

Machine learning for RNA (secondary) structure prediction

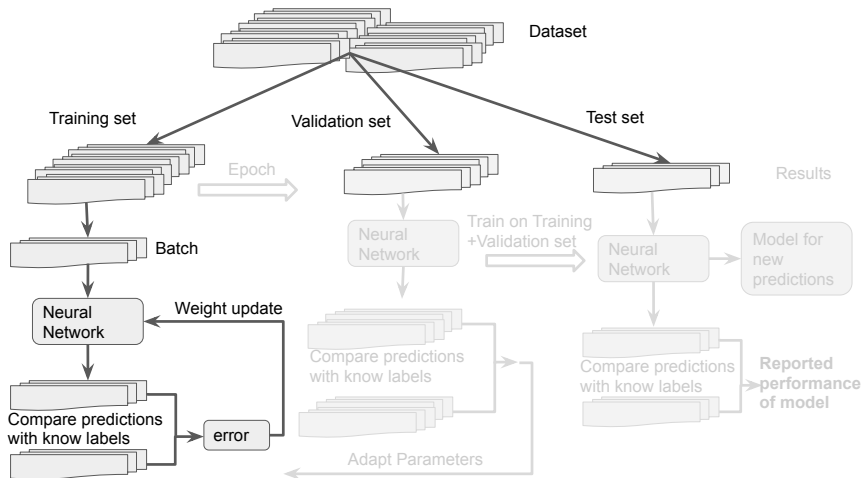
Julia Wielach



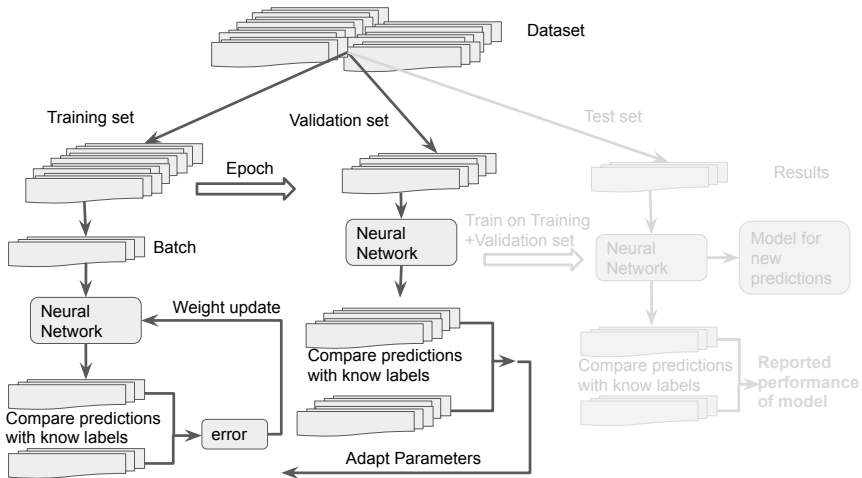
12th February 2020

- ▷ Machine Learning
 - Training Process
 - Network Types
- ▷ ML for RNA structure prediction
 - Introduction
 - Aim of Project
 - Input
 - Output
 - Networks

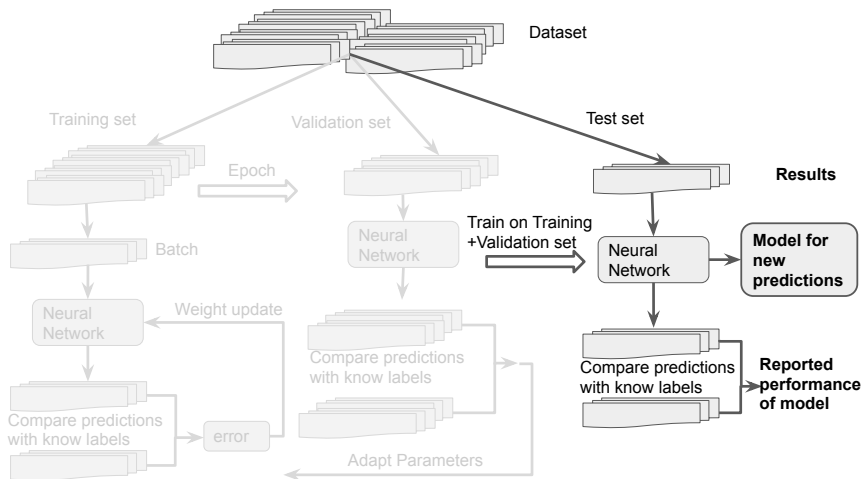
Training



Validation



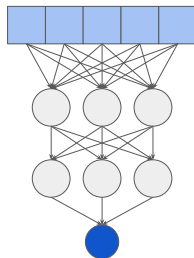
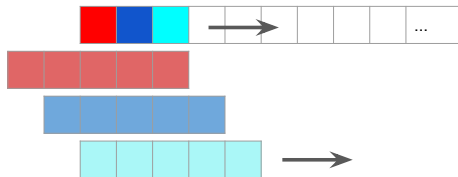
Testing



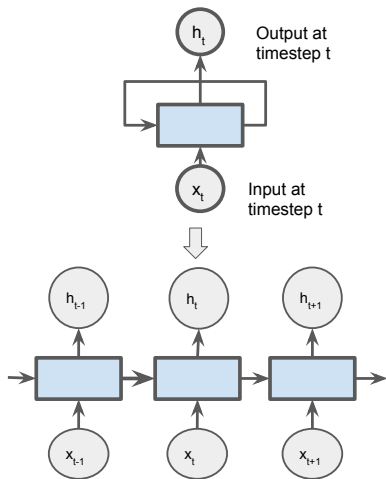
Sliding Window

Benchmark

Only local dependencies (within window)

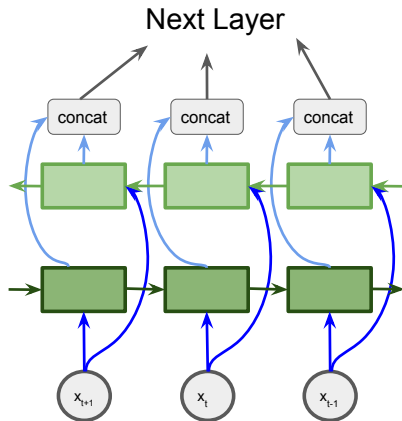


Recurrent Neural Network (RNN)



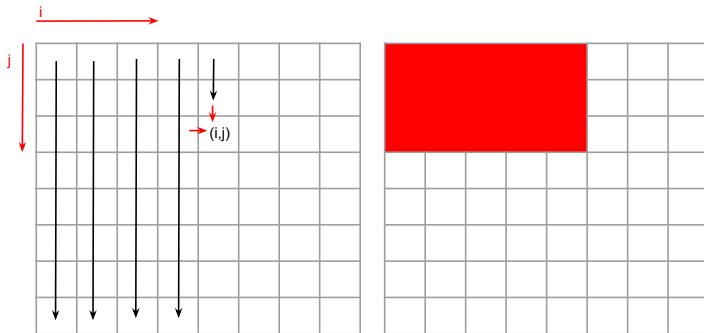
Graphic based on [colah's blog](#)

Bidirectional RNN



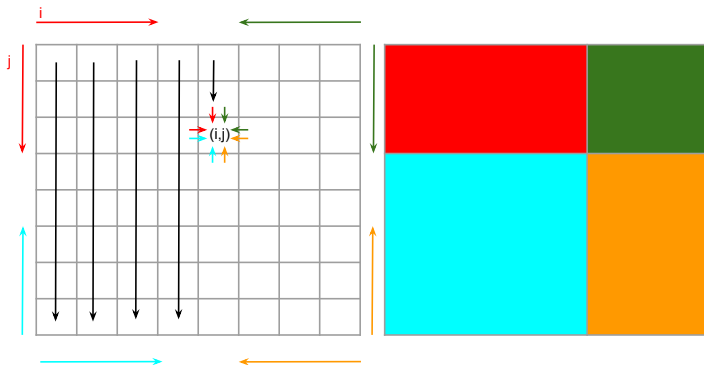
Multidimensional RNN (2D)

Input from two directions instead of one
(i,j) only reached after (i,j-1) and (i-1,j)



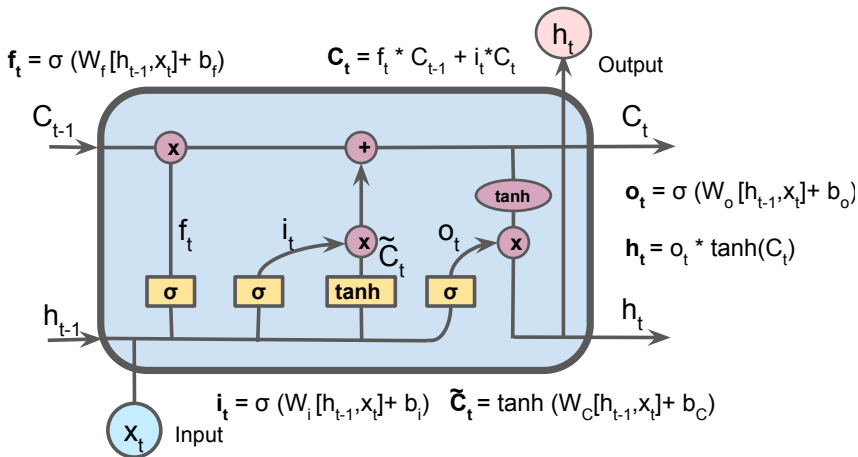
Multidimensional RNN (2D)

At every point context from 4 directions



Graves A., Fernández S., Schmidhuber J. (2007)

Long Short Term Memory (LSTM) Network



Hochreiter, S. & Schmidhuber, J. (1997)

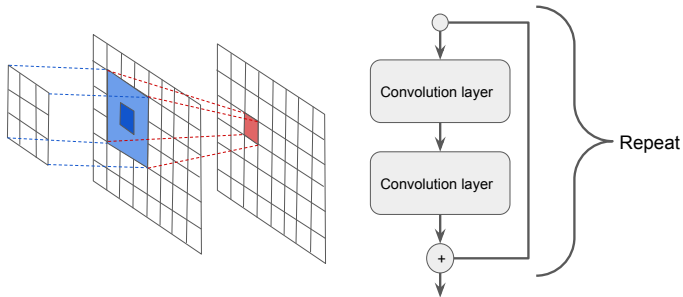
Graphic based on [colah's blog](#)

Convolution NNs & ResNets

Capture local surroundings

Problems with very deep networks → performance decrease

Skip Connections



What?

```
AACGCAUUGGAUACCUGUGUAUGAUUAUUACACGGUAGAGUACGCGCUCGCGGA
                                ↓
001110000001111011111000000011111111111110000111111100
..(((.....(((.((((.....)))))))))(((.....))))))..
```

Why?

- ▷ Used for similar tasks
- ▷ Possibly capture complicated dependencies
- ▷ Already used in publications
- ▷ No performance comparisons

Current situation + project

- ▷ First Networks published
 - Very good performances reported
 - No comparisons between new models
 - Possible bias introduced by datasets
 - Predictions specific for different types of RNAs
- ▷ Tests with artificial random RNA sequences predicted by RNAfold
 - Similar performance reachable ?
 - Dataset size dependency
 - Capacity dependency

One hot encoding

A → [1, 0, 0, 0]
C → [0, 1, 0, 0]
G → [0, 0, 1, 0]
U → [0, 0, 0, 1]



ACUGAC ...

↓ ↓ ↓ ↓ ↓ ↓ ...

1 0 0 0 1 0 ...
0 1 0 0 0 1 ...
0 0 0 1 0 0 ...
0 0 1 0 0 0 ...

Input

Concatenation of one hot encodings
Matrix of size $L \times L \times 8$

	A	C	U	...
A	$\begin{bmatrix} 1, 0, 0, 0 \\ , 1, 0, 0, 0 \end{bmatrix}$	$\begin{bmatrix} 0, 1, 0, 0 \\ , 1, 0, 0, 0 \end{bmatrix}$	$\begin{bmatrix} 0, 0, 0, 1 \\ , 1, 0, 0, 0 \end{bmatrix}$	
C		$\begin{bmatrix} 0, 1, 0, 0 \\ , 0, 1, 0, 0 \end{bmatrix}$	$\begin{bmatrix} 0, 0, 0, 1 \\ , 0, 1, 0, 0 \end{bmatrix}$	
U			$\begin{bmatrix} 0, 0, 0, 1 \\ , 0, 0, 0, 1 \end{bmatrix}$	
...				

Input

Matrix with value "rating" basepair
Matrix of size $L \times L \times 1$

	A	C	G	U
A	0	0	0	2
C		0	3 +2 * b	0
G			0	0.8
U				0

Diagram illustrating a matrix with values and annotations:

- The cell containing '2' (row A, column U) is highlighted with a cyan border.
- The cell containing '3 + 2 * b' (row C, column G) is highlighted with a green border.
- The cell containing '0.8' (row G, column U) is highlighted with a purple border.
- An arrow points from the '3 + 2 * b' cell to the '2' cell, labeled '+3 * b'.

Output

- ▷ Problem of incorrect secondary structure

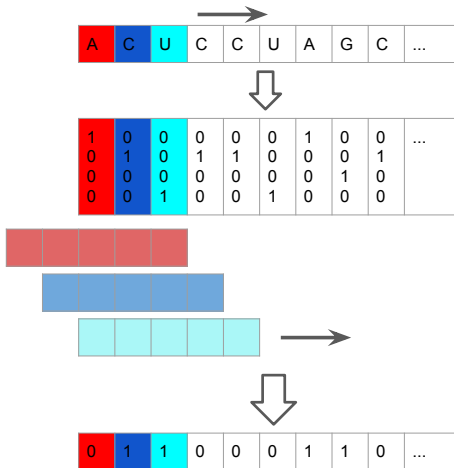
Incorrect pairing partners

Opening and closing brackets → different number

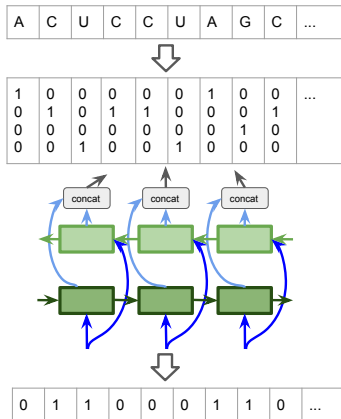
Separate label for every base → conflicting basepairs
eg. $i+j$, $j+k$

- ▷ Classify additional interactions (Pseudoknots, Triple and Non-canonical basepairs)
- ▷ Additional Corrections
- ▷ Simplification (1/0 for Paired/Unpaired)

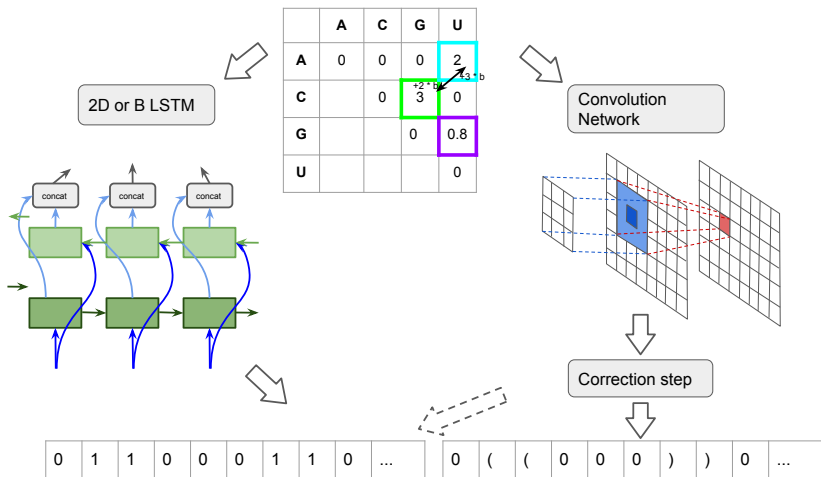
Sliding Window



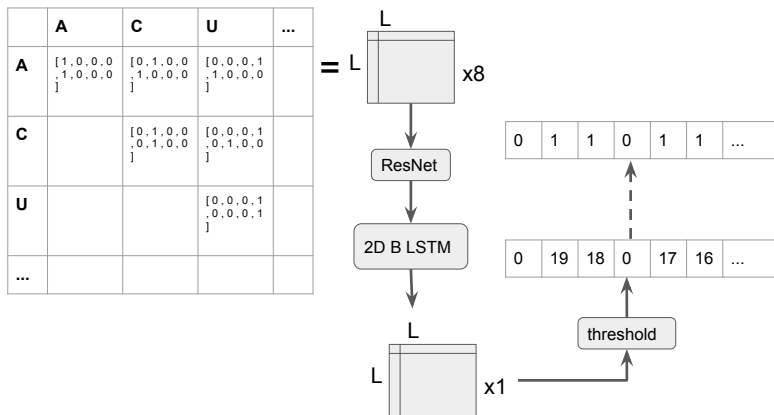
B-LSTM



LSTM / CNN



ResNet & 2D BLSTM



Acknowledgements:

- ▷ Christoph Flamm
- ▷ Ivo Hofacker
- ▷ Michael Wolfinger

- ▷ Gregor Entzian
- ▷ Irene Katharina Beckmann
- ▷ Maria Waldl



<https://xkcd.com/1838/>

TBI



universität
wien

Supplement 1 (Metrics)

	Actually Positive	Actually negative	
Predicted Positive	True positive (T_P)	False positive (F_P)	Precision/Positive Predictive Value: $T_P / T_P + F_P$
Predicted Negative	False negative (F_N)	True negative (T_N)	Negative Predictive Value: $T_N / T_N + F_N$
	Sensitivity: $T_P / T_P + F_N$	Specificity: $T_N / T_N + F_P$	

$$F_1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

$$\text{Accuracy} = \frac{T_P + T_N}{T_P + T_N + F_P + F_N}$$

$$\text{MCC} = \frac{T_P \times T_N - F_P \times F_N}{\sqrt{(T_P + F_P)(T_P + F_N)(T_N + F_P)(T_N + F_N)}}$$