

Recurrent Neural Networks

15th November, 2022



Alexey Zaytsev,

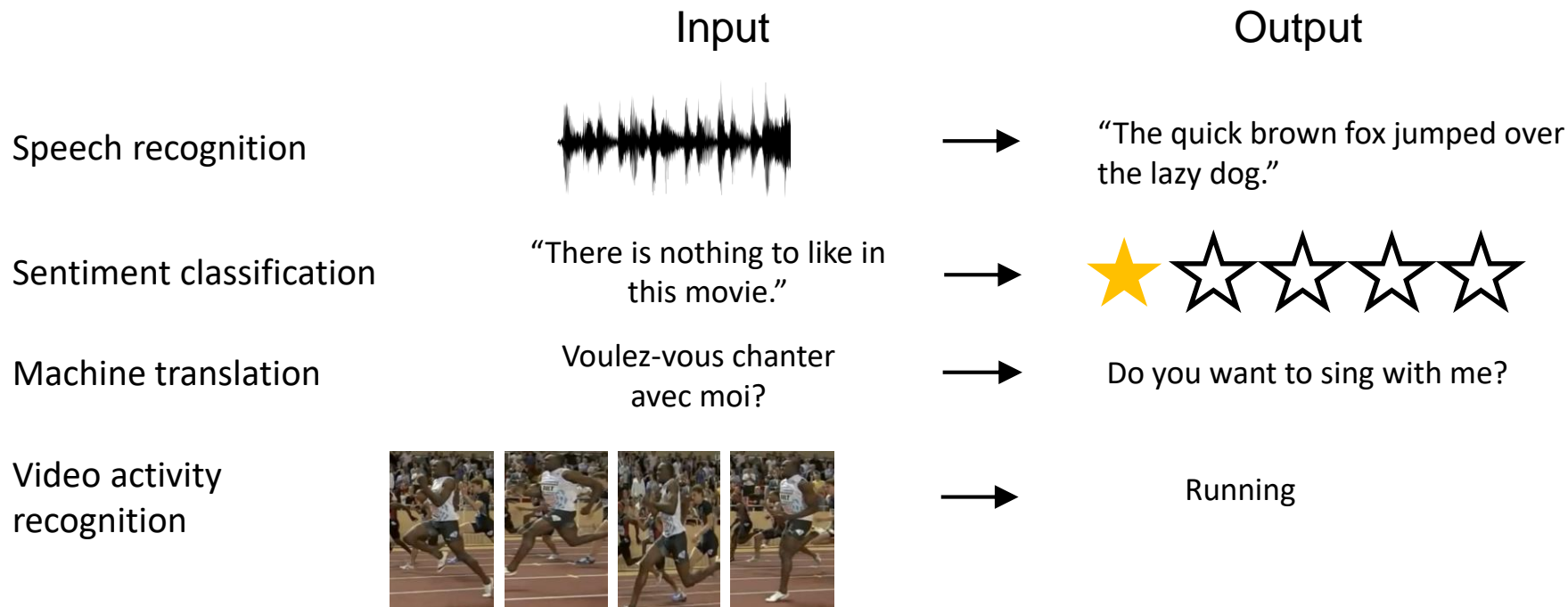
Assistant professor, Skoltech

Common ways for classic ML application for time series data

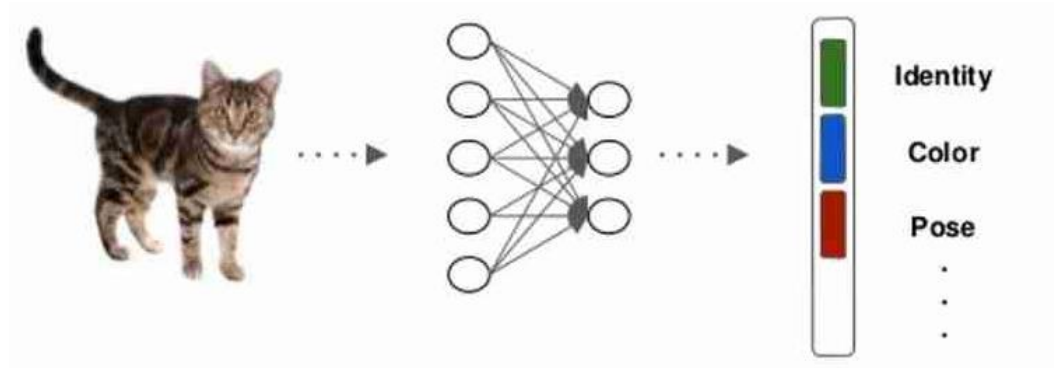
1. Take input data including history for the target variable
 2. Add differences, combinations, rolling means, medians, etc.
 3. Add one-hot-encoding for important categorical features (day of week, holiday or not)
- Now we have input features for all points
 - Let's apply our favorite ML regression algorithm

dmlc
XGBoost

Deep Learning problems with sequential data: we need representations



One example problem: why do we need Neural networks?



Textbook example: next word prediction

The most complicated and difficult part of it was only just beginning.

Textbook example: next word prediction

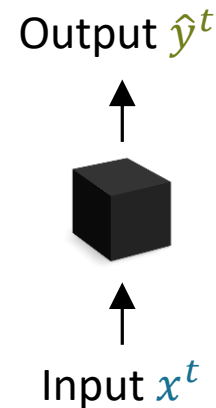
Idea 1: use previous word(s)

Problem 1: long-term dependencies

“**France** is where I grew up, but I now live in Boston. I speak fluent ____.”

The most complicated and difficult part of it was only just beginning.

Feature representation: [0, 0, 0, 1, 0, 0]



Textbook example: next word prediction

Idea 1: use previous word(s)

Problem 1: long-term dependencies

Idea 2: use bag of words model

The most complicated and difficult part of it was only just beginning.



Feature representation: [0, 3, 0, 2, 0, 0]

Bag of words: number of occurrences of each word

Textbook example: next word prediction

Idea 1: use previous word(s)

Idea 2: use bag of words model

Problem 1: long-term dependencies

Problem 2: order preservation

The food was good, not bad at all.

vs.

The food was bad, not good at all.

The most complicated and difficult part of it was only just beginning.

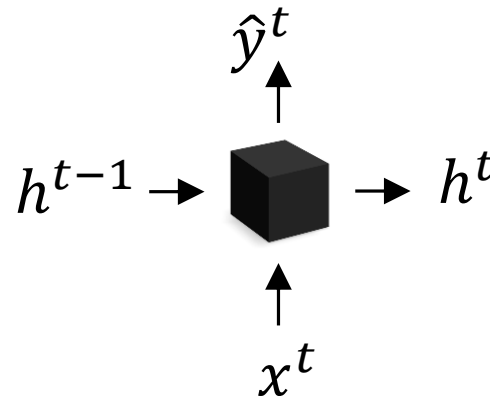


Feature representation: [0, 3, 0, 2, 0, 0]

Bag of words: number of occurrences of each word
(see also TF-IDF features)

Model Design Criteria

1. Variable-length sequences processing
2. Long-term memory
3. Maintain order information
4. Natural preprocessing



Recurrent Neural Networks are the solution!

Sequence processing with classic ML models

1. Variable-length sequences processing

YES (if one to one)

2. Long-term memory

NO 

3. Maintain order information

NO 

4. Natural preprocessing

a kind of

Sequence data examples

Speech recognition



Sentiment classification

“There is nothing to like in
this movie.”

DNA sequence analysis

AGCCCCTGTGAGGAACTAG

Machine translation

Voulez-vous chanter
avec moi?

Video activity
recognition



Large data set

Semi-
structured data

+

+

+

+

+

+

+

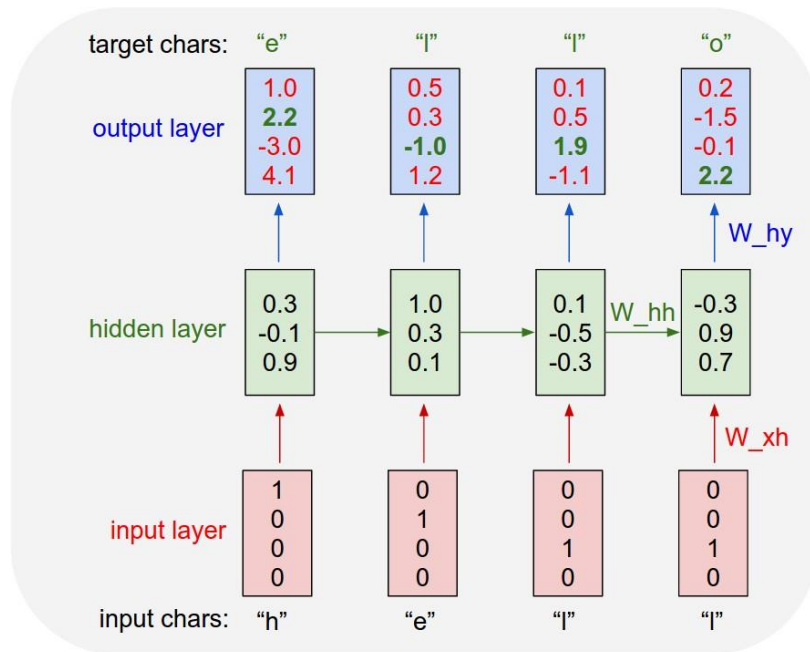
+

+

+

Examples of texts generated by LSTM (Long Short Term Memory NN)

- Shakespeare
- Wiki
- Algebraic geometrics articles
- Linux Source Code
- Dinosaurs names



For $\bigoplus_{i=1}^n \mathcal{O}_{X_i}$ where $\mathcal{L}_{X_i} = 0$, hence we can find a closed subset \mathcal{H} in \mathcal{H} and any sets \mathcal{F} on X , U is a closed immersion of S , then $U \rightarrow T$ is a separated algebraic space.

Proof. Proof of (1). It also start we get

$$S = \text{Spec}(R) = U \times_X U \times_X U$$

and the comparably in the fibre product covering we have to prove the lemma generated by $\coprod_{i \in I} U_i \rightarrow V$. Consider the maps M along the set of points Sch_{aff} and $U \rightarrow U$ is the fibre category of S in U in Section 77 and the fact that any U affine, we Morphisms, Lemma 77. Hence we obtain a scheme S and any open subset $W \subset U$ in $\text{Sch}(G)$ such that $\text{Spec}(R) \rightarrow S$ is smooth or an

$$U = \bigcup U_i \times_S U_i$$

which has a nonzero morphism we may assume that f_i is of finite presentation over S . We claim that $\mathcal{O}_{X_{\mathcal{F}}}$ is a scheme where $x, x', x'' \in S'$ such that $\mathcal{O}_{X_{\mathcal{F}}'} \rightarrow \mathcal{O}_{X_{\mathcal{F}}''}$ is separated. By Algebra, Lemma 77 we can define a map of complexes $\text{GL}_n(x'/x'')$ and we win. \square

To prove study we see that \mathcal{F}_{ij} is a covering of X' , and T_i is an object of $\mathcal{F}_{X_{ij}}$ for $i > 0$ and \mathcal{F}_p exists and let \mathcal{F}_i be a preabel of $\mathcal{O}_{X_{ij}}$ -modules on C as a \mathcal{F} -module. In particular $\mathcal{F} = U/\mathcal{F}$ we have to show that

$$\hat{\mathcal{F}}^* = \mathcal{F}^* \otimes_{\mathcal{O}_{X_{ij}}} \mathcal{O}_{X_{ij}} = \mathcal{F}^* \otimes \mathcal{F}$$

is a unique morphism of algebraic stacks. Note that

$$\text{Atroves} = (\text{Sch}/S)_{\text{aff}}^{\text{op}} / (\text{Sch}/S)_{\text{aff}}$$

and

$$V = \Gamma(S, \mathcal{O}) \rightarrow \Gamma(U, \text{Spec}(A))$$

is an open subset of X . Thus U is affine. This is a continuous map of X is the inverse, the groupoid scheme S .

Proof. See discussion of sheaves of sets. \square

The result for prove any open covering follows from the less of Example 77. It may replace S by $X_{\text{pro-étale}}$ which gives an open subspace of X and T equal to $S_{2\text{-étale}}$ see Descent, Lemma 77. Namely, by Lemma 77 we see that R is geometrically regular over S .

Lemma 0.1. Assume (1) and (2) by the construction in the description.

Suppose $X = \text{lim}[X]$ (by the formal open covering X and a single map $\text{Proj}_X(A) = \text{Spec}(B)$ over U compatible with the complex

$$\text{Set}(A) = \Gamma(X, \mathcal{O}_{X, \mathcal{O}_{X, \mathcal{O}_X}}).$$

When in this case of to show that $Q \rightarrow G_{2, X}$ is stable under the following result in the second conditions of (1), and (2). This finishes the proof. By Definition 77 (without element is when the closed subschemes are categorical. If T is surjective we may assume that T is connected with residue fields of S . Moreover there exists a closed subspace $Z \subset X$ of X where U is proper (some defining as a closed subset of the uniqueness it suffices to check the fact that the following theorem (1) I is locally of finite type. Since $S = \text{Spec}(R)$ and $Y = \text{Spec}(R)$.

Proof. This is form all sheaves of sheaves on X . But given a scheme U and a surjective étale morphism $U \rightarrow X$. Let $U \cap U_i = \bigcup_{j=1}^n U_{ij}$ be the scheme X over S at the schemes $X_i \rightarrow X$ and $U = \text{lim}_i U_i$. \square

The following lemma surjective restrocomposes of this implies that $\mathcal{F}_{X_i} = \mathcal{F}_{X_i} = \mathcal{F}_{X_i}$.

Lemma 0.2. Let X be a locally Noetherian scheme over S , $E = \mathcal{F}_{X_{ij}}$. Set $Z = \mathcal{F}_i \subset \mathcal{F}_i$. Since $\mathcal{F}^* \subset \mathcal{F}^*$ are nonzero over $u_0 \leq p$ is a subset of $\mathcal{F}_{i,0} \circ \mathcal{F}_i$ works.

Lemma 0.3. In Situation 77. Hence we may assume $q = 0$.

Proof. We will use the property we see that p is the next functor (77). On the other hand, by Lemma 77 we see that

$$D(\mathcal{O}_X) = \mathcal{O}_X(D)$$

where K is an F -algebra where δ_{i+1} is a scheme over S . \square

Conclusions

- Classic ML can't handle semi-structured data common in sequential data processing
- We should *learn representations* via Neural Networks
- Results are nice even for relatively simple models

Problem statements

We want to model the conditional distribution:

$$p(y_t | y_1, \dots, y_{t-1}, x_1, \dots, x_t)$$



Our goal would be to propose such model

Problem statements

We also look at similar models and problems.

(a) Many-to-one



Classification

(b) One-to-one



Anomaly detection

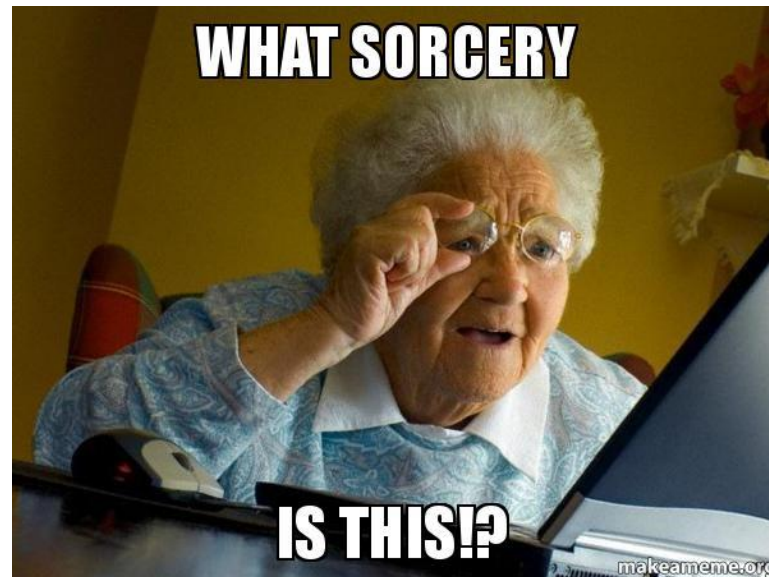
(c) Many-to-many



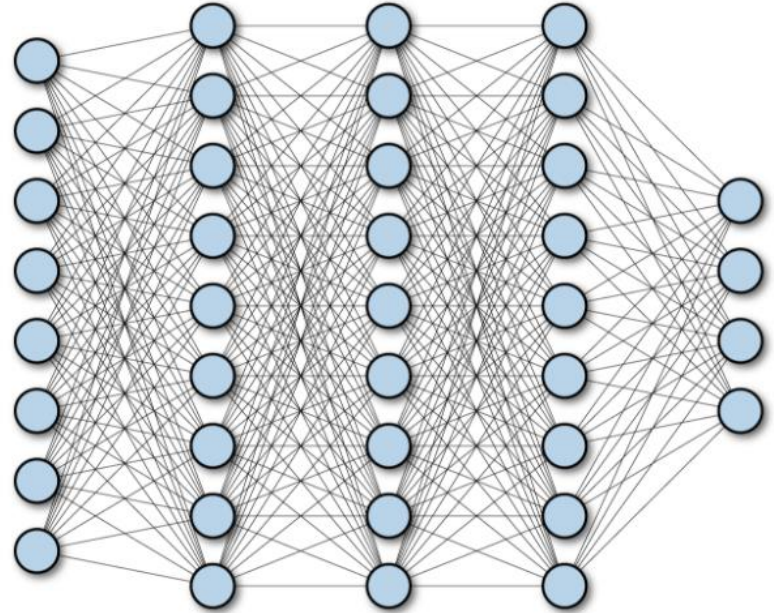
Machine translation

Model Design Criteria

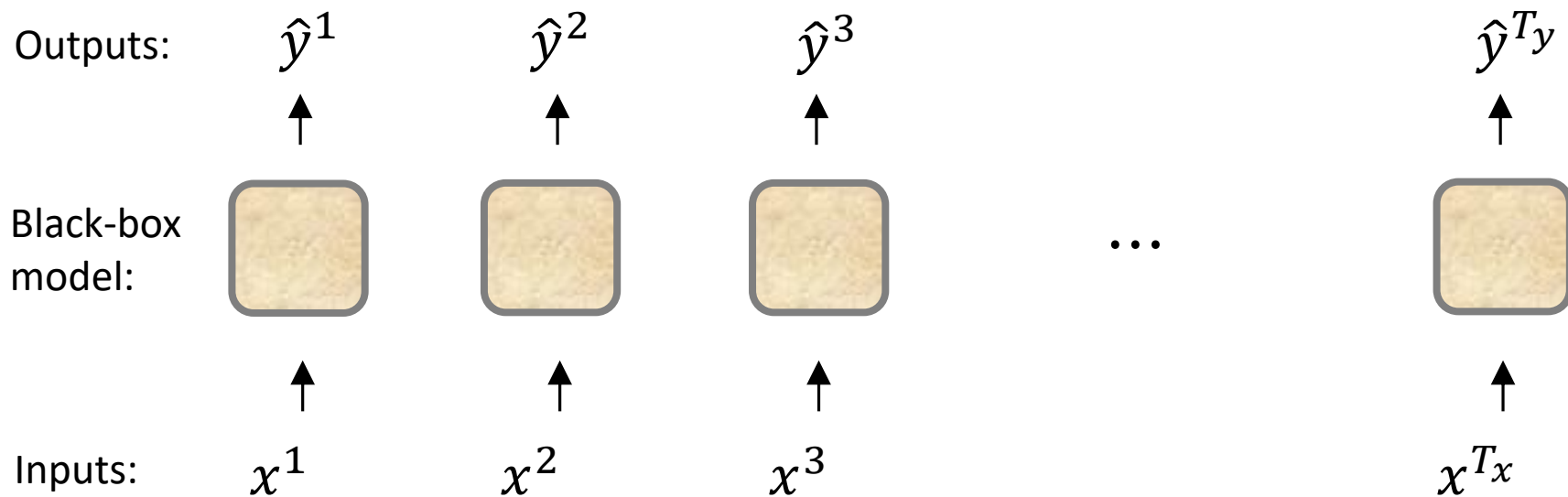
1. Variable-length sequences processing
2. Long-term memory
3. Maintain order information
4. Natural preprocessing
5. **General end2end approach**



Fully connected neural networks?



Separate Fully-Connected Neural Networks or other separate models



We model $p(y_t|x_t)$ instead of $p(y_t|y_1, \dots, y_{t-1}, x_t)$

Why not a standard fully-connected network?



- Inputs, outputs can be different lengths in different examples
- Huge number of parameters (for length 512 it is about 256 000)
- Doesn't share features learned across different positions of text
- Not easy to separate contribution from different y -s ...*think about it as a problem to think about e.g. $p(y$*

SOTA options

Recurrent Neural
Networks

Today

SOTA in some
problem

One dimensional
convolutions (1D
CNN)

Next lecture (also
see ROCKET)

SOTA in another
problems

Transformers

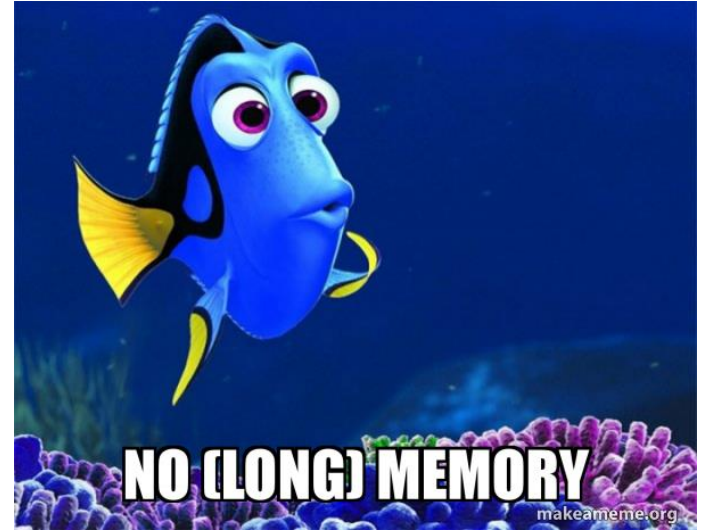
After them

SOTA in most NLP
problems

SOTA results for scoring problem

		Micro Precision
Two RNN architectures	LSTM	<u>0.762</u>
	GRU	0.765
	1D CNN	0.745
	Transformer	0.755
	Gradient boosting	0.697

A vanilla Recurrent Neural Network



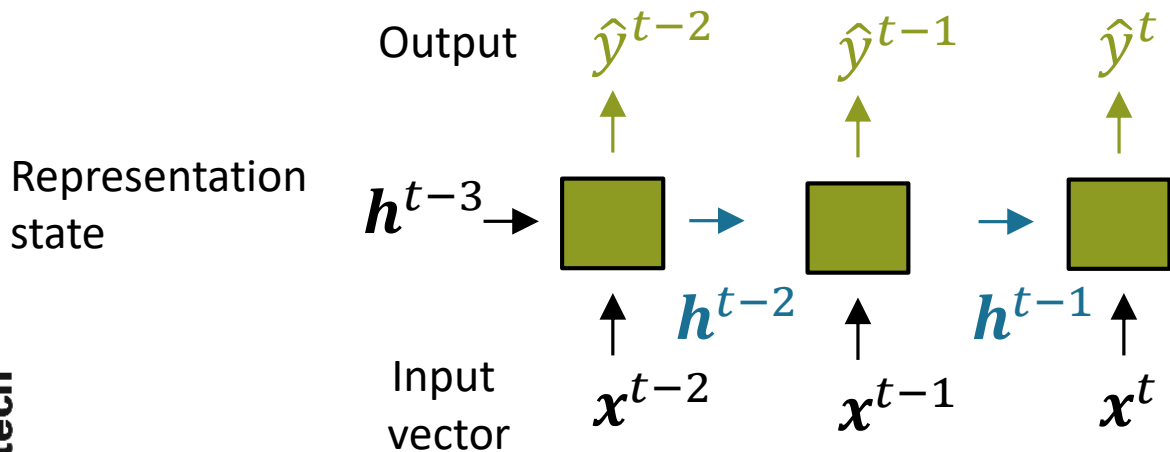
Simple RNN model

At each step:

\hat{y}^t - output / model prediction

\mathbf{x}^t - input vector / new information

\mathbf{h}^t - cell / hidden state



Main idea:

The model inside dark green block is the same for all time moments!

$$p(y_t | y_1, \dots, y_{t-1}, \mathbf{x}_1, \dots, \mathbf{x}_t) = p(y_t | \mathbf{h}_{t-1}, \mathbf{x}_t) = p(y_t | \mathbf{h}_t)$$

Simple RNN model block

General form:

$$\mathbf{h}^t = f_h(\mathbf{x}^t, \mathbf{h}^{t-1})$$

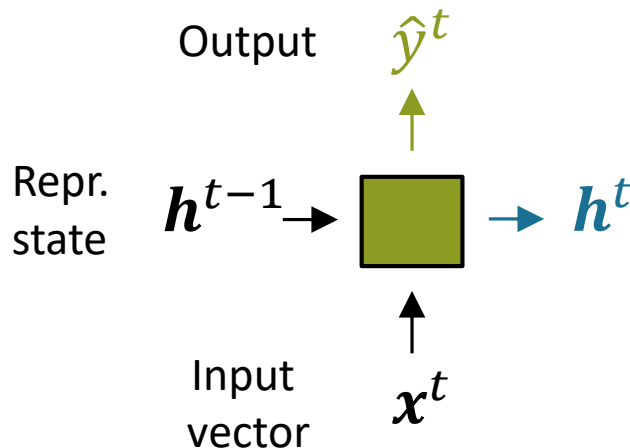
$$\hat{\mathbf{y}}^t = f_y(\mathbf{h}^t)$$

Vanilla RNN model for classification:

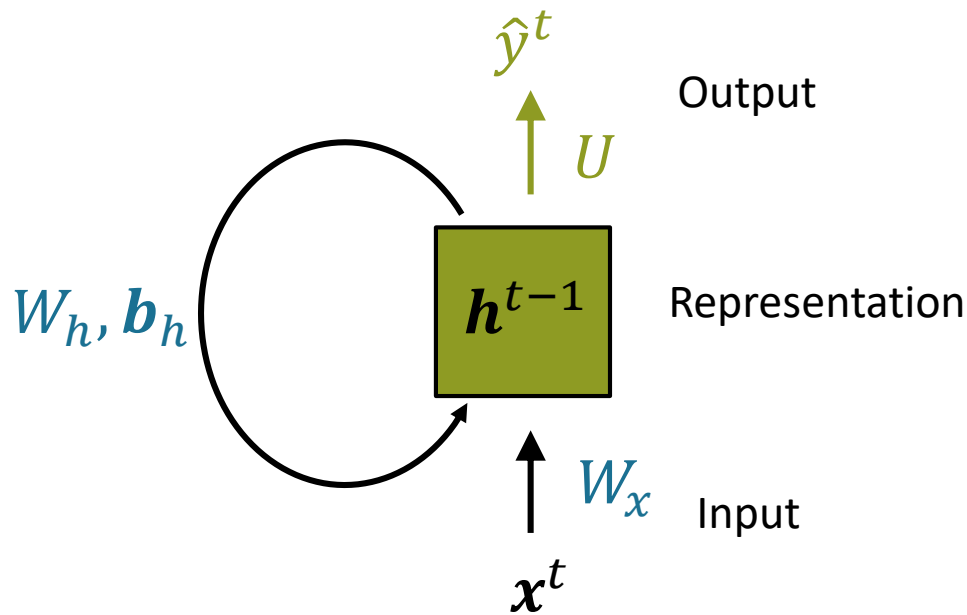
$$\mathbf{h}^t = \tanh(V\mathbf{x}^t + W\mathbf{h}^{t-1} + \mathbf{b}_h)$$

$$\hat{\mathbf{y}}^t = \text{softmax}(U\mathbf{h}^t + b_y)$$

Parameters: V , W , U



Compressed form of RNN



$$h^t = \tanh(Vx^t + Wh^{t-1} + b_h)$$

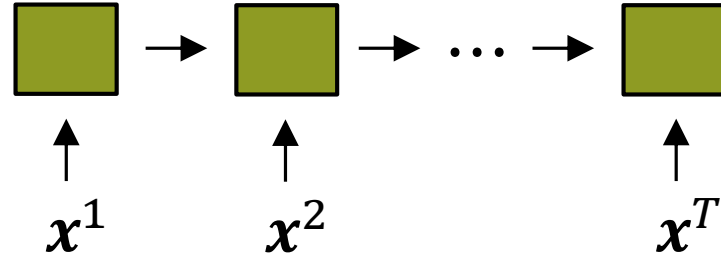
$$\hat{y}^t = \text{softmax}(Uh^t + b_y)$$

Other problems are solved as well

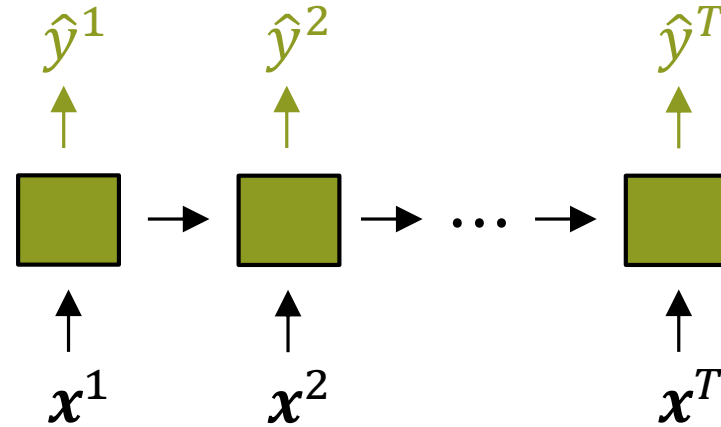


No recurrent structure

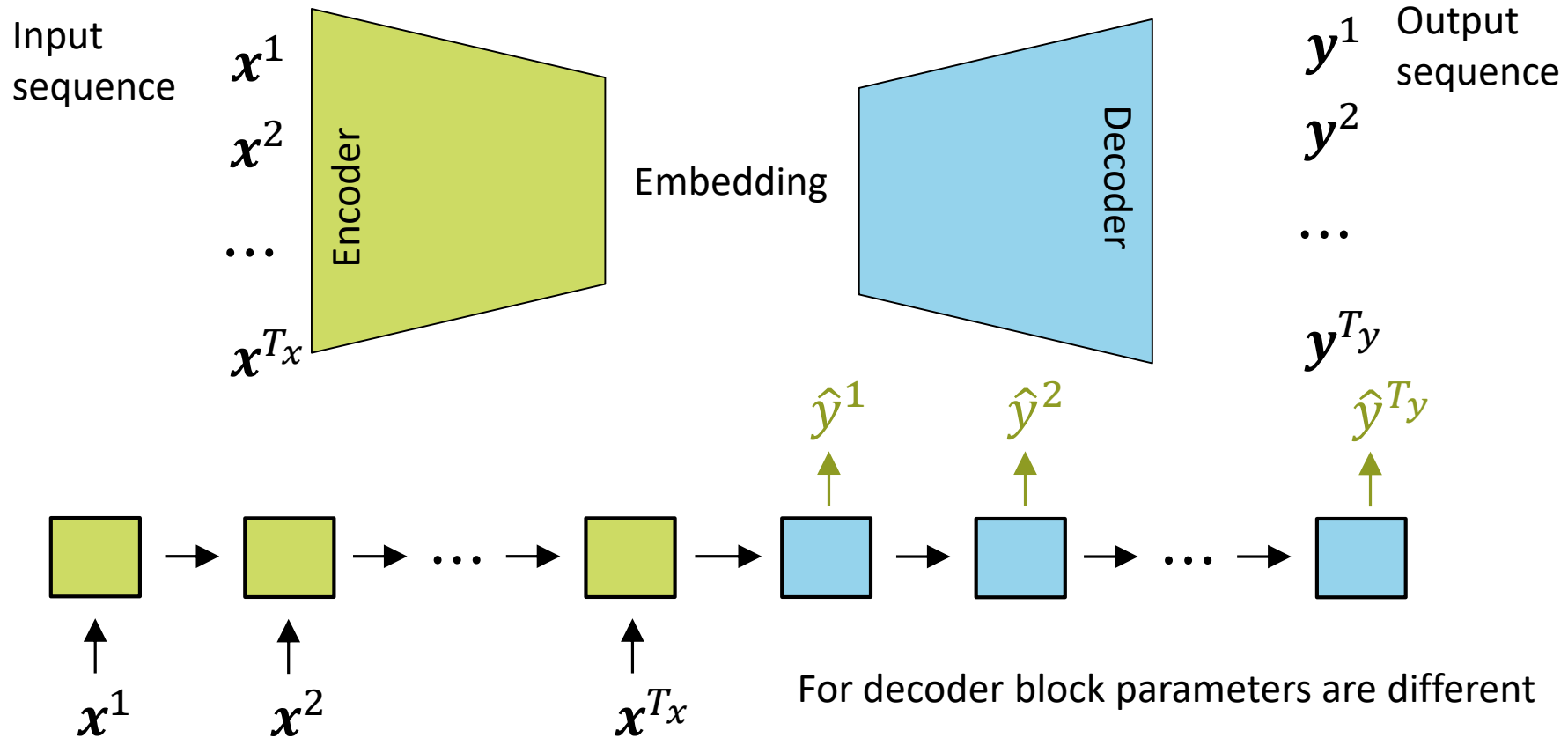
Sequence to a single output



Sequence to a sequence of outputs



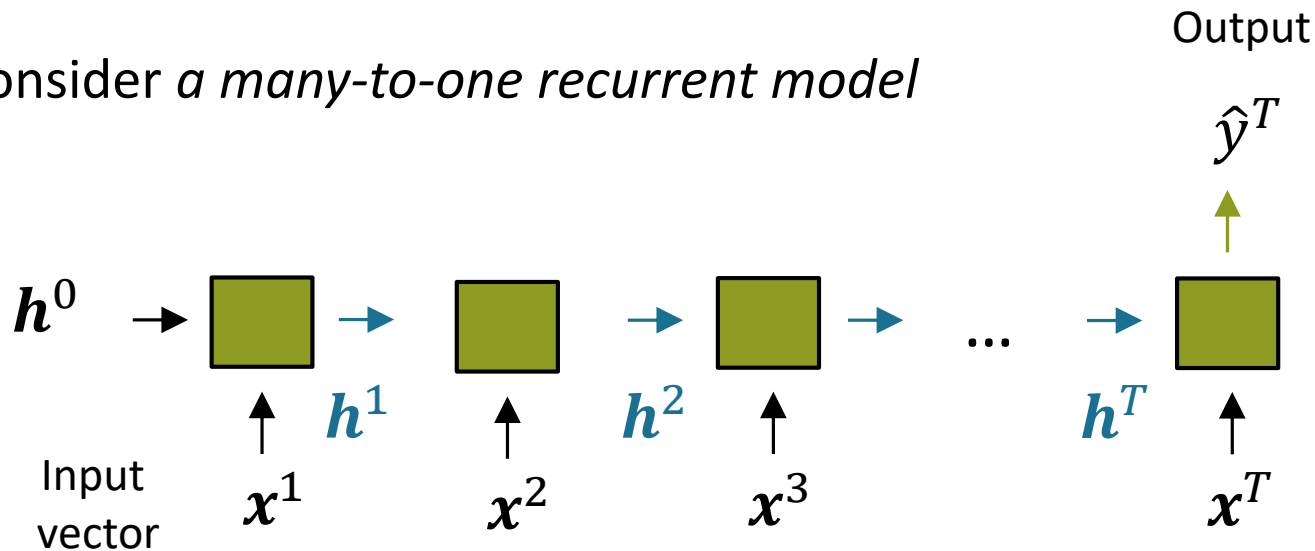
A sequence to sequence model



Details on Recurrent Neural Networks

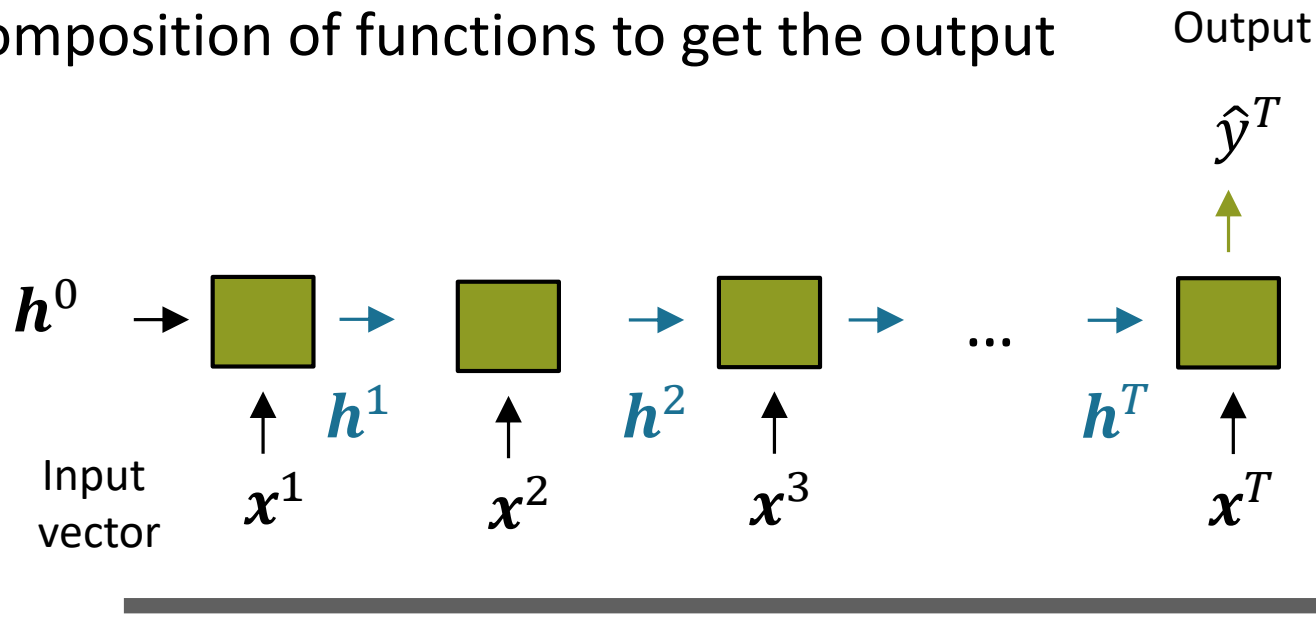
An RNN model

Let us consider *a many-to-one recurrent model*



Forward propagation through RNN

Composition of functions to get the output



Processing block is the same for all time moments.

The processing time is linear in T : it is $O(T)$.

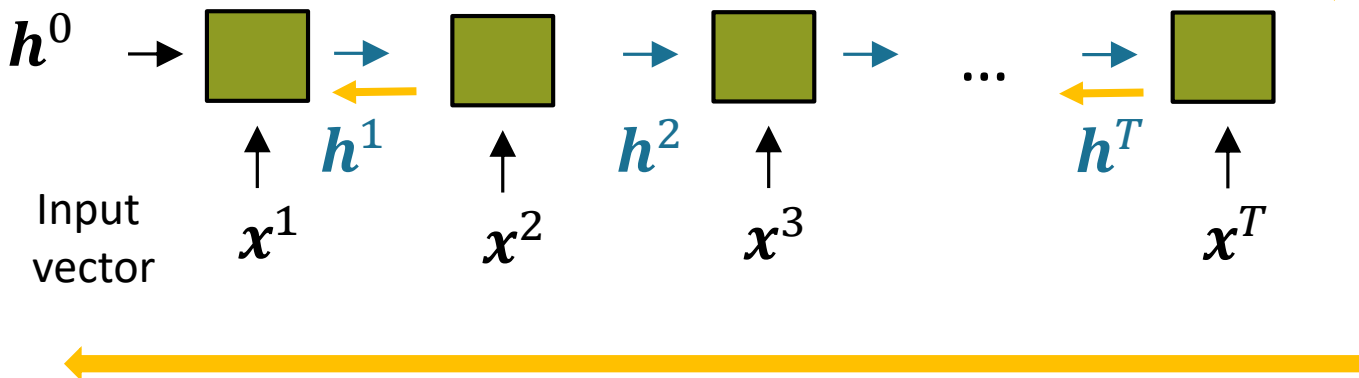
Backward propagation through RNN

Composition of functions to get the output and calculate loss and its derivatives

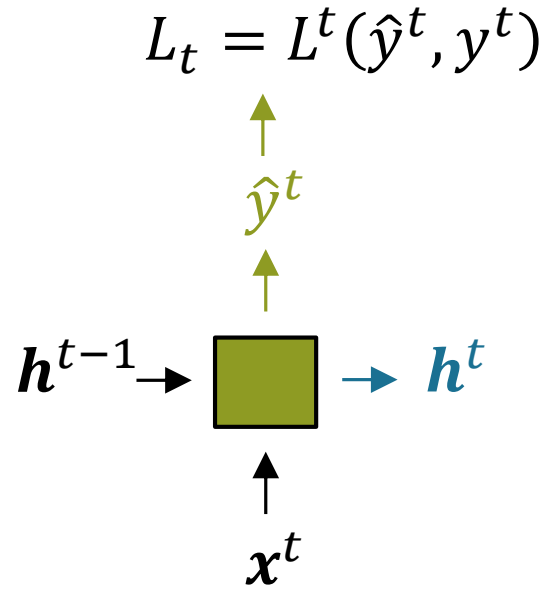
Let's discuss pros and cons for this model

Loss function

$$L(y^T, \hat{y}^T)$$



Backpropagation w.r.t. U

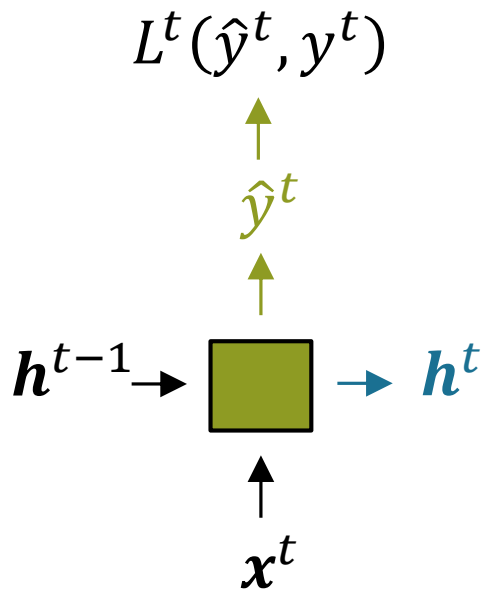


$$\frac{\partial L_T}{\partial U} = \frac{\partial L_T}{\partial \hat{y}^T} \frac{\partial \hat{y}^T}{\partial U}$$

Recall the formula for the output

$$\hat{y}^t = \text{softmax}(U\mathbf{h}^t + \mathbf{b}_y)$$

Backpropagation w.r.t. W



$$\frac{\partial L_T}{\partial W} = \frac{\partial L_T}{\partial \hat{\mathbf{y}}^T} \frac{\partial \hat{\mathbf{y}}^T}{\partial W}$$

$$\frac{\partial L_t}{\partial \hat{\mathbf{y}}^t} \frac{\partial \hat{\mathbf{y}}^t}{\partial W} = \frac{\partial L_t}{\partial \hat{\mathbf{y}}^t} \frac{\partial \hat{\mathbf{y}}^t}{\partial \mathbf{h}^t} \left(\frac{\partial \mathbf{h}^t}{\partial W} + \frac{\partial \mathbf{h}^t}{\partial \mathbf{h}^{t-1}} \frac{\partial \mathbf{h}^{t-1}}{\partial W} + \dots \right)$$

$$\frac{\partial L_t}{\partial \hat{\mathbf{y}}^t} \frac{\partial \hat{\mathbf{y}}^t}{\partial W} = \frac{\partial L_t}{\partial \hat{\mathbf{y}}^t} \frac{\partial \hat{\mathbf{y}}^t}{\partial \mathbf{h}^t} \sum_{i=0}^t \frac{\partial \mathbf{h}^i}{\partial W} \left(\prod_{j=i+1}^t \frac{\partial \mathbf{h}^j}{\partial \mathbf{h}^{j-1}} \right)$$

We have a product of derivatives
It can be a problem.

$$\mathbf{h}^t = \tanh(V\mathbf{x}^t + W\mathbf{h}^{t-1} + \mathbf{b}_h)$$

Problem: gradient explosion

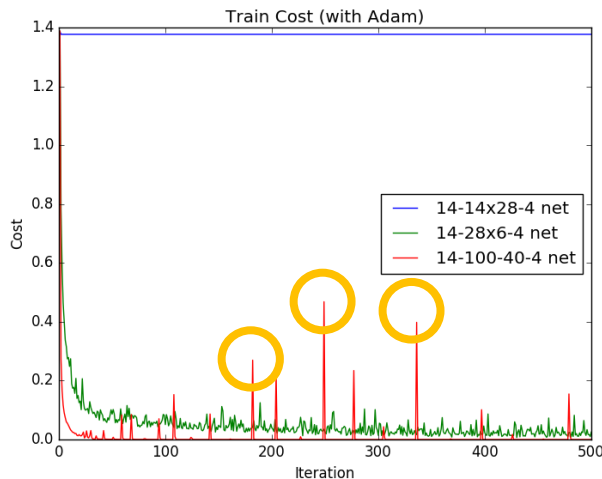
$$\frac{\partial L_t}{\partial W} = \frac{\partial L_t}{\partial \hat{y}^t} \frac{\partial \hat{y}^t}{\partial h^t} \sum_{i=0}^t \frac{\partial h^i}{\partial W} \left(\prod_{j=i+1}^t \frac{\partial h^j}{\partial h^{j-1}} \right)$$

$$\left\| \frac{\partial h^j}{\partial h^{j-1}} \right\| > 1, \text{ nonstationary (c.t. ARIMA lecture)}$$

Solution:

- Gradient clipping to scale big gradients

1. $g = \frac{\partial L}{\partial W}$
2. If $g > t$ for some threshold t :
$$g = \frac{t}{\|g\|} g$$

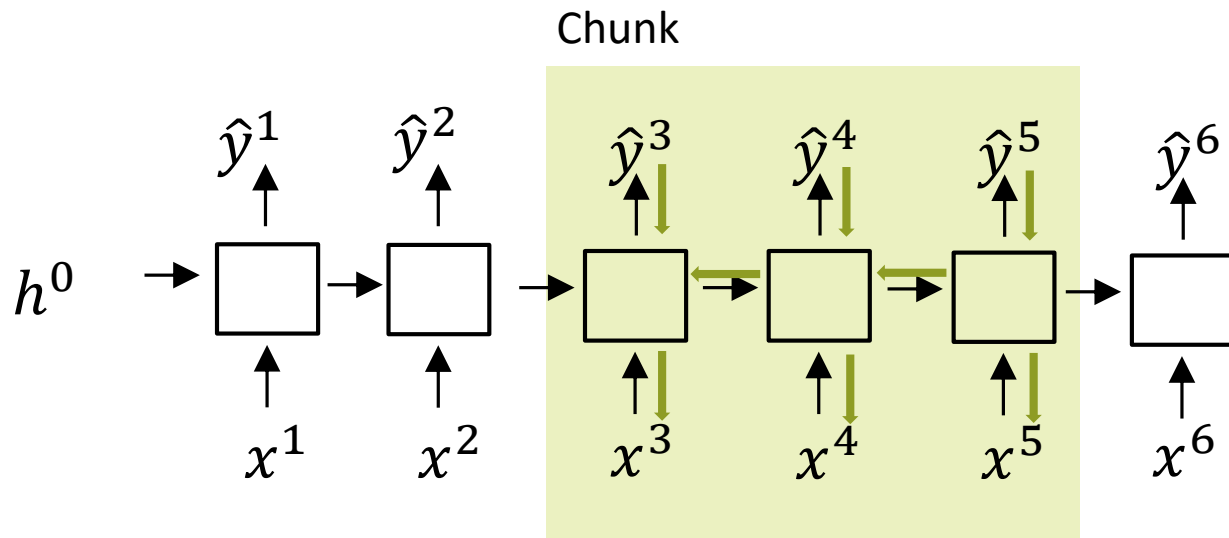


Threshold t is selected given the dynamic of loss function over iterations

Problem: gradient explosion

Solutions:

- Gradient clipping to scale big gradients
- Truncated backpropagation through time



Problems of classic RNN: gradient vanishing

A more serious problem

Many values < 1

Product $\ll 1$

$$\left\| \frac{\partial h^j}{\partial h^{j-1}} \right\| < 1$$

Can't learn long range dependences

Bias parameters to capture
long-term dependencies

Hard to detect!

Tricks:

Activation functions

- Use ReLU

Parameter initialization

- Initialize weights to identity matrix
- Initialize biases to zero

Another big problem of classic RNN

Problem: Neural networks forget fast, and it is hard to learn long-term dependencies

Solution: Gated architectures

More complex recurrent units with gates to control what information is passed through

- GRU: Gated Recurrent Unit
- LSTM Long-Short Term Memory
- Other?

Selection of RNN architecture

LSTM

GRU



Better RNN units: LSTM and GRU

- LSTM: long short term memory [1]
- GRU: Gated recurrent unit [2]

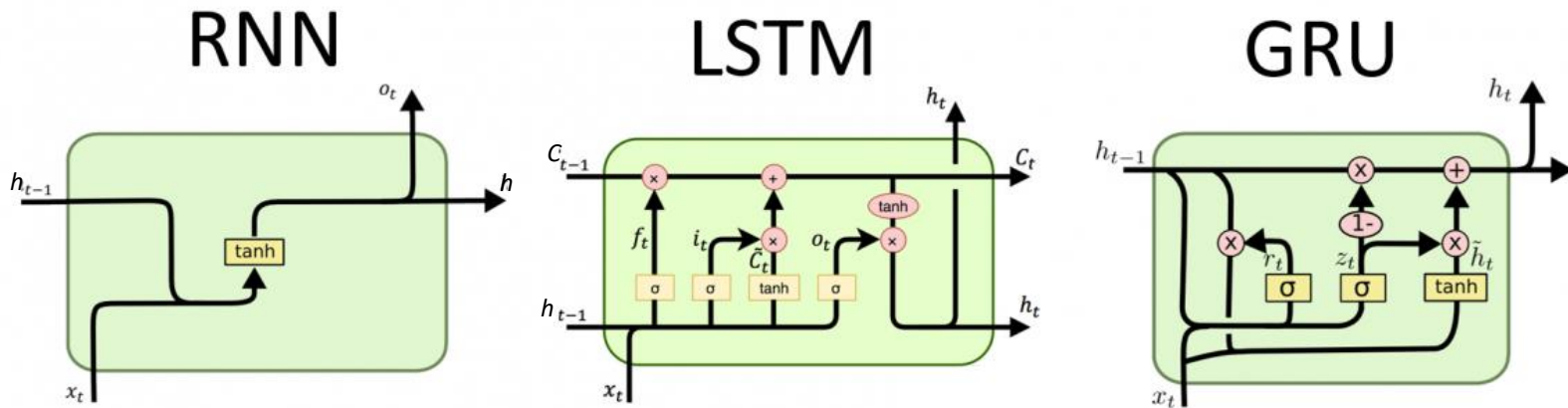


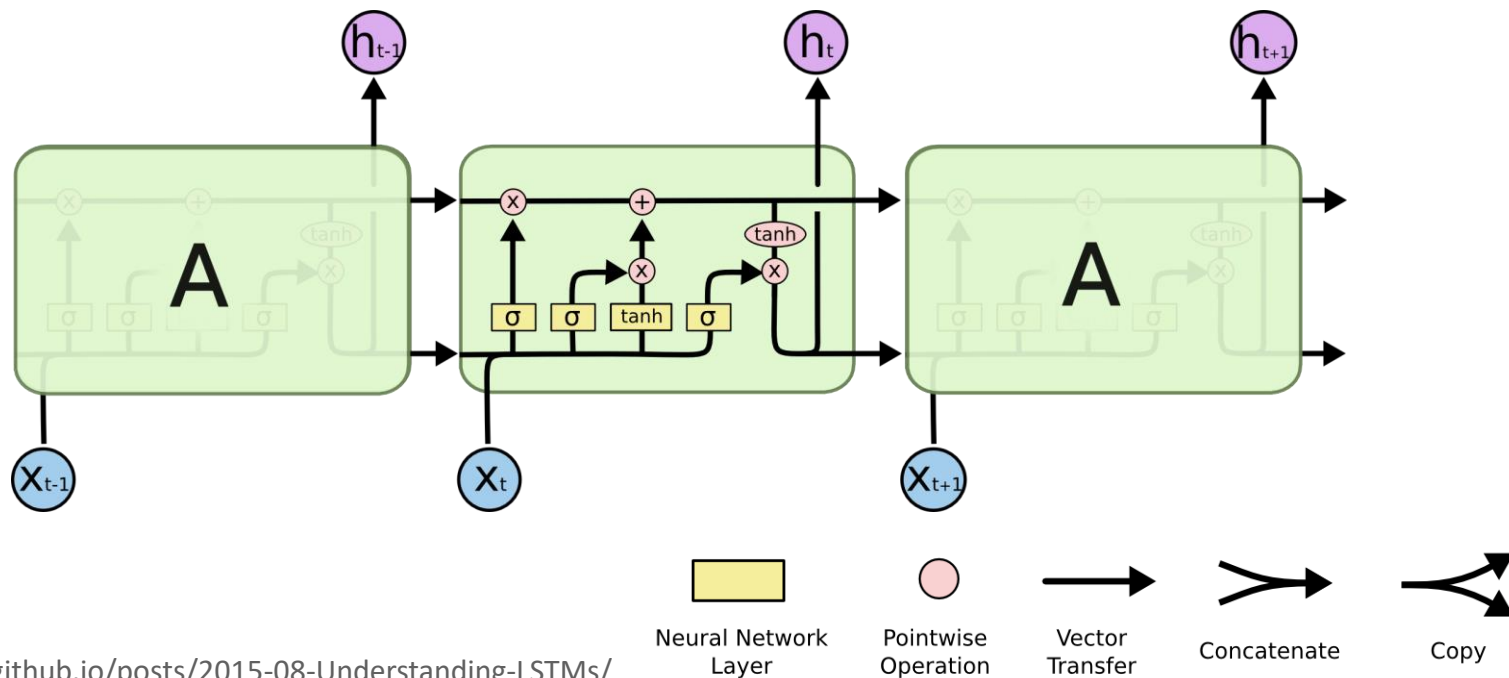
Image source: <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

1. Schmidhuber, J., & Hochreiter, S. Long short-term memory. Neural Computations. 1997.
2. Cho, K., Van Merriënboer et al. Learning phrase representations using RNN encoder-decoder for statistical machine translation. EMNLP. 2014.

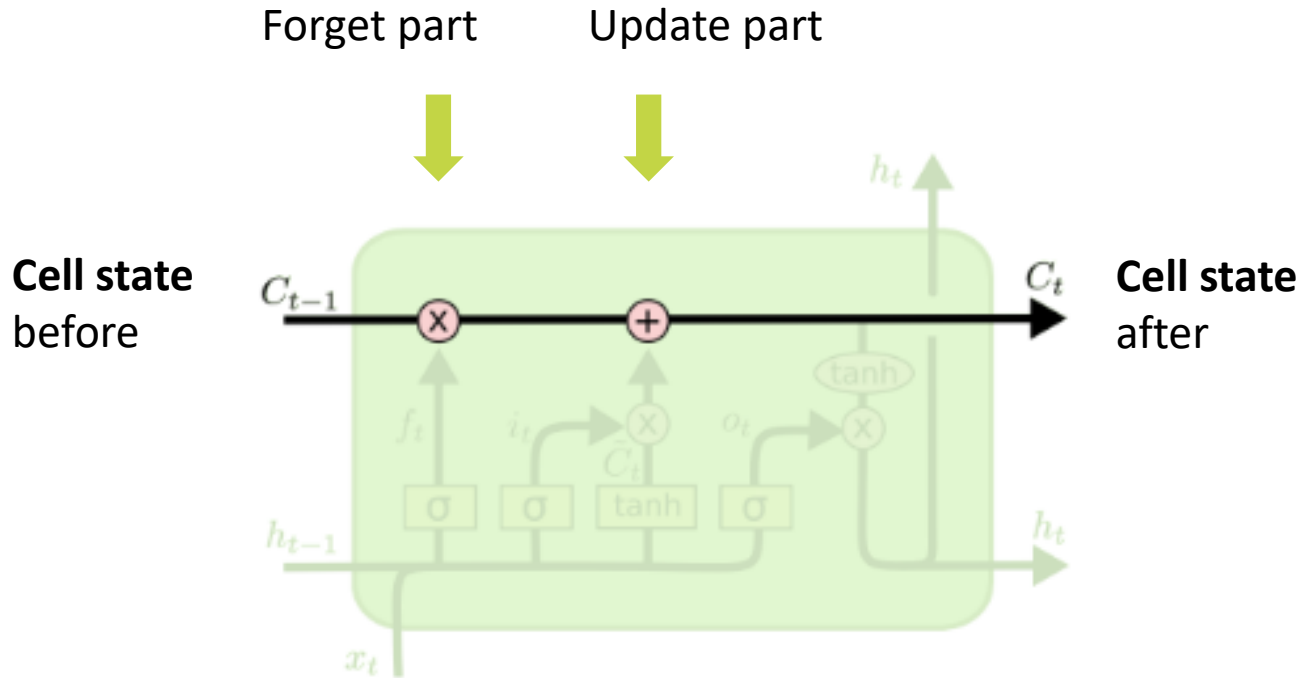
Details on how LSTM works

“Remembering information for long periods of time is practically the default behavior of LSTM”

Temporal representation =
Cell representation (long term) +
Hidden state (short term)

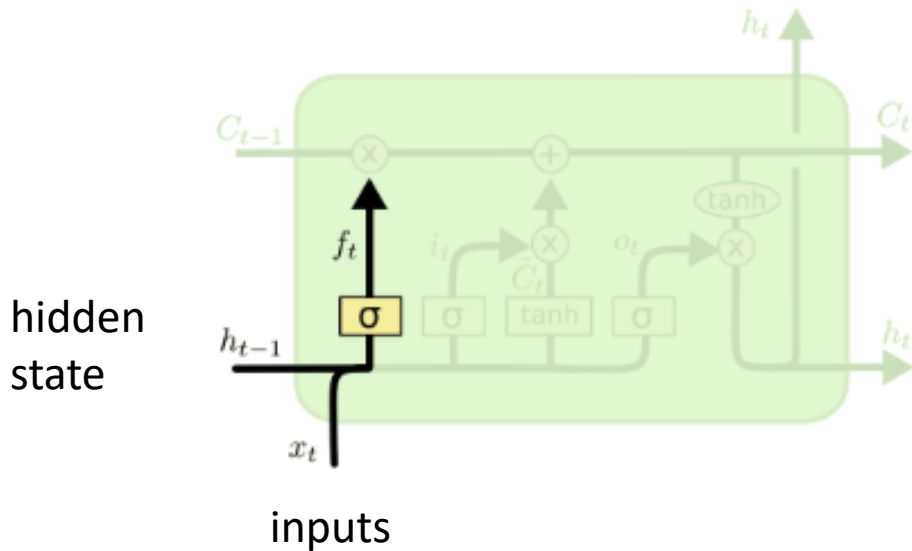


Long term memory part – Cell state



Forget part

Identify how much should we forget



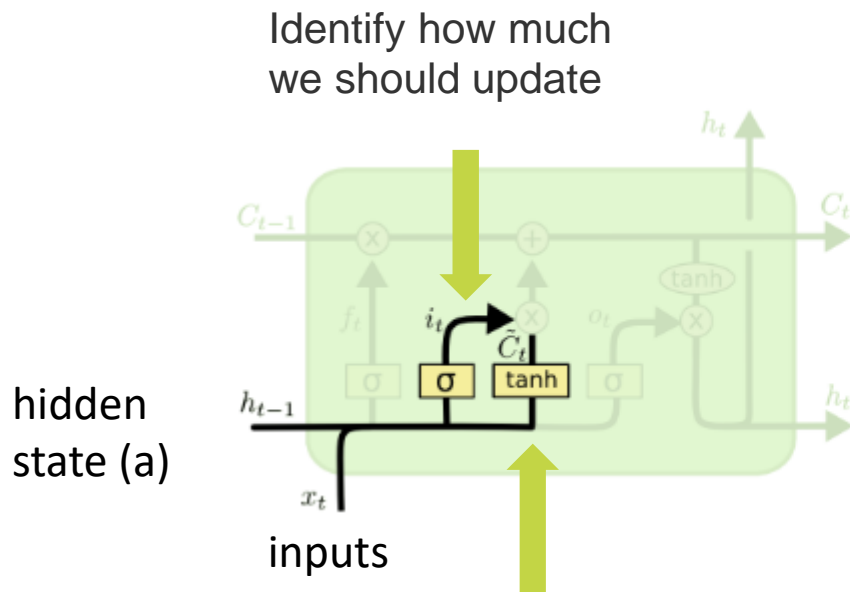
$$f_t = \sigma(W_f[\mathbf{h}_{t-1}; \mathbf{x}_t] + \mathbf{b}_f)$$

$[\mathbf{h}_{t-1}; \mathbf{x}_t]$ is the concatenation

W_f, \mathbf{b}_f are forget parameters

σ is the sigmoid function, so the output is between 0 and 1

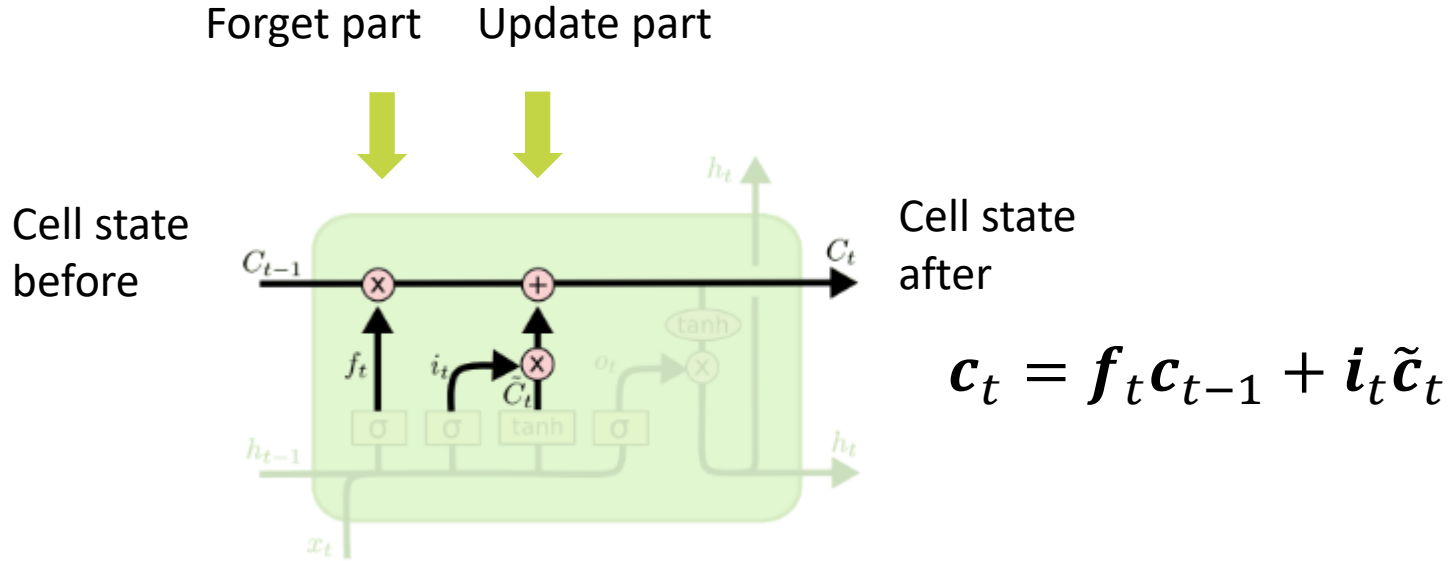
Update part



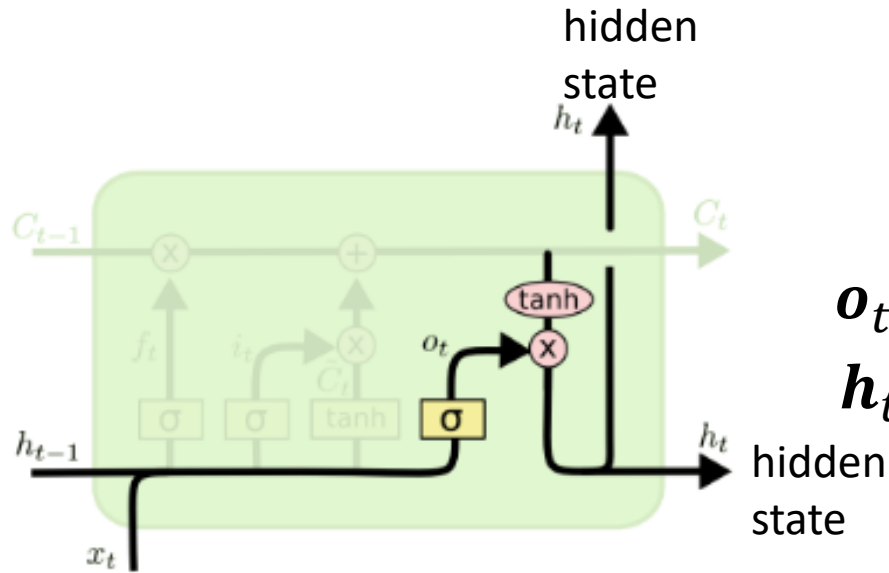
$$i_t = \sigma(W_i[h_{t-1}; x_t] + b_i)$$
$$\tilde{c}_t = \tanh(W_c[h_{t-1}; x_t] + b_c)$$

What update we do

Long term memory part – Cell state



Update everything else

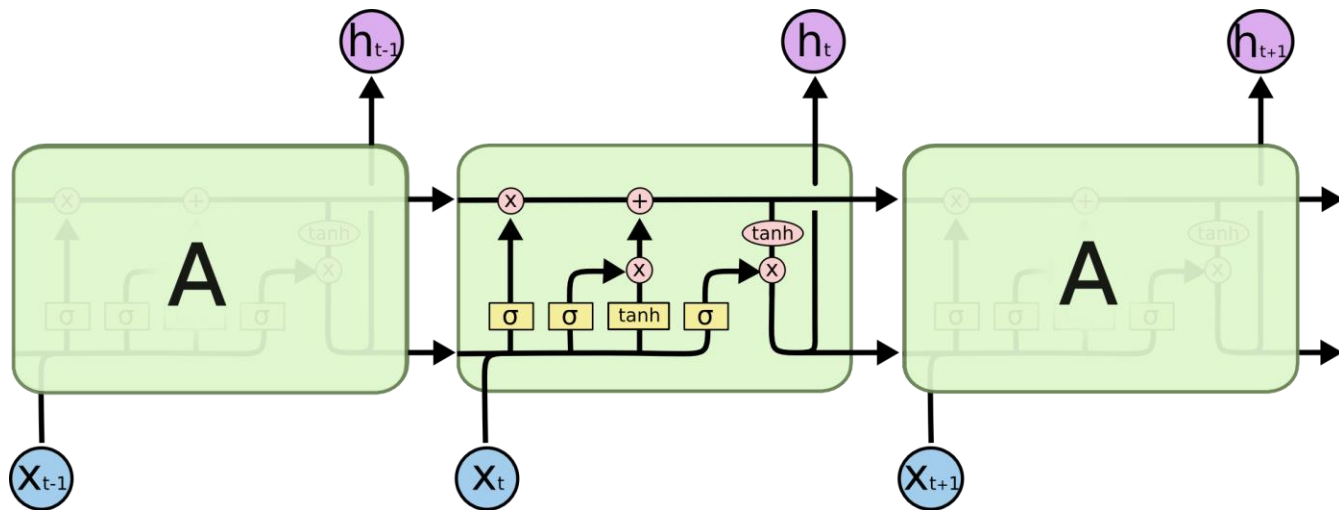


$$o_t = \sigma(W_o[h_{t-1}; x_t] + b_o)$$

$$h_t = o_t \tanh(c_t)$$

Details on how LSTM works

- There are cell and hidden (activation) states
- LSTM block forgets and updates cell state during processing at one block



GRU – Gated Recurrent Unit

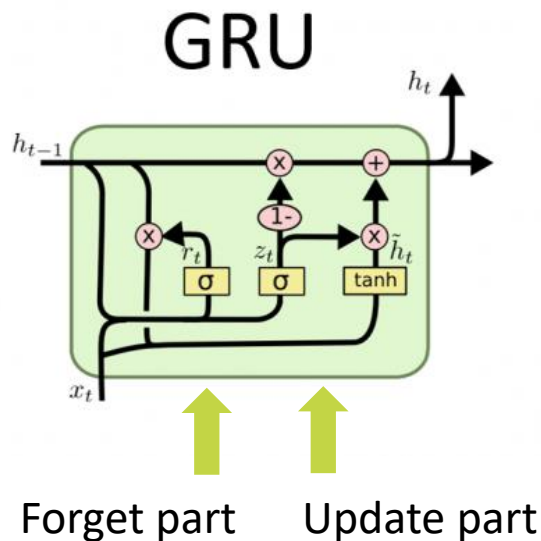
- Update gate – what to pay attention to
- Reset gate – what to forget

$$\mathbf{r}^t = \sigma(W_{xr}\mathbf{x}^t + W_{hr}\mathbf{h}^{t-1} + b_r)$$

$$\mathbf{z}^t = \sigma(W_{xz}\mathbf{x}^t + W_{hz}\mathbf{h}^{t-1} + b_z)$$

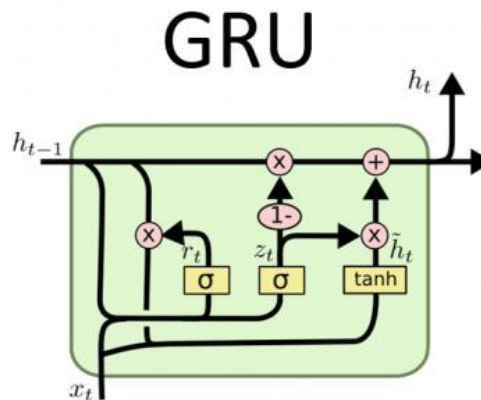
$$\tilde{\mathbf{h}}^t = \tanh(W_{xh}\mathbf{x}^t + W_{hr}(\mathbf{r}^t \odot \mathbf{h}^{t-1}) + b_h)$$

$$\mathbf{h}^t = \mathbf{z}^t \odot \mathbf{h}^{t-1} + (1 - \mathbf{z}^t) \odot \tilde{\mathbf{h}}^t$$

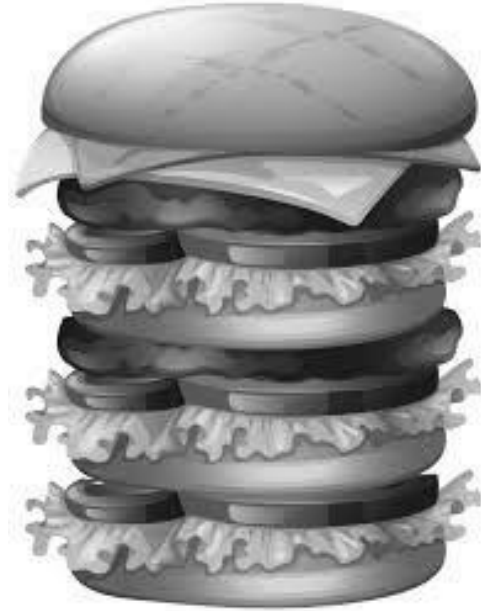


GRU – Gated Recurrent Unit

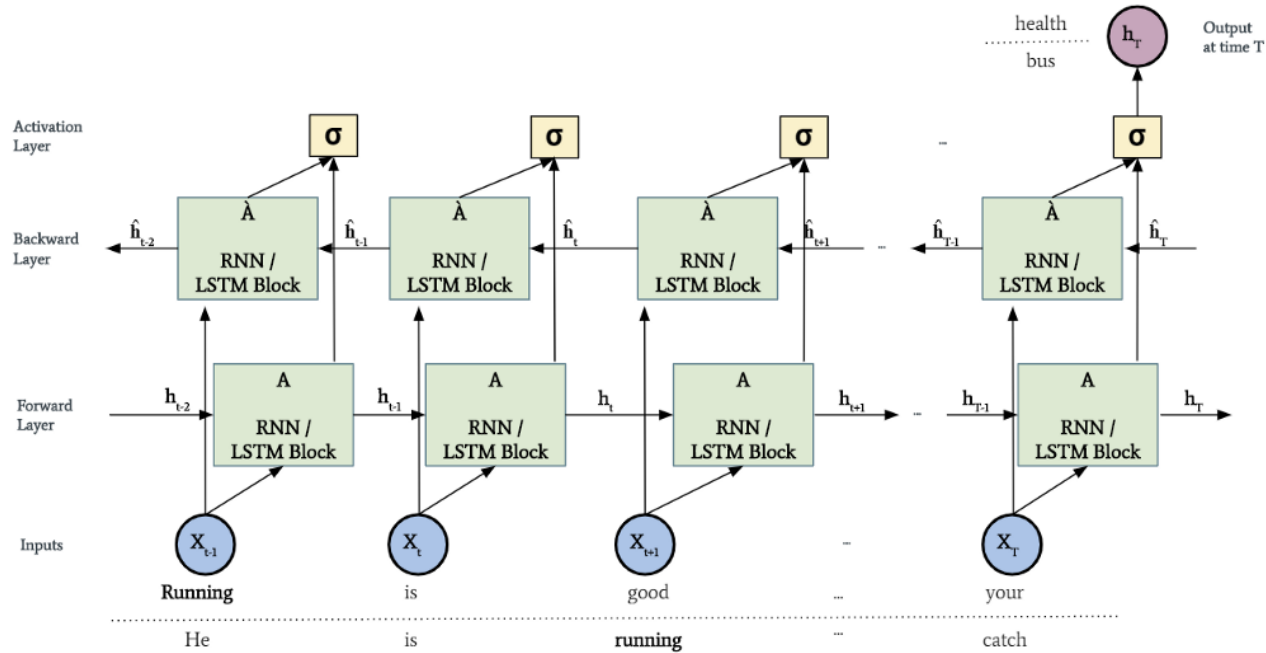
- Update gate – what to pay attention to
- Reset gate – what to forget
- Slightly worse than LSTM for NLP – but not in all problems
- Simpler and cheaper than LSTM



Multilayer architectures



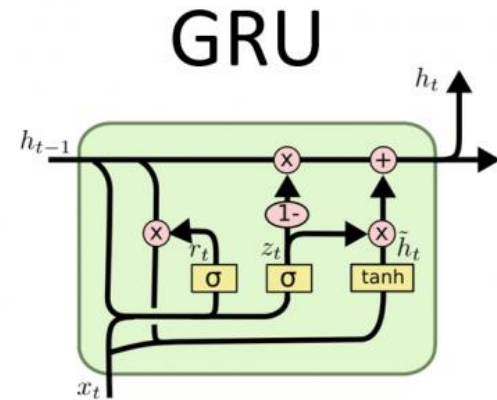
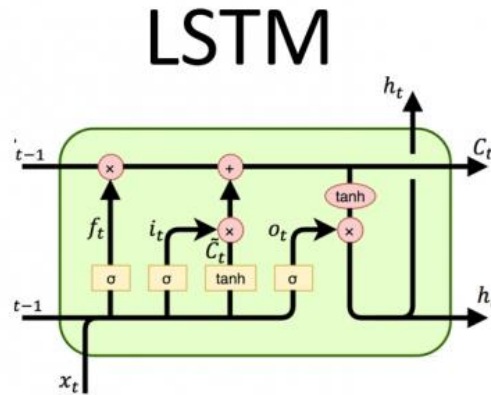
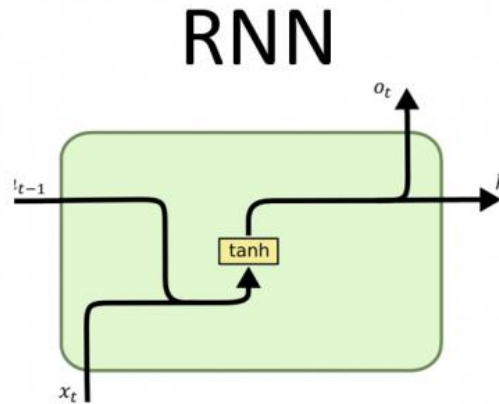
Other architectures: bidirectional LSTM



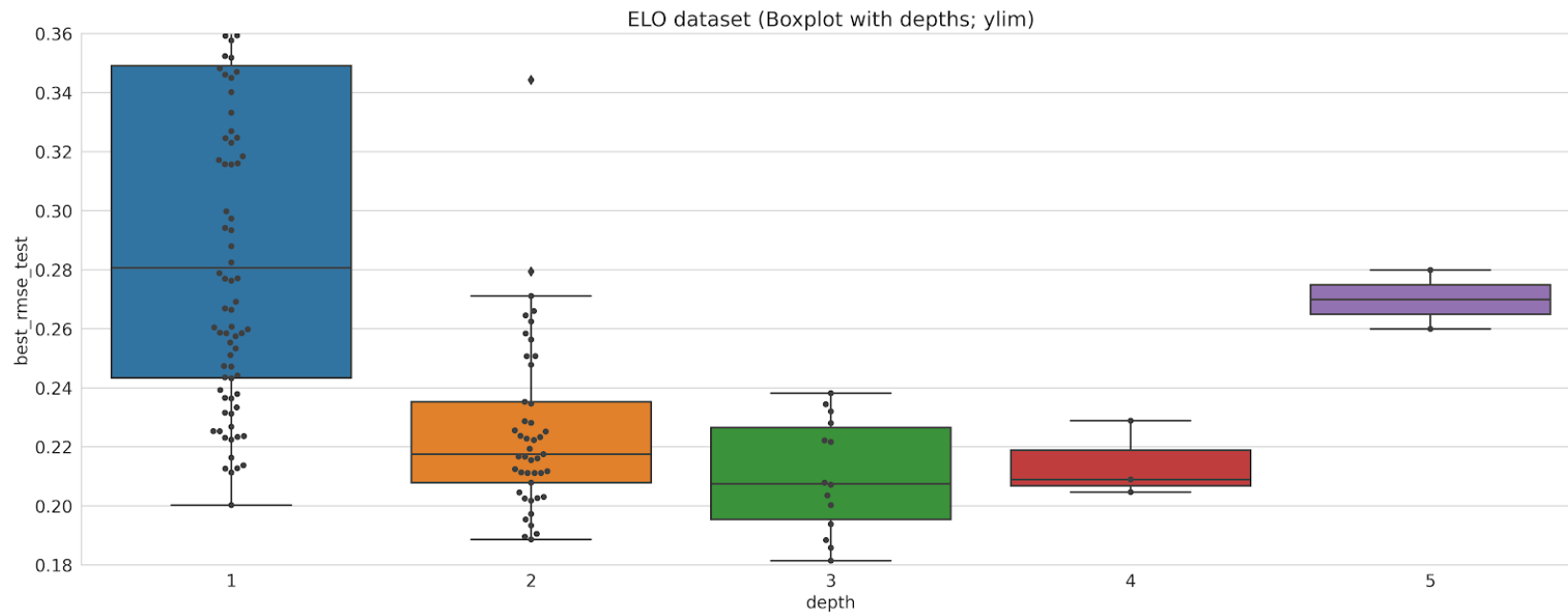
Other architectures: bidirectional LSTM

Bidirectional LSTM are useful when we benefit from the future data or can use it:

- Handwriting Recognition
- Speech Recognition
- Protein Structure Prediction (Bioinformatics)



Multiple layers RNN performance



**Does a better
recurrent block
exist?**



Towards a better recurrent block

- LSTM architecture is ad-hoc and has a substantial number of components whose purpose is not immediately apparent
- Like the LSTM, it is hard to tell, at a glance, which part of the GRU is essential for its functioning.
- Let's compare 10 000 different architectures on 3 problems with 1 000 of them pass the initial filtering stage: genetic algorithm
- Each architecture has been evaluated on about 220 hyperparameter settings.
- 230 000 hyperparameter configurations in total!



Best found architectures MUTx are close to GRU

GRU:

$$\begin{aligned}
 r_t &= \text{sigm}(W_{\text{xr}}x_t + W_{\text{hr}}h_{t-1} + b_r) \\
 z_t &= \text{sigm}(W_{\text{xz}}x_t + W_{\text{hz}}h_{t-1} + b_z) \\
 \tilde{h}_t &= \tanh(W_{\text{xh}}x_t + W_{\text{hh}}(r_t \odot h_{t-1}) + b_h) \\
 h_t &= z_t \odot h_{t-1} + (1 - z_t) \odot \tilde{h}_t
 \end{aligned}$$

Arch.	N	N-dropout	P
Tanh	3.612	3.267	6.809
LSTM	3.492	3.403	6.866
LSTM-f	3.732	3.420	6.813
LSTM-i	3.426	3.252	6.856
LSTM-o	3.406	3.253	6.870
LSTM-b	3.419	3.345	6.820
GRU	3.410	3.427	6.876
MUT1	3.254	3.376	6.792
MUT2	3.372	3.429	6.852
MUT3	3.337	3.505	6.840

Table 2. Negative Log Likelihood on the music datasets. N stands for Nottingham, N-dropout stands for Nottingham with nonzero dropout, and P stands for Piano-Midi.

MUT1:

$$\begin{aligned}
 z &= \text{sigm}(W_{\text{xz}}x_t + b_z) \\
 r &= \text{sigm}(W_{\text{xr}}x_t + W_{\text{hr}}h_t + b_r) \\
 h_{t+1} &= \tanh(W_{\text{hh}}(r \odot h_t) + \tanh(x_t) + b_h) \odot z \\
 &+ h_t \odot (1 - z)
 \end{aligned}$$

MUT2:

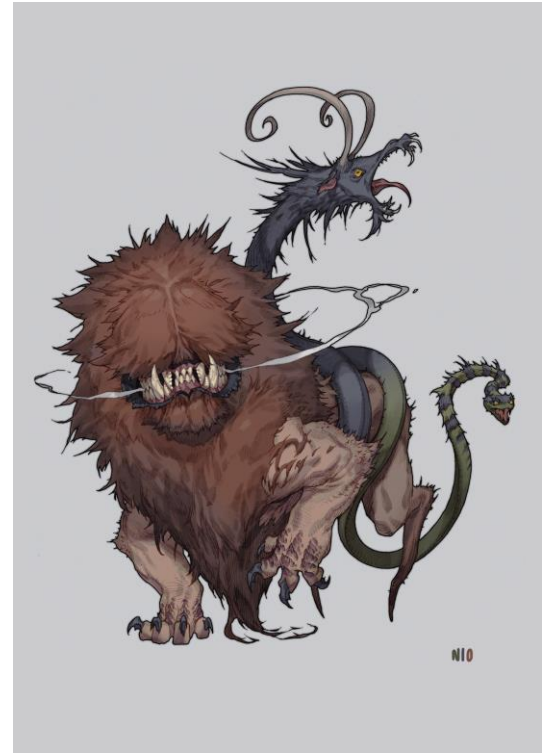
$$\begin{aligned}
 z &= \text{sigm}(W_{\text{xz}}x_t + W_{\text{hz}}h_t + b_z) \\
 r &= \text{sigm}(x_t + W_{\text{hr}}h_t + b_r) \\
 h_{t+1} &= \tanh(W_{\text{hh}}(r \odot h_t) + W_{\text{xh}}x_t + b_h) \odot z \\
 &+ h_t \odot (1 - z)
 \end{aligned}$$

MUT3:

$$\begin{aligned}
 z &= \text{sigm}(W_{\text{xz}}x_t + W_{\text{hz}}\tanh(h_t) + b_z) \\
 r &= \text{sigm}(W_{\text{xr}}x_t + W_{\text{hr}}h_t + b_r) \\
 h_{t+1} &= \tanh(W_{\text{hh}}(r \odot h_t) + W_{\text{xh}}x_t + b_h) \odot z \\
 &+ h_t \odot (1 - z)
 \end{aligned}$$

Towards a better recurrent block

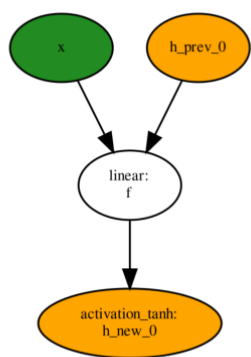
- LSTM architecture is ad-hoc and has a substantial number of components whose purpose is not immediately apparent
- Like the LSTM, it is hard to tell, at a glance, which part of the GRU is essential for its functioning.
- Let's compare 8 LSTM variants and hope for the best by search over the space of hyperparameters with 5400 runs in total
- No significant improvement over the common LSTM
- Some advices on hyperparameters selection



Neural architecture search for a better Recurrent block

- Linear: $f(x_1, \dots, x_n) = W_1x_1 + \dots + W_nx_n + b$,
- Blending (element wise): $f(z, x, y) = z \odot x + (1 - z) \odot y$,
- Element wise product and sum,
- Activations: Tanh, Sigmoid, and LeakyReLU.

- Number of nodes < 25
- Number of hidden states < 4
- Number of linear input vectors < 4



(a) Simple RNN Cell



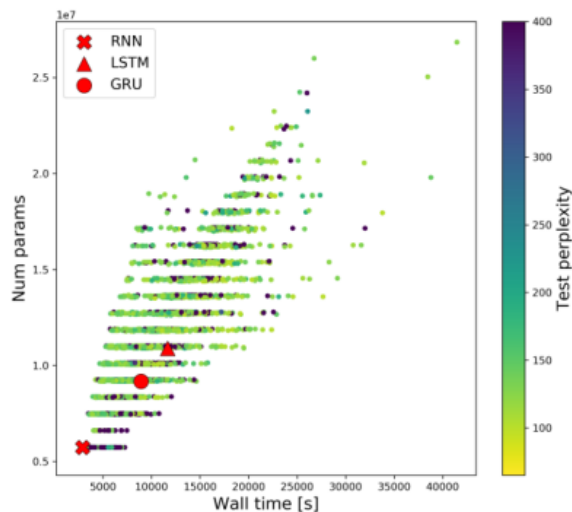
(b) LSTM Cell



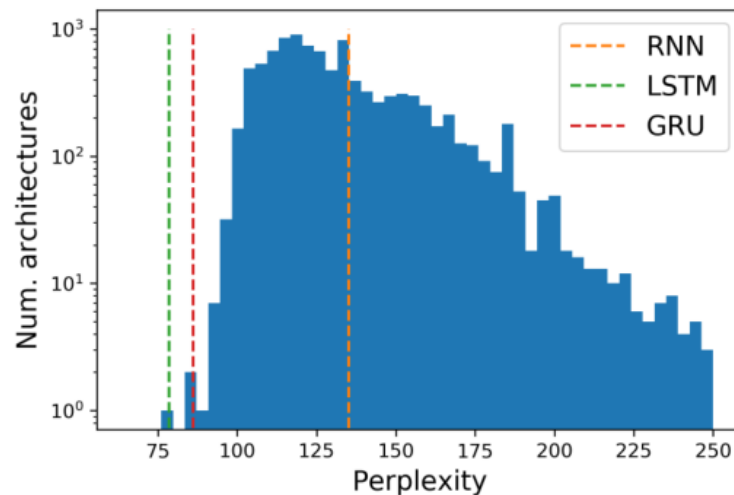
(c) GRU Cell

Figure 2: Examples of conventional RNN cells. Colors of nodes highlight the corresponding previous and new hidden states, green color also highlights the input vector. Black dashed, blue and red edges indicate blending arguments z , x and y respectively.

Neural architecture search for a better Recurrent block



(a) Joint distribution of metrics.

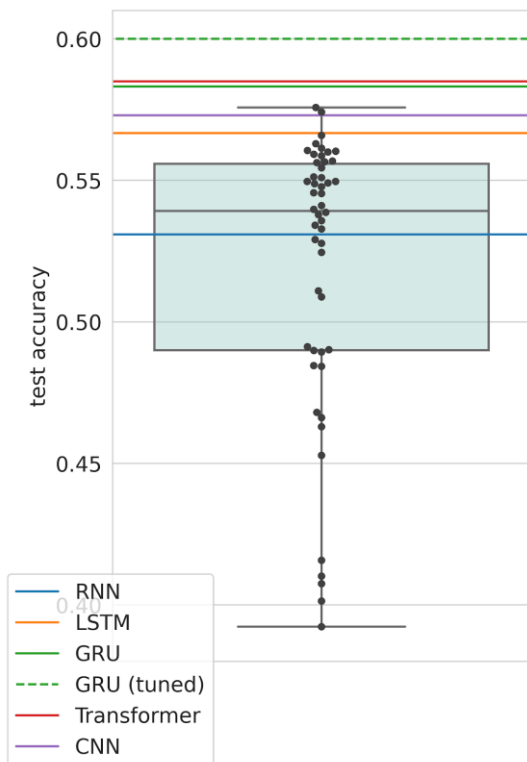


(b) Best test perplexity distribution.

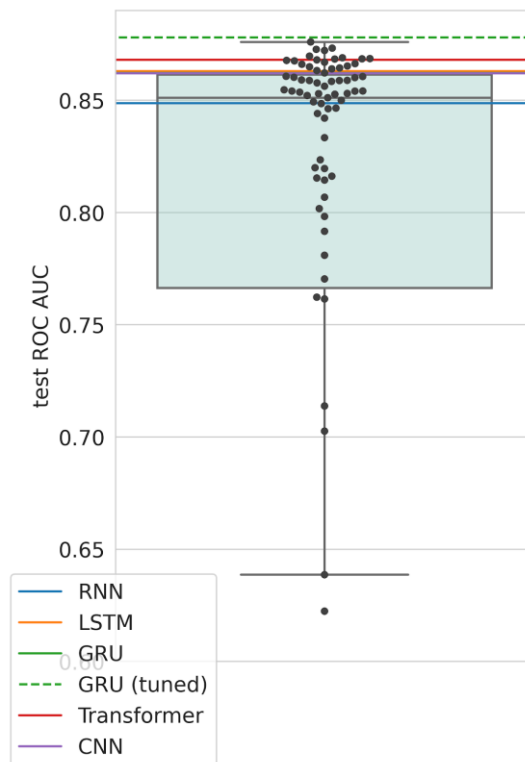
Figure 4: Architectures metrics on PTB.

Best architecture for Transactions data

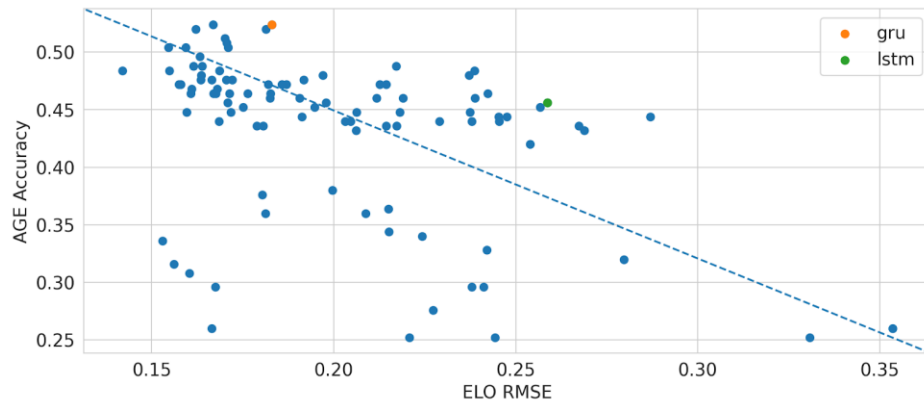
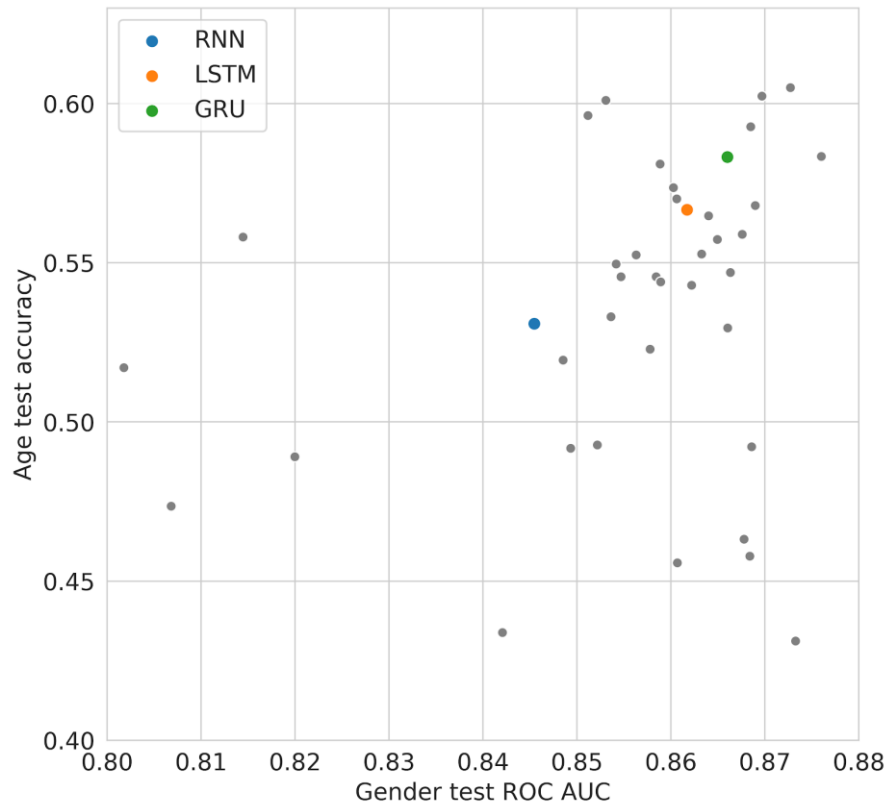
Age prediction




Gender prediction



Architecture performance transferability

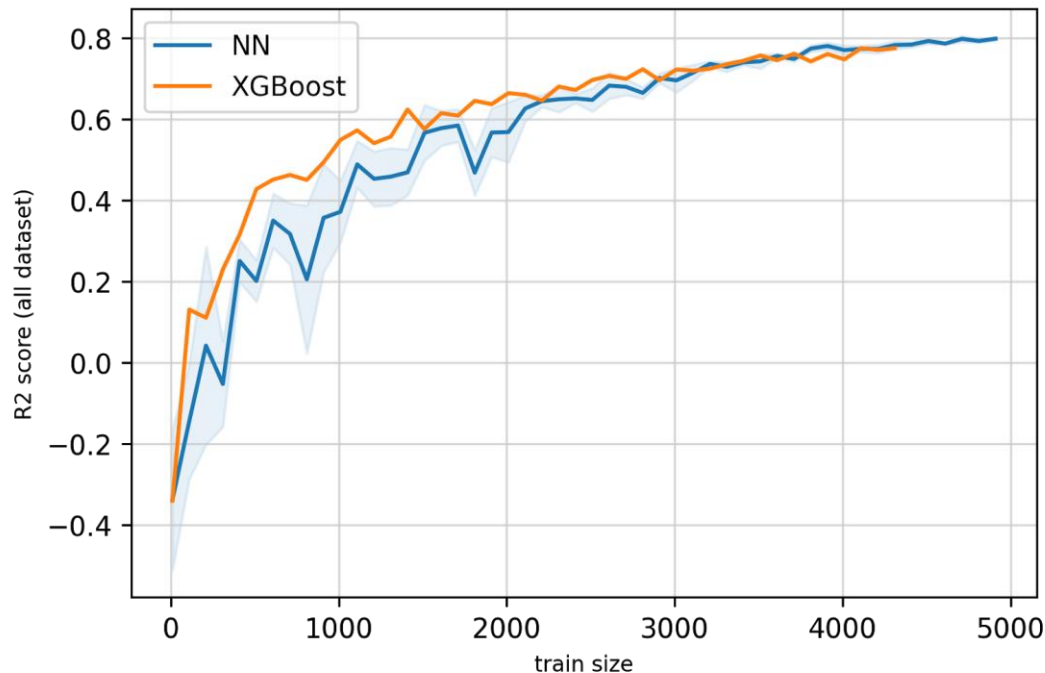




**Can we predict the
performance of an
architecture?**

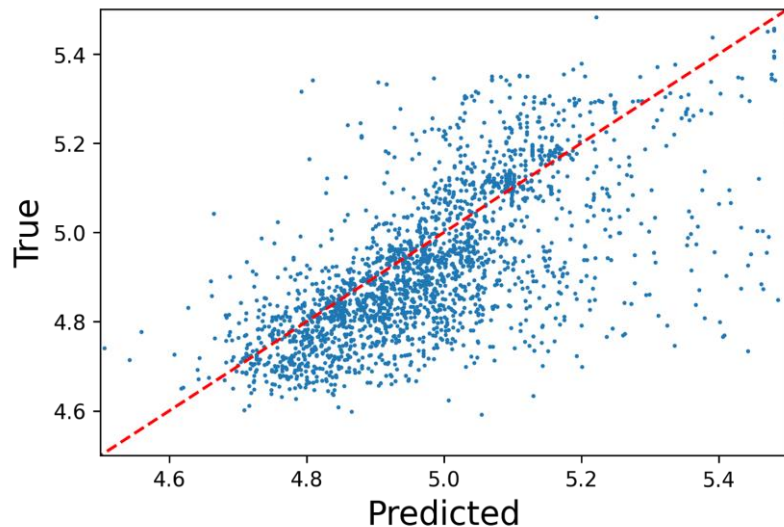
We can, but we need large samples

For path features and NLP NAS bench with 11000 graph2vec features



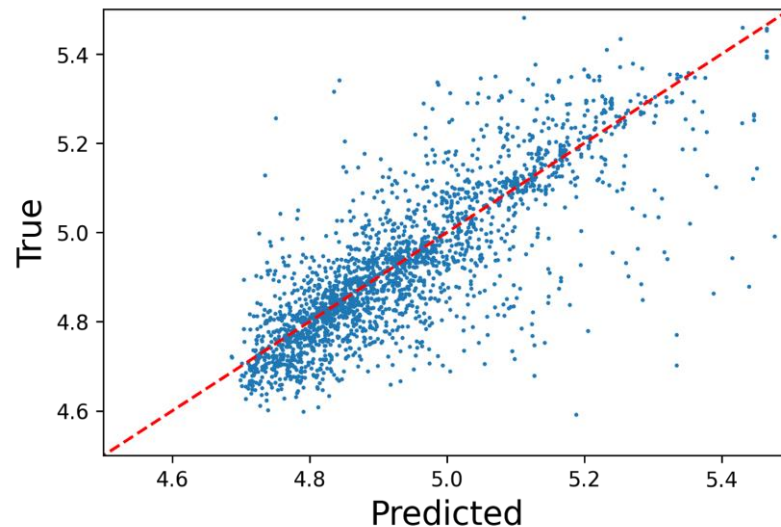
We can, but we need large samples

For path features and NLP NAS bench with 11000 features



$R^2 = 0.7455$

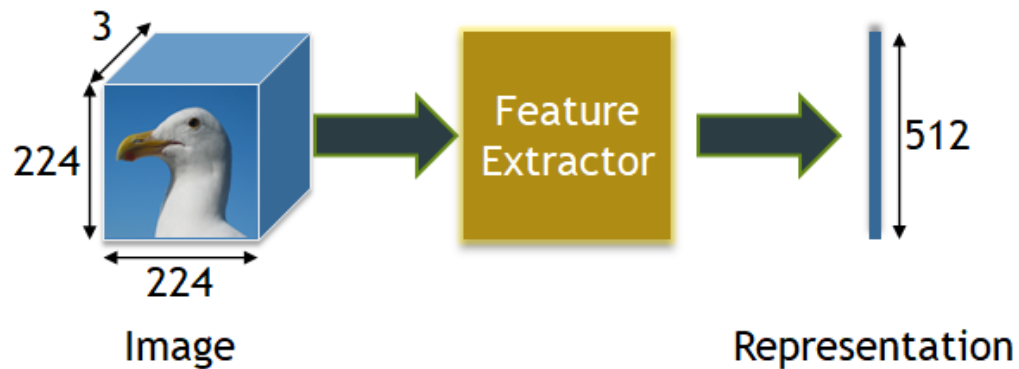
XGBRegressor 1000 estimators,
lr=0.05, default parameters



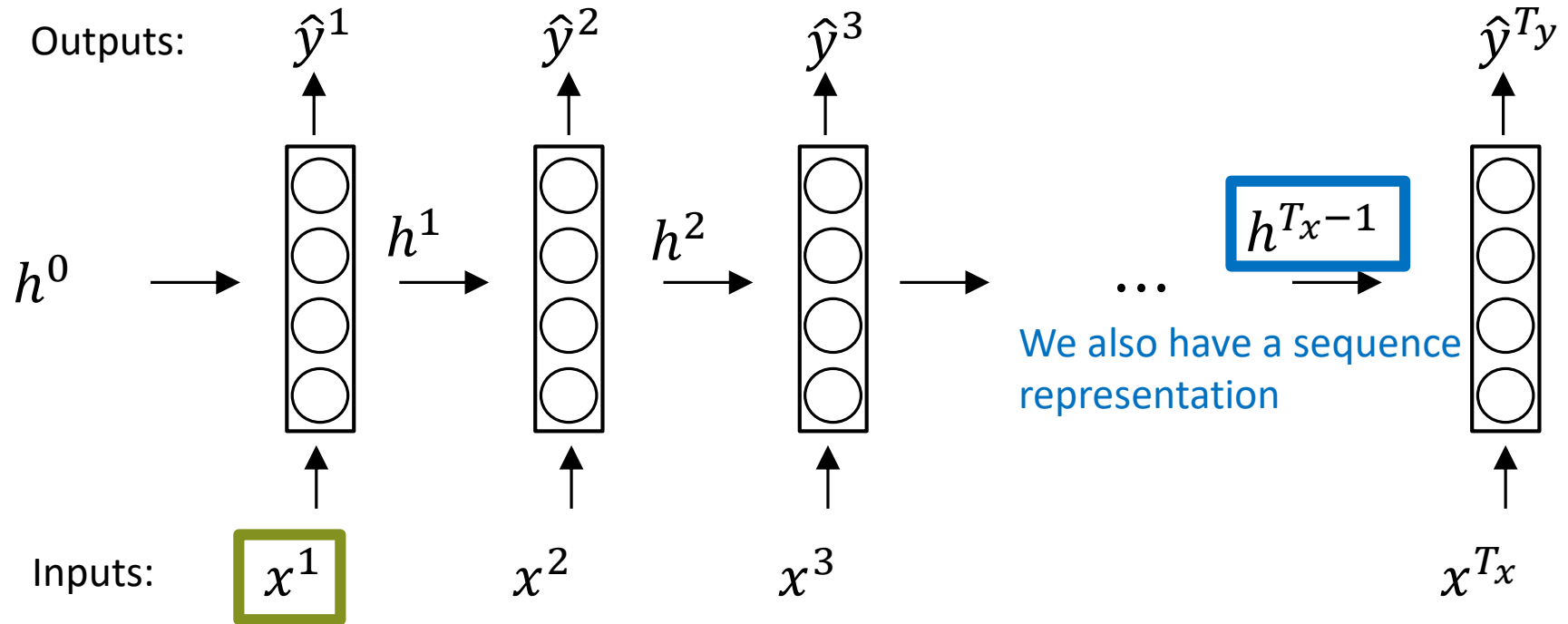
$R^2 = 0.7866$

Dense NN, 10 layers with dim=20,
ELU lr=0.005, 120 epochs

Representation learning is the core feature of Neural Networks



Representation learning is still here for Recurrent Neural Networks



Most of the time we also learn representations of objects in an end2end manner with backpropagation

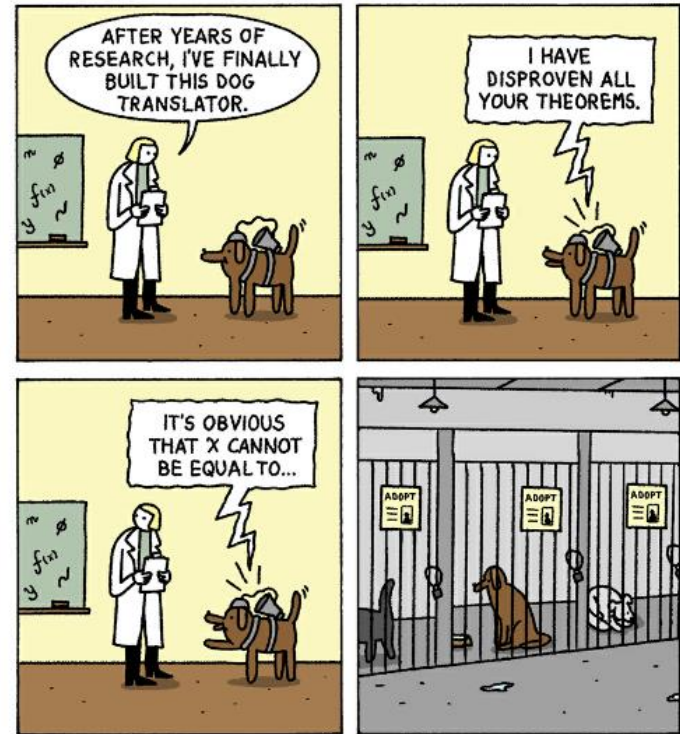
Take-home messages

- For some types of data classic methods fail:
we need to learn a representation i.e. extract features automatically
- Neural Networks provide enough flexibility for this problem for various data types
- The basic architecture is Recurrent Neural Network (RNN)
- But we can do better in terms of keeping the necessary information with LSTM and GRU blocks/architectures

Sources

- Recurrent Neural Networks | in Deep learning course by MIT 6.S191
- Coursera course on Sequence models
<https://www.coursera.org/learn/nlp-sequence-models>
- <http://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Machine translation: application example



War and Peas

Dog Translation Machine

Machine translation, 50-s

Cold war child: translator from Russian to English IBM 701 Translator

Doctor Dostert predicted that “five, perhaps three years hence, interlingual meaning conversion by electronic process in important functional areas of several languages may well be an accomplished fact.” (1954)

Rule-based approach that uses English-Russian dictionary

Project shut-down in 5 years: no significant progress

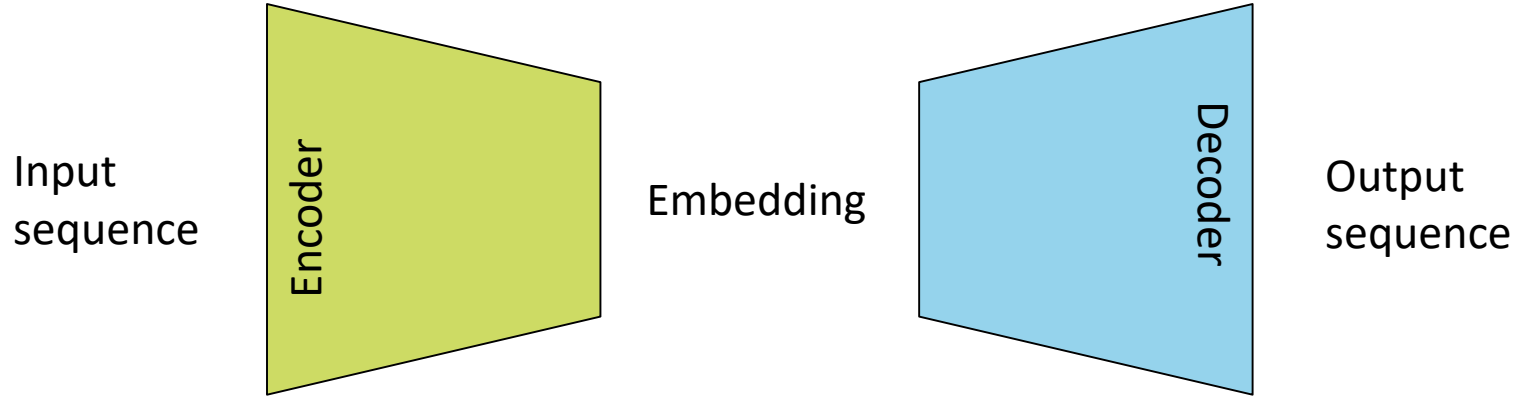


https://www.ibm.com/ibm/history/exhibits/701/701_translator.html
<https://youtu.be/8ZtdVUB007A>

Statistical approach – the leading one before 2014

- Complicated and heavy model
- Many separate components
- Complex generation of inputs
- Support of a sophisticated system
- Quality is not great

Neural network for machine translation



seq2seq (sequence to sequence) architecture

Not a perfect solution with little control over result



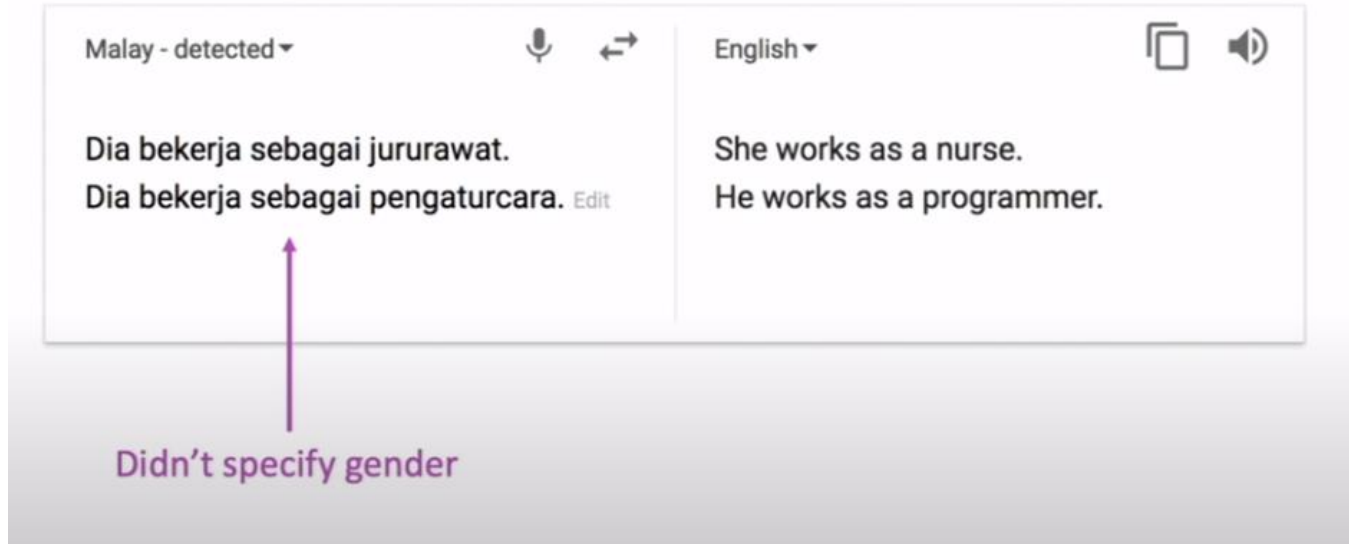
[Open in Google Translate](#)

[Feedback](#)

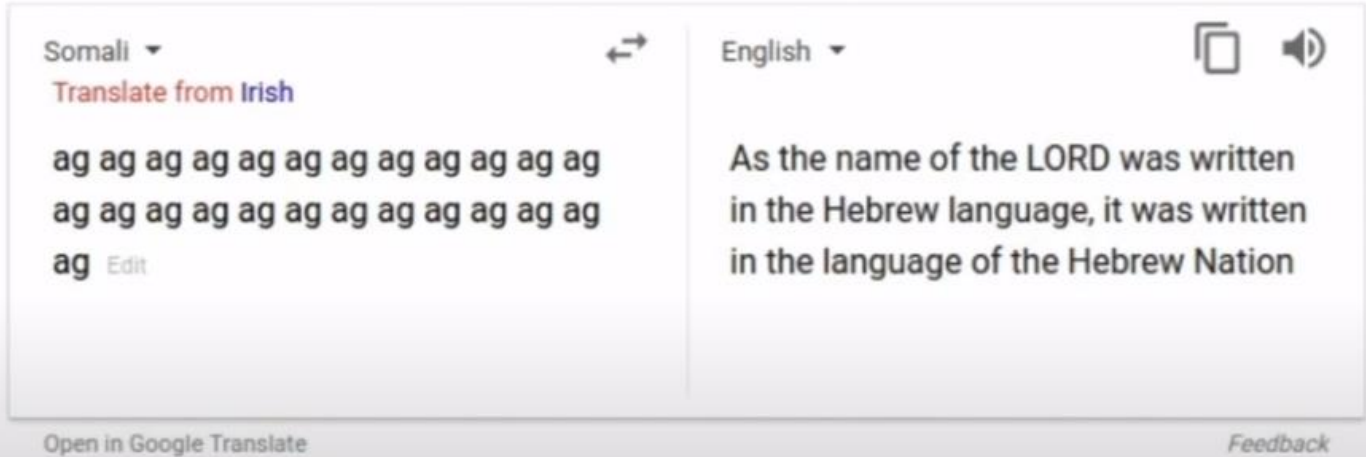


?

Contains biases from learning sample

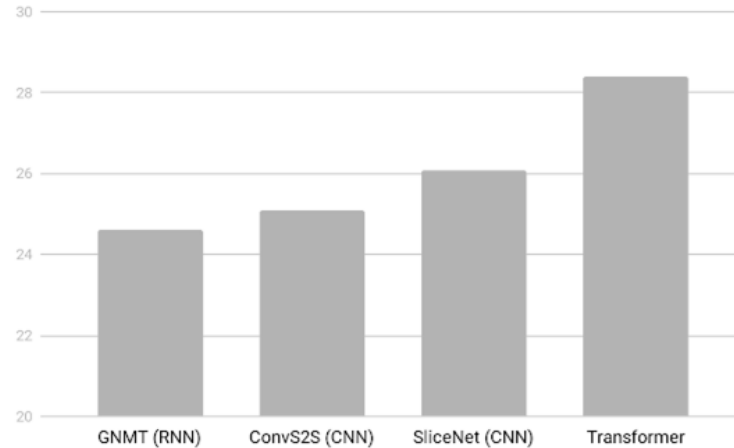


Output strange results

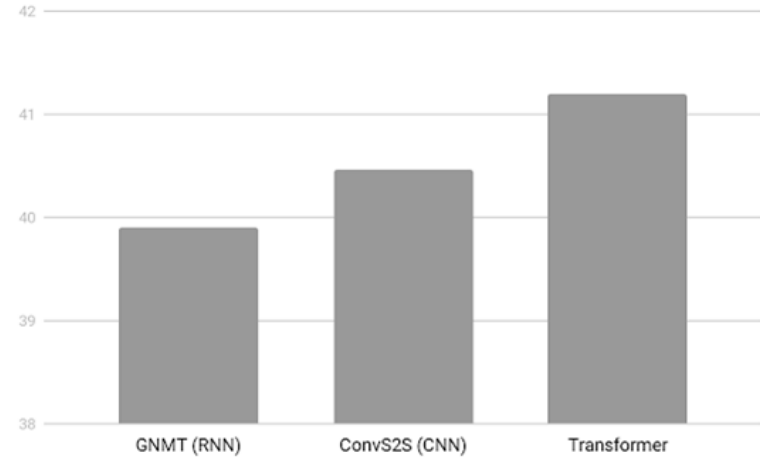


Attention/transformer – new state-of-the art for Machine translation

English German Translation quality



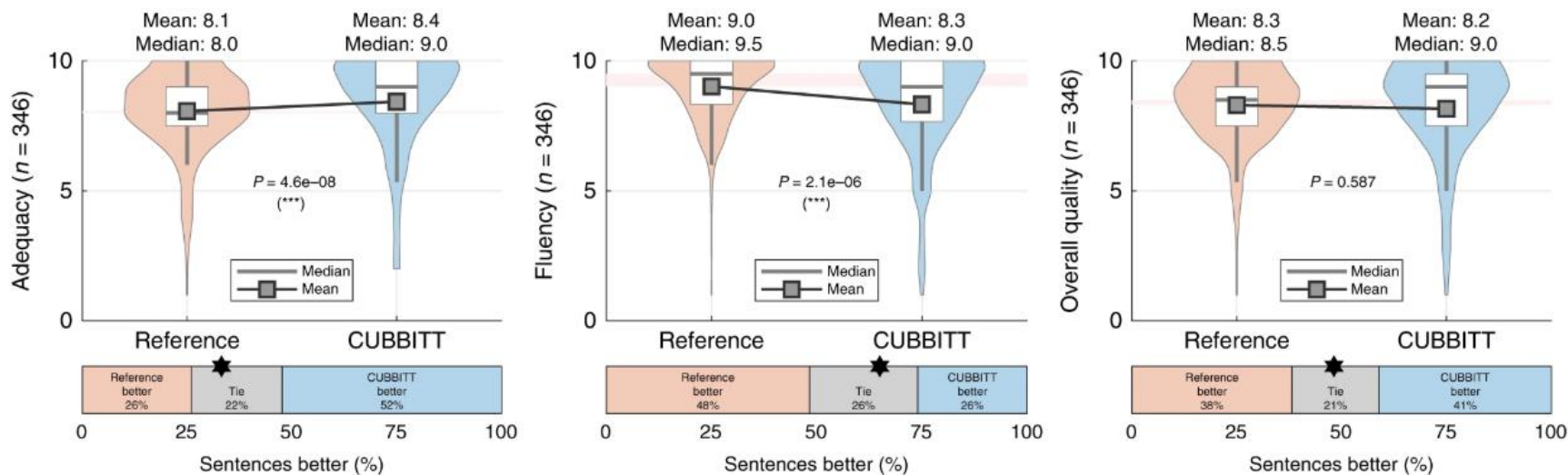
English French Translation Quality



- Maximize BLEU translation quality metric
- The best model is **Transformer**

NLP deep learning models: CUBBITT is comparable to reference human translation for a selected dataset on News translation

Context-aware evaluation: 6 non-professionals



Conclusions

- Machine translation problem was attacked by many problems during a long period of time
- Right now the state-of-the-art is Neural Networks Transformers
- They can even improve over human-translation

Another semi-structured data problems

Credit scoring:
default prediction

Money
Transactions data



Fraud detection
in healthcare insurance

A history of visits



Is there a fraud?