

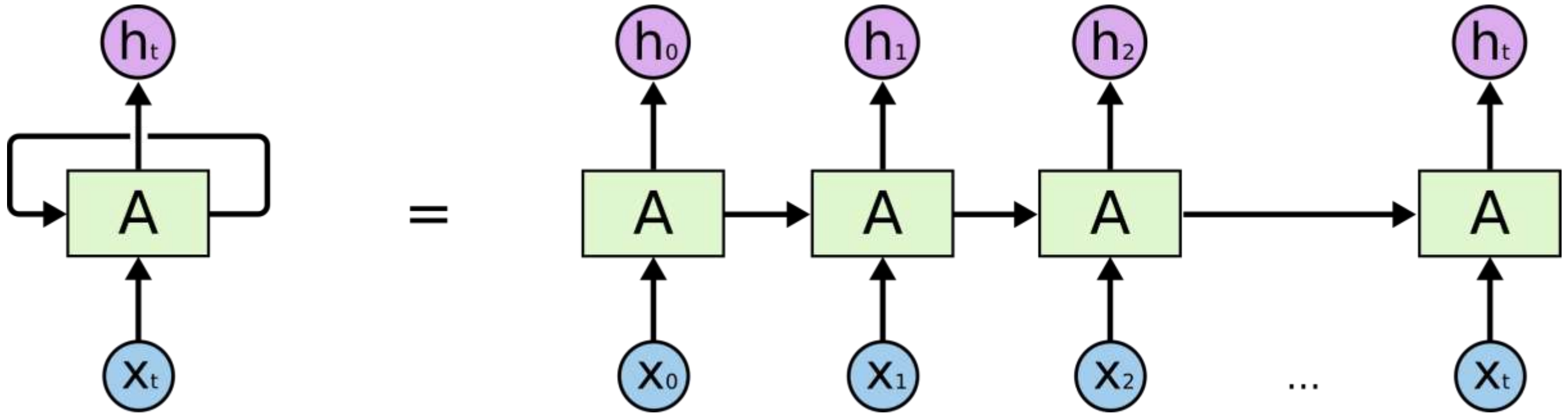
# xLSTM

---

Extended  
Long Short-Term  
Memory

# RNN

---

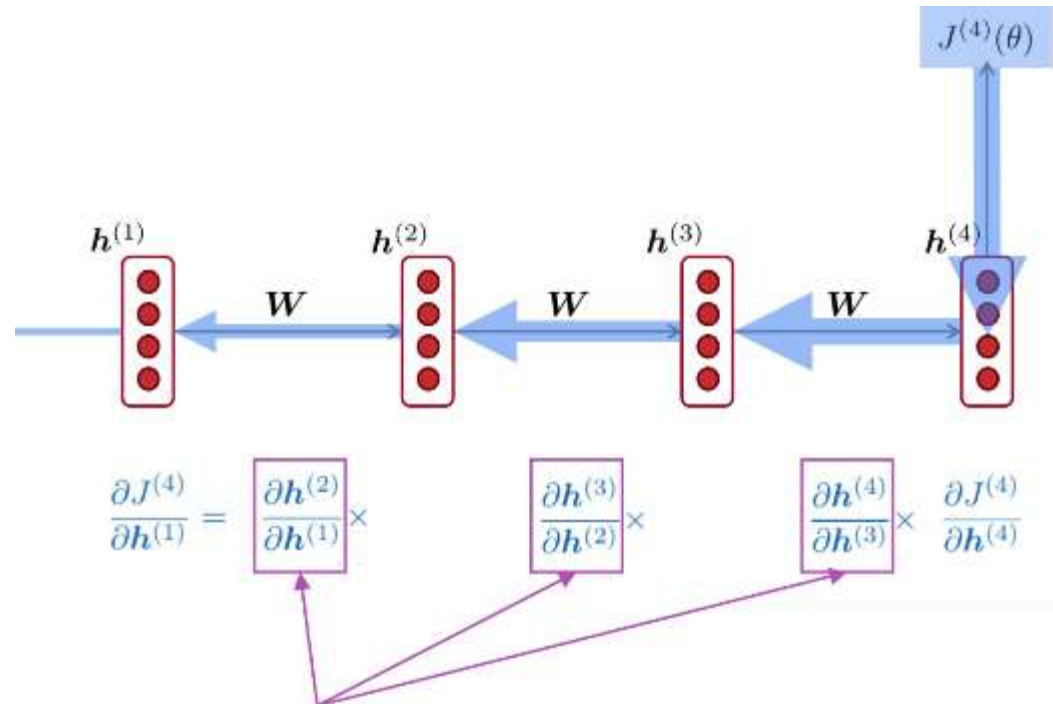
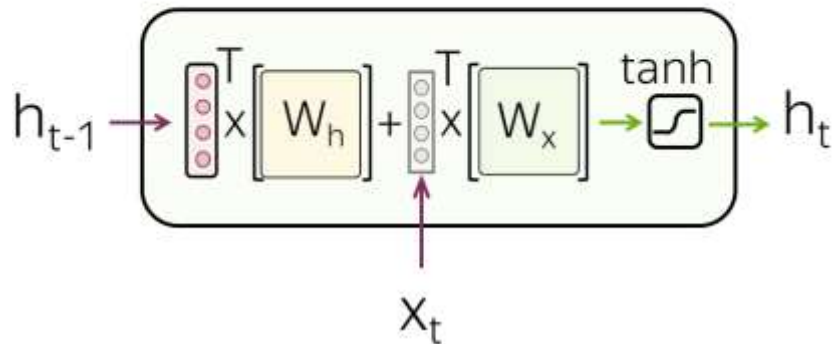


# Проблемы

- Вектор контекста меняется на каждом шаге, поэтому не имеет возможности пройти на некотором этапе неизменным
- Затухание градиента

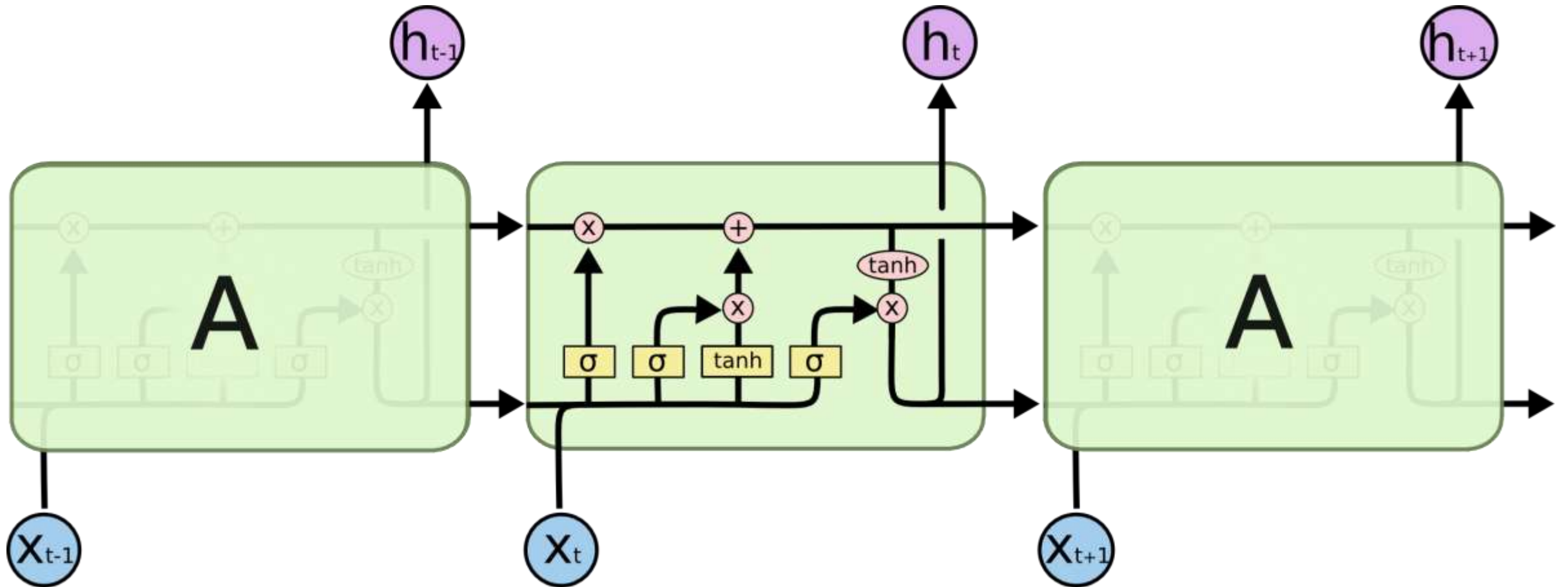
## Vanilla RNN

$$h_t = \tanh(h_{t-1}W_h + x_tW_x)$$

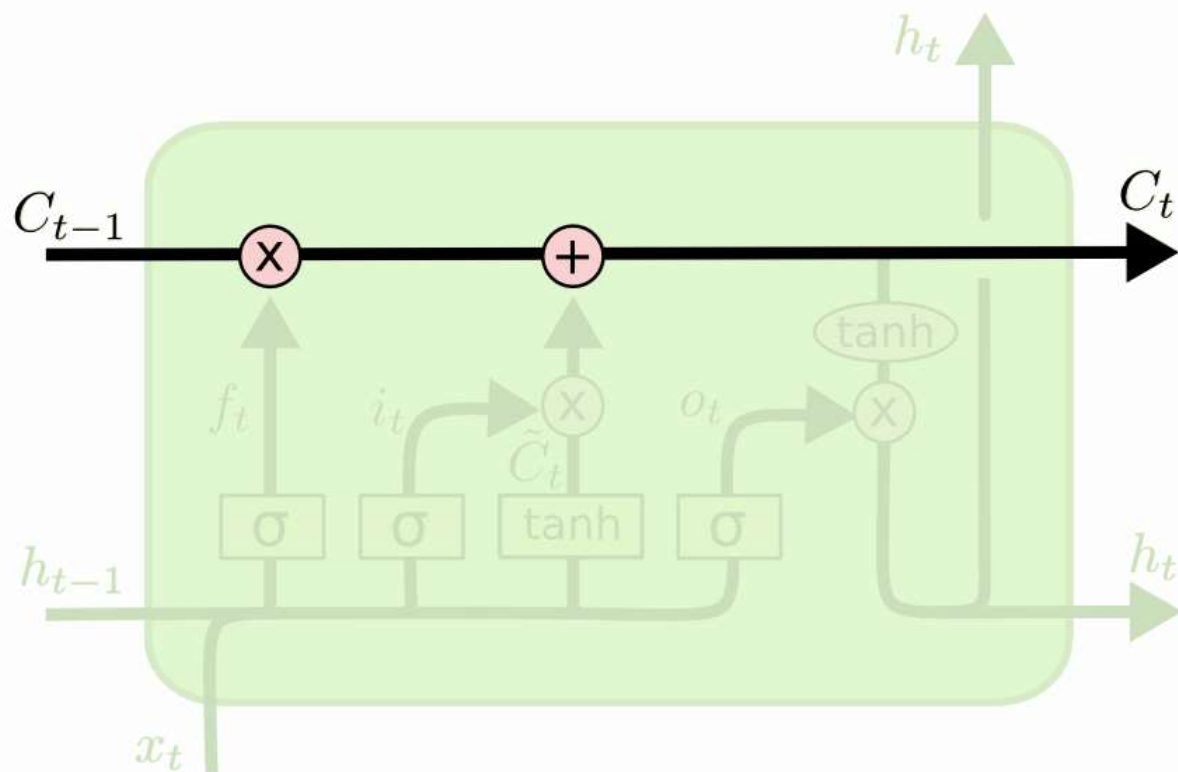


# LSTM

---

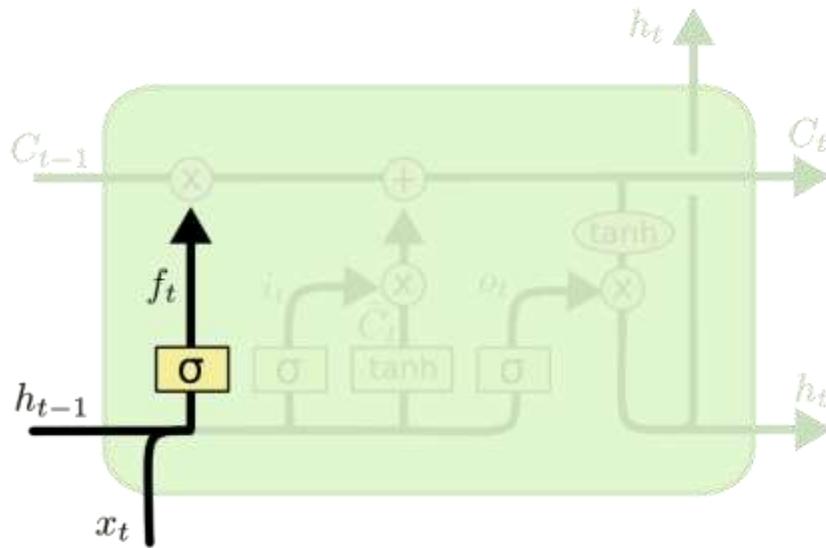


# Ключевая идея



# Ворота забывания

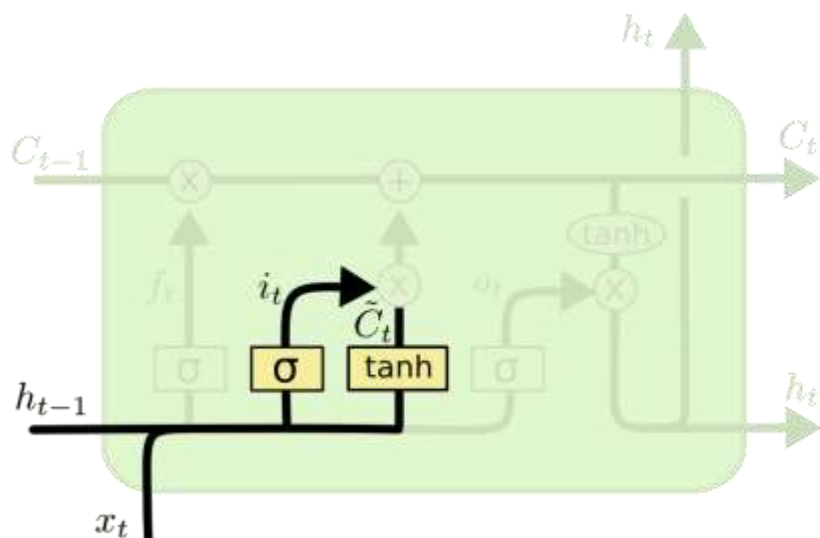
Нужно понять, что необходимо забыть, а что оставить в памяти. Эту подзадачу можно интерпретировать как бинарную классификацию. Берётся вектор краткосрочной памяти, добавляется новый ввод, всё это умножается на матрицу забывания. В результате будет получен вектор из чисел между 0 и 1. Далее путем поэлементного умножения на вектор долгосрочной памяти часть информации удалится (умножение на ноль), часть останется (умножение на 1).



$$f_t = \sigma (W_f \cdot [h_{t-1}, x_t] + b_f)$$

# Ворота входа

Нужно решить, что сохранить в памяти: имеется два потока (для первого используется бинарная классификация для регулирования того, какую часть второго потока добавить. После этого данные потоки перемножаются и добавляются в долгосрочную память.

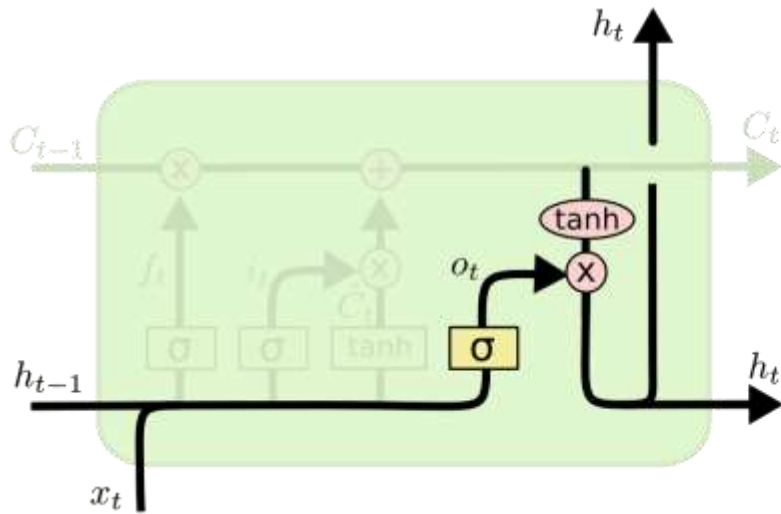


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] + b_C)$$

# Ворота выхода

Сколько информации из cell state следует отдавать на выход из LSTM-блока. В краткосрочную память добавляется новый вход, а затем переносится часть информации из долгосрочной памяти в краткосрочную.



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

$$h_t = o_t * \tanh (C_t)$$



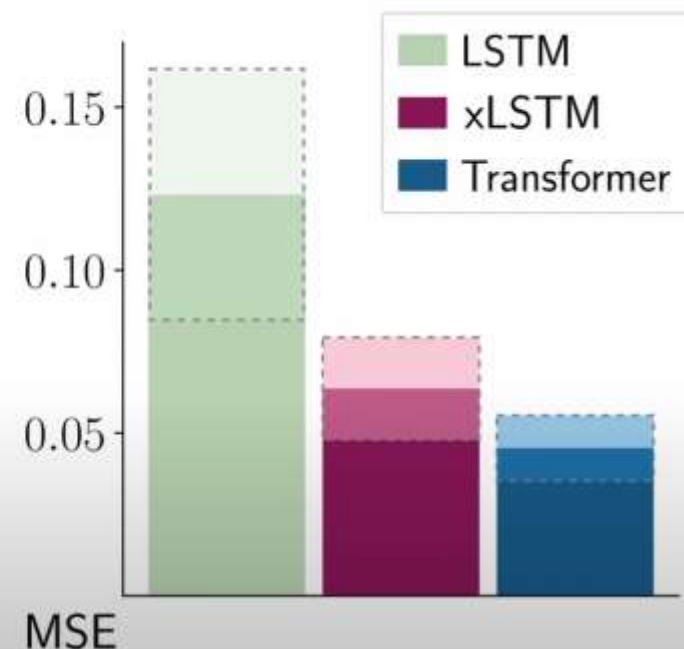
# Недостатки LSTM

Неспособность пересмотреть решение о сохранении данных внутри ячейки памяти. Это демонстрируют на простой задаче поиска ближайшего соседа (Nearest Neighbor Search), где сначала даётся референсный вектор, а далее сканируется последовательность других векторов, и модель должна найти наиболее похожий вектор и вернуть связанное с ним значение, когда последовательность закончится. Когда в последовательности попадает ещё более подходящий вектор, то модель не справляется.

## Nearest Neighbor Search Problem:

Search Key	Search Values				
5	2	1.5	4	5.5	9
Price Tag:	12	17	14	11	16
Prediction:	12	12	14	11	11

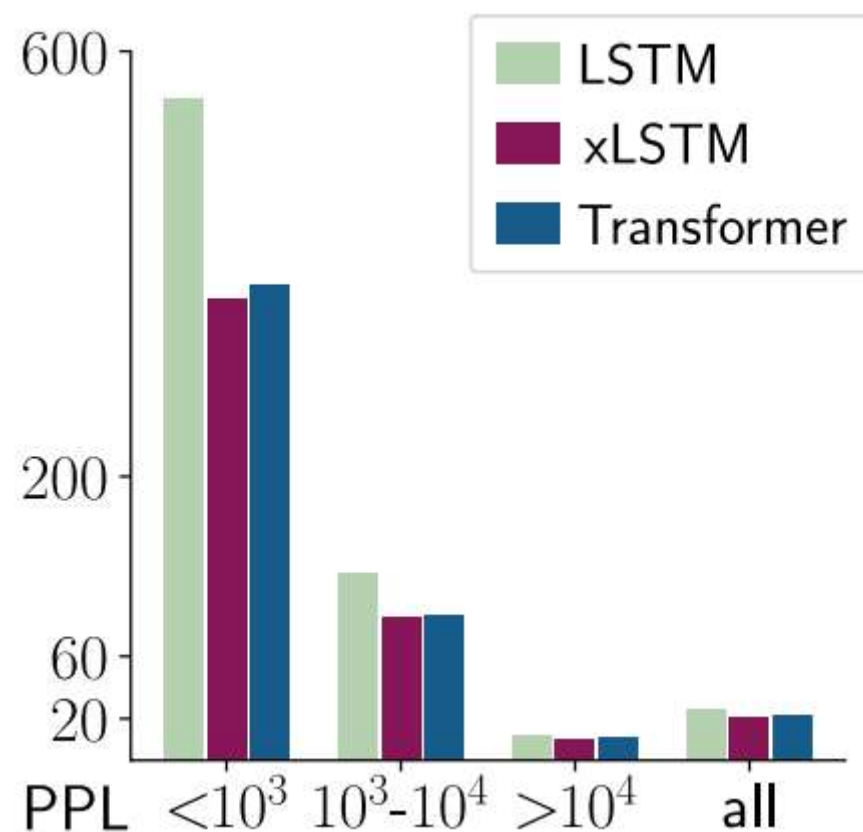
Goal: "Predict the price of the closest value to the key"



# Недостатки LSTM

Ограниченная память — всё надо впихнуть внутрь скаляра, который хранится в ячейке памяти LSTM. Это демонстрируют на задаче предсказания редкого токена (Rare Token Prediction)

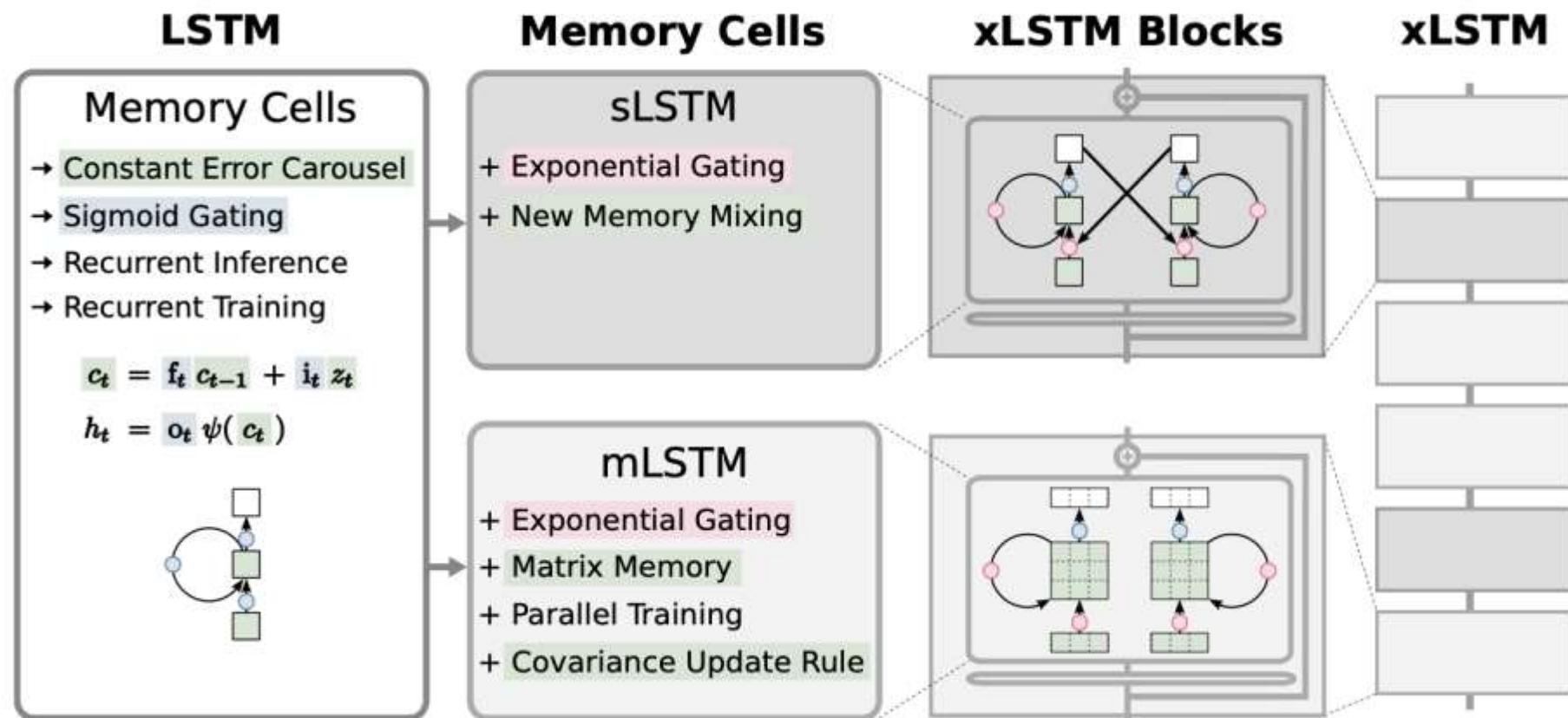
Плохая параллелизация в силу последовательной обработки скрытых состояний ячейки между соседними временными отсчётами (состояние зависит от предыдущего через hidden-hidden связи)





# xLSTM

---



# Решения в xLSTM

Exponential Gates: добавляются экспоненциальные функции активации на input и forget gate. Также появляется отдельное состояние нормализатора, и чтобы от экспоненты всё не взлетело в космос - состояние стабилизатора.

$$c_t = f_t c_{t-1} + i_t z_t \quad \text{cell state} \quad (8)$$

$$n_t = f_t n_{t-1} + i_t \quad \text{normalizer state} \quad (9)$$

$$h_t = o_t \tilde{h}_t, \quad \tilde{h}_t = c_t / n_t \quad \text{hidden state} \quad (10)$$

$$z_t = \varphi(\tilde{z}_t), \quad \tilde{z}_t = \mathbf{w}_z^\top \mathbf{x}_t + r_z h_{t-1} + b_z \quad \text{cell input} \quad (11)$$

$$i_t = \exp(\tilde{i}_t), \quad \tilde{i}_t = \mathbf{w}_i^\top \mathbf{x}_t + r_i h_{t-1} + b_i \quad \text{input gate} \quad (12)$$

$$f_t = \sigma(\tilde{f}_t) \text{ OR } \exp(\tilde{f}_t), \quad \tilde{f}_t = \mathbf{w}_f^\top \mathbf{x}_t + r_f h_{t-1} + b_f \quad \text{forget gate} \quad (13)$$

$$o_t = \sigma(\tilde{o}_t), \quad \tilde{o}_t = \mathbf{w}_o^\top \mathbf{x}_t + r_o h_{t-1} + b_o \quad \text{output gate} \quad (14)$$

$$m_t = \max(\log(f_t) + m_{t-1}, \log(i_t)) \quad \text{stabilizer state} \quad (15)$$

$$i'_t = \exp(\log(i_t) - m_t) = \exp(\tilde{i}_t - m_t) \quad \text{stabil. input gate} \quad (16)$$

$$f'_t = \exp(\log(f_t) + m_{t-1} - m_t) \quad \text{stabil. forget gate} \quad (17)$$

# Решения в xLSTM

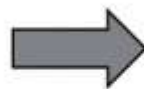
Scalar Pre-activations:

$$\tilde{z}_t = w_z^\top x_t + r_z h_{t-1} + b_z$$

$$\tilde{i}_t = w_i^\top x_t + r_i h_{t-1} + b_i$$

$$\tilde{f}_t = w_f^\top x_t + r_f h_{t-1} + b_f$$

$$\tilde{o}_t = w_o^\top x_t + r_o h_{t-1} + b_o$$



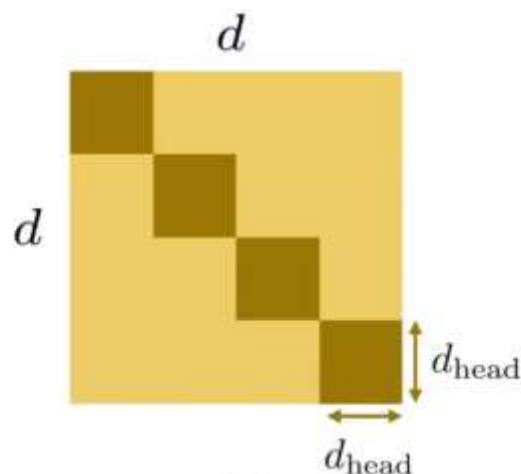
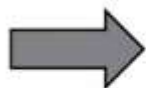
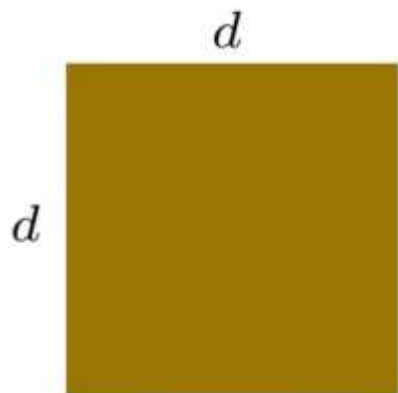
Vector Pre-activations:

$$\tilde{z}_t = W_z x_t + R_z h_{t-1} + b_z$$

$$\tilde{i}_t = W_i x_t + R_i h_{t-1} + b_i$$

$$\tilde{f}_t = W_f x_t + R_f h_{t-1} + b_f$$

$$\tilde{o}_t = W_o x_t + R_o h_{t-1} + b_o$$

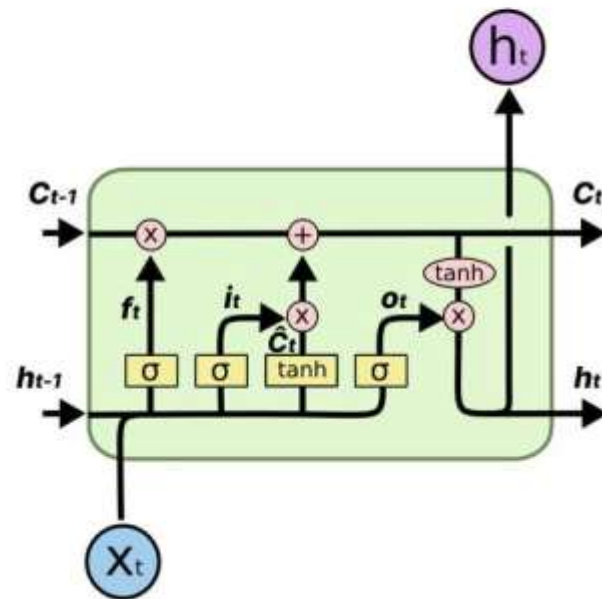


$R$

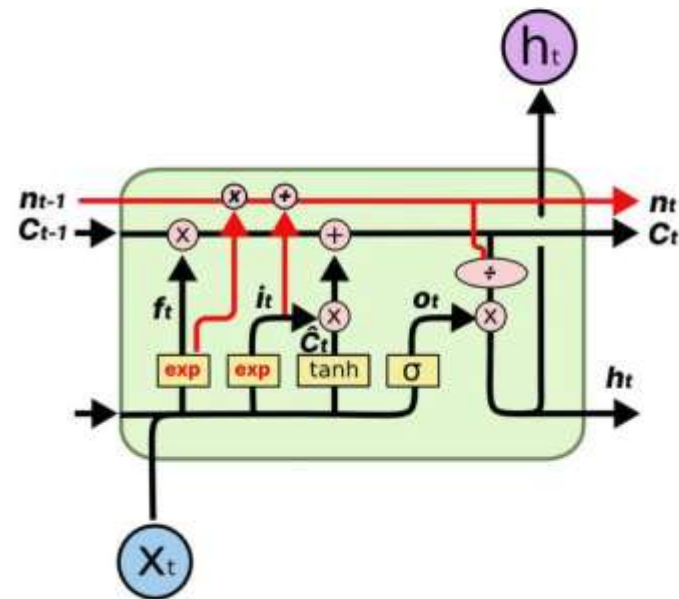
$R$

# Решения в xLSTM

Блок LSTM



Блок sLSTM



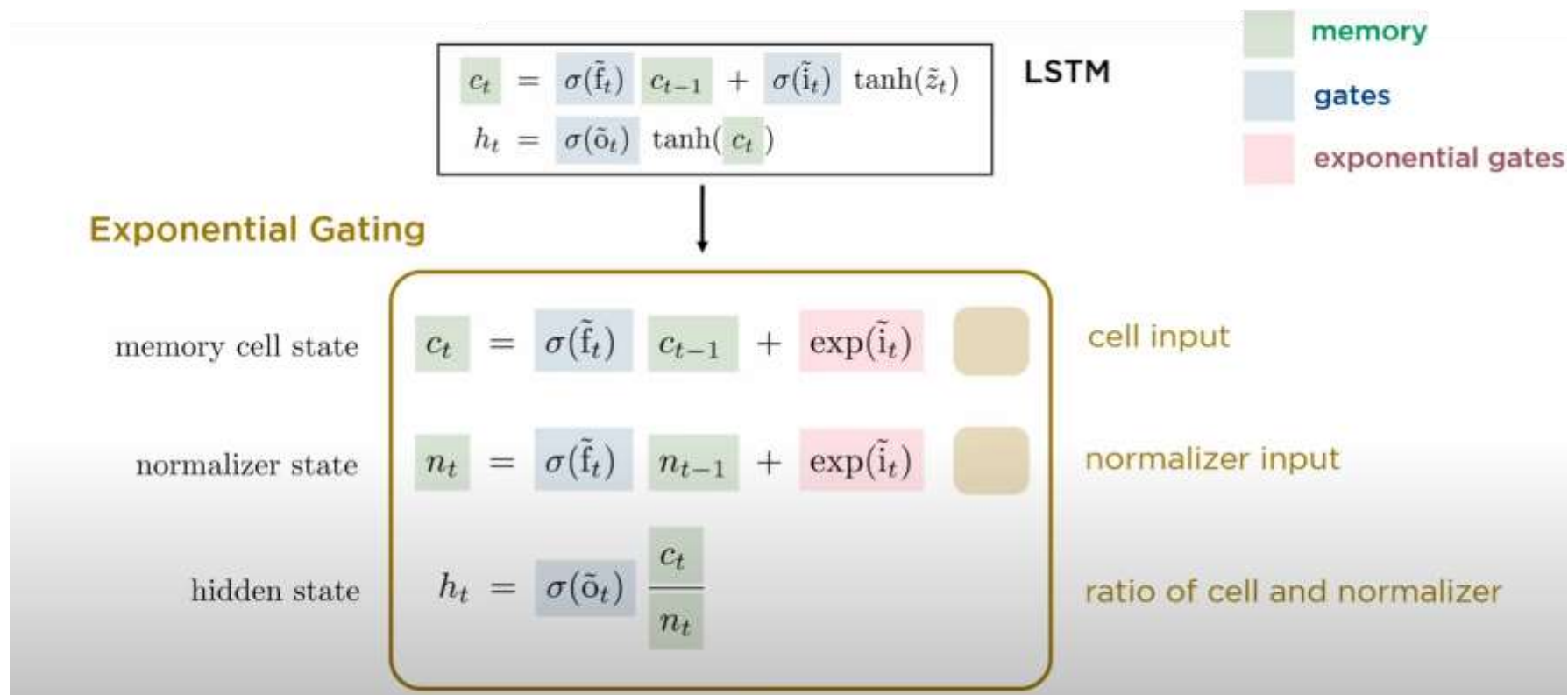
# Решения в xLSTM

Matrix memory: не что иное, как ассоциативная память; аналогично терминологии трансформера, создают запросы, ключи и значения qkv и получают hidden state в качестве значений v наиболее близких к запросу q. В формулу обновления ассоциативной памяти так же включены input и forget gates чтобы сделать ее более "lstm-like" это все параллелизуется

$$\begin{aligned}
 C_t &= f_t C_{t-1} + i_t v_t k_t^\top && \text{cell state (19)} \\
 n_t &= f_t n_{t-1} + i_t k_t && \text{normalizer state (20)} \\
 h_t &= o_t \odot \tilde{h}_t, & \tilde{h}_t &= C_t q_t / \max \left\{ \left| n_t^\top q_t \right|, 1 \right\} && \text{hidden state (21)} \\
 q_t &= W_q x_t + b_q && \text{query input (22)} \\
 k_t &= \frac{1}{\sqrt{d}} W_k x_t + b_k && \text{key input (23)} \\
 v_t &= W_v x_t + b_v && \text{value input (24)} \\
 i_t &= \exp(\tilde{i}_t), & \tilde{i}_t &= w_i^\top x_t + b_i && \text{input gate (25)} \\
 f_t &= \sigma(\tilde{f}_t) \text{ OR } \exp(\tilde{f}_t), & \tilde{f}_t &= w_f^\top x_t + b_f && \text{forget gate (26)} \\
 o_t &= \sigma(\tilde{o}_t), & \tilde{o}_t &= W_o x_t + b_o && \text{output gate (27)}
 \end{aligned}$$



# Решения в xLSTM





# Решения в xLSTM

Exponential  
Gating

$$\begin{aligned} c_t &= \sigma(\tilde{f}_t) c_{t-1} + \exp(\tilde{i}_t) \\ n_t &= \sigma(\tilde{f}_t) n_{t-1} + \exp(\tilde{i}_t) \\ h_t &= \sigma(\tilde{o}_t) \frac{c_t}{n_t} \end{aligned}$$

sLSTM

$$\begin{aligned} c_t &= \sigma(\tilde{f}_t) c_{t-1} + \exp(\tilde{i}_t) \tanh(\tilde{z}_t) & c_t \in \mathbb{R} \\ n_t &= \sigma(\tilde{f}_t) n_{t-1} + \exp(\tilde{i}_t) & n_t \in \mathbb{R} \\ h_t &= \sigma(\tilde{o}_t) \frac{c_t}{n_t} & h_t \in \mathbb{R} \end{aligned}$$

“scalar” cell state

mLSTM

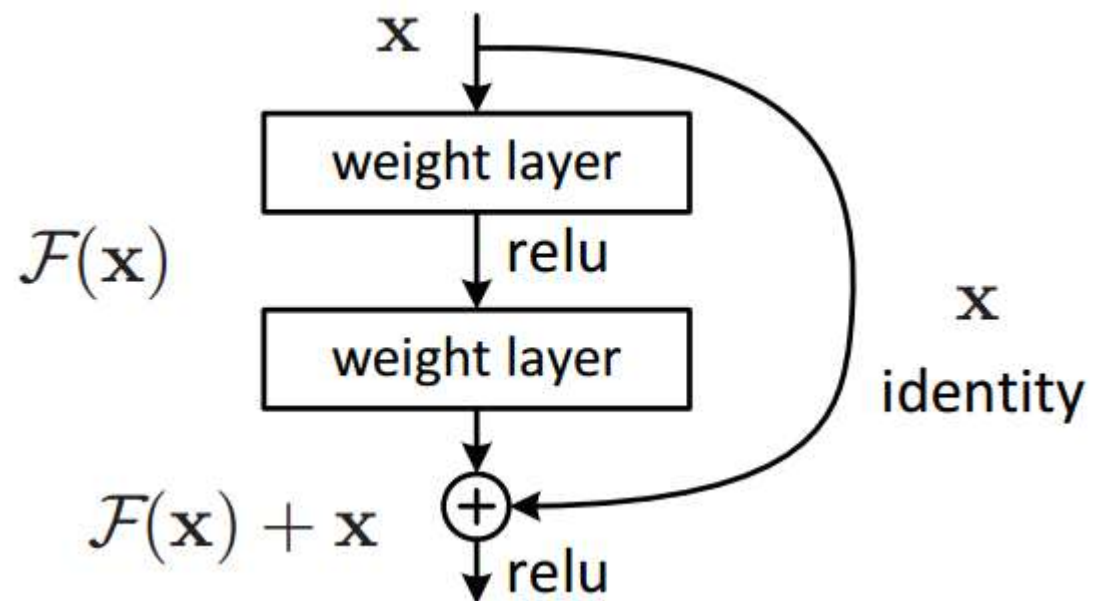
$$\begin{aligned} C_t &= \sigma(\tilde{f}_t) C_{t-1} + \exp(\tilde{i}_t) v_t k_t^\top & C_t \in \mathbb{R}^{d \times d} \\ n_t &= \sigma(\tilde{f}_t) n_{t-1} + \exp(\tilde{i}_t) k_t & n_t \in \mathbb{R}^d \\ h_t &= \sigma(\tilde{o}_t) \odot \frac{C_t q_t}{\max\{|n_t^\top q_t|, 1\}} & h_t \in \mathbb{R}^d \end{aligned}$$

“matrix” cell state

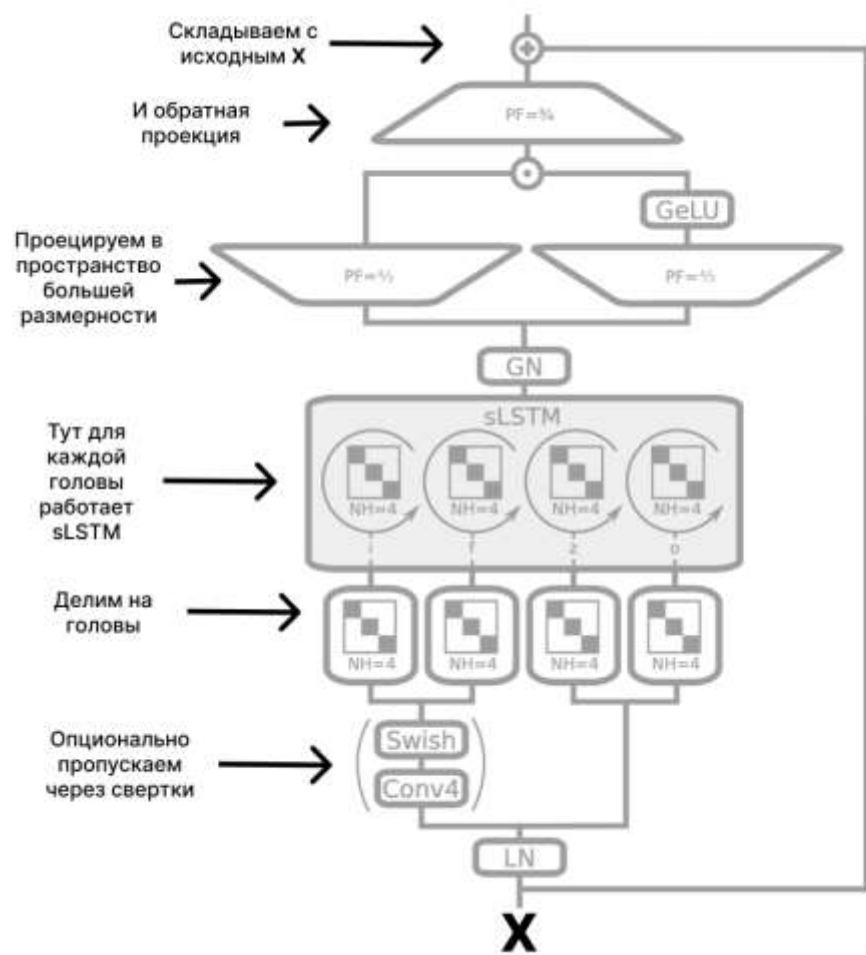
# Как все собрать в xLSTM?

---

Дополнительно обернули  
каждую из структур в  
residual блоки

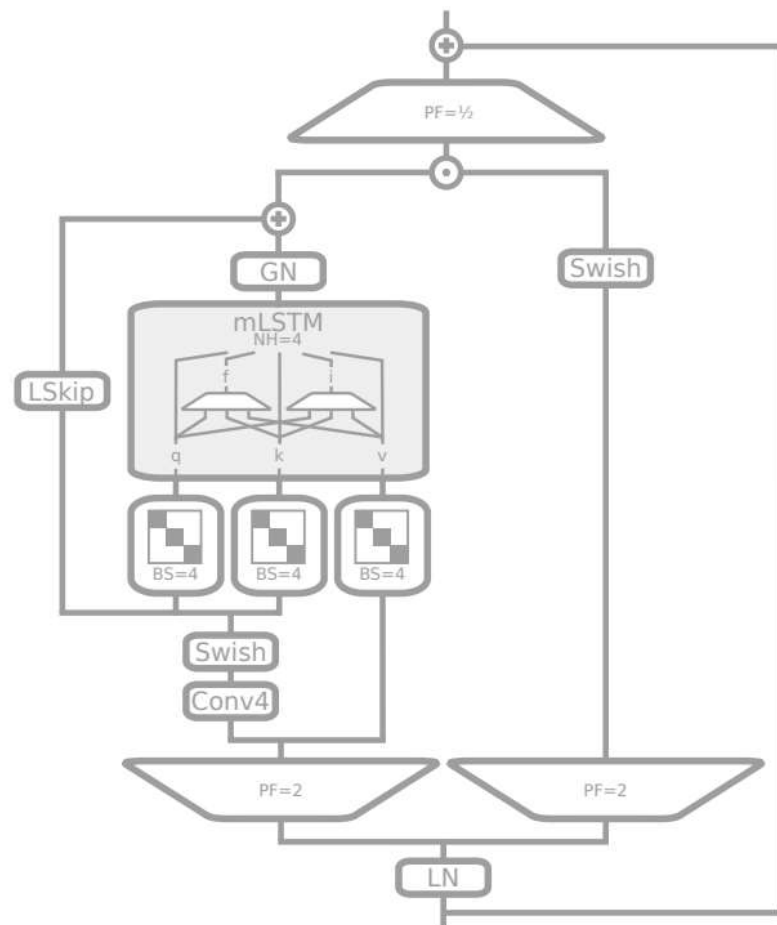


# Как все собрать в xLSTM?



Входы  $X$  делятся на головы – на рисунке ниже их четыре. Перед этим входы опционально могут быть пропущены через несколько сверточных слоев. После деления в каждой из голов отрабатывает sLSTM. Выходы с голов затем объединяются с помощью GroupNorm, а потом проецируются в пространство большей размерности и обратно. Эта последняя часть названа **up-projection** и проделывается для того, что повысить качество histories seperating. Это умение сети разделять "линии повествования": например, понимать, к какому из персонажей относится то или иное действие.

# Как все собрать в xLSTM?



Все то же самое, что и в sLSTM, но в другом порядке. Отражаем входы в пространство большей размерности  $\rightarrow$  делим на головы  $\rightarrow$  пропускаем через mLSTM  $\rightarrow$  объединяем по GroupNorm  $\rightarrow$  проецируем обратно в родную размерность  $\rightarrow$  складываем со входами, чтобы получится residual  $\rightarrow$  готово!

# Результаты

Ablation studies on the new xLSTM components.

Model	Modification	Exponential Gating	Matrix Memory	#Params M	SlimPajama (15B) ppl ↓
LSTM	Vanilla Multi-Layer LSTM	✗	✗	607.8	2417.86
	Adding Resnet Backbone	✗	✗	506.1	35.46
	Adding Up-Projection Backbone	✗	✗	505.9	26.01
xLSTM[0:1]	Adding Exponential Gating	✓	✗	427.3	17.70
xLSTM[7:1]	Adding Matrix Memory	✓	✓	408.4	<b>13.48</b>

# Результаты

Model	#Params M	SlimPajama (15B) ppl ↓
<b>GPT-3</b>	356	14.26
Llama	407	<u>14.25</u>
H3	420	18.23
Mamba	423	<u>13.70</u>
Hyena	435	17.59
RWKV-4	430	<u>15.62</u>
RWKV-5	456	16.53
RWKV-6	442	17.40
RetNet	431	16.23
HGRN	411	21.83
GLA	412	19.56
HGRN2	411	16.77
<b>xLSTM[1:0]</b>	409	<b><u>13.43</u></b>
<b>xLSTM[7:1]</b>	408	<b><u>13.48</u></b>

