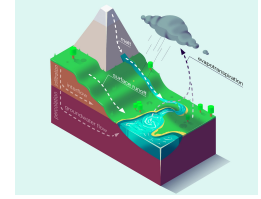
Model	MC ^a	FHV ^b	NSE ^c
MC-LSTM	✓	-14.7 _{-7.0 -23.4}	0.744 _{0.814 0.641}
LSTM	✗	-15.7 _{-8.6 -23.8}	0.763 _{0.835 0.676}
mHM	✓	-18.6 _{-9.5 -27.7}	0.666 _{0.730 0.588}
...
HBVub	✓	-18.5 _{-8.5 -27.8}	0.676 _{0.749 0.578}

^a: Mass conservation (MC).

^b: Top 2% peak flow bias: $(-\infty, \infty)$, values closer to zero are desirable.

^c: Nash-Sutcliffe Efficiency: $(-\infty, 1]$, values closer to one are desirable.

MC-LSTM ON HYDROLOGY



TL ;DR: MC-LSTM

- LSTM + inductive bias
- stochastic matrices for conservation
- Better generalisation
- Cell states easier to interpret