



LSTM - Сети с Долгой краткосрочной памятью

Преподаватель: Герард Костин

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

Проблемы RNN

- Ошибки, распространяющиеся в обратном направлении, умножаются на каждом слое, что приводит к экспоненциальному убыванию (если производная мала) или росту (если производная большая).
- Очень сложно реализовать обучение глубоких сетей или простых RNN сетей на большой временной глубине
- Очень сложно выучить зависимости на большом расстоянии, такие как согласование субъект-глагол.

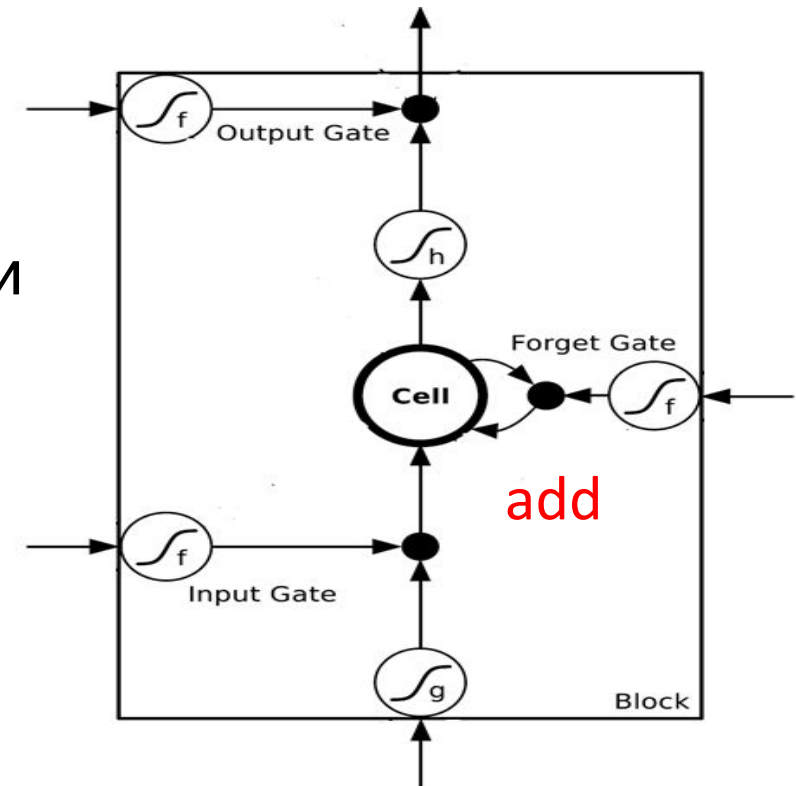
Как побороть проблемы RNN

- Long Short-term Memory (LSTM)

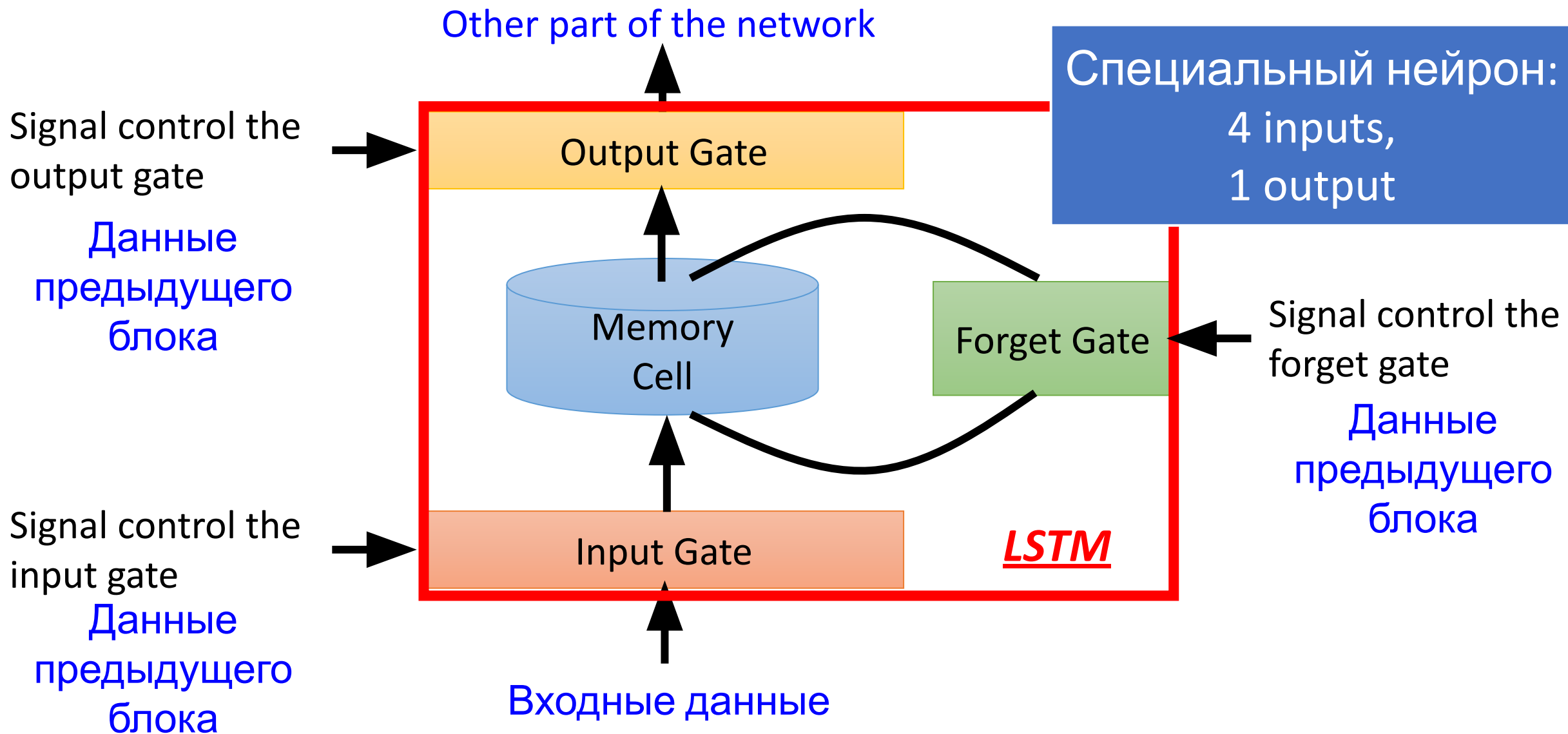
- Могут работать с исчезающим градиентом (НЕ! gradient explode)

- Память и Вход складываются (*added*)
- Информация никогда не теряется, если вентиль “забвения” закрыт

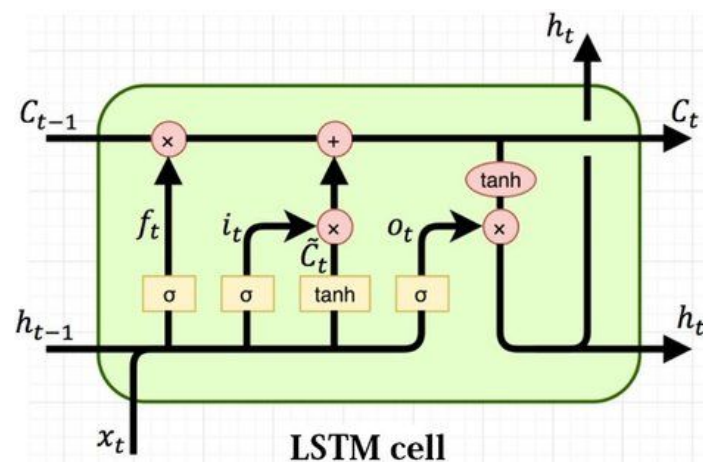
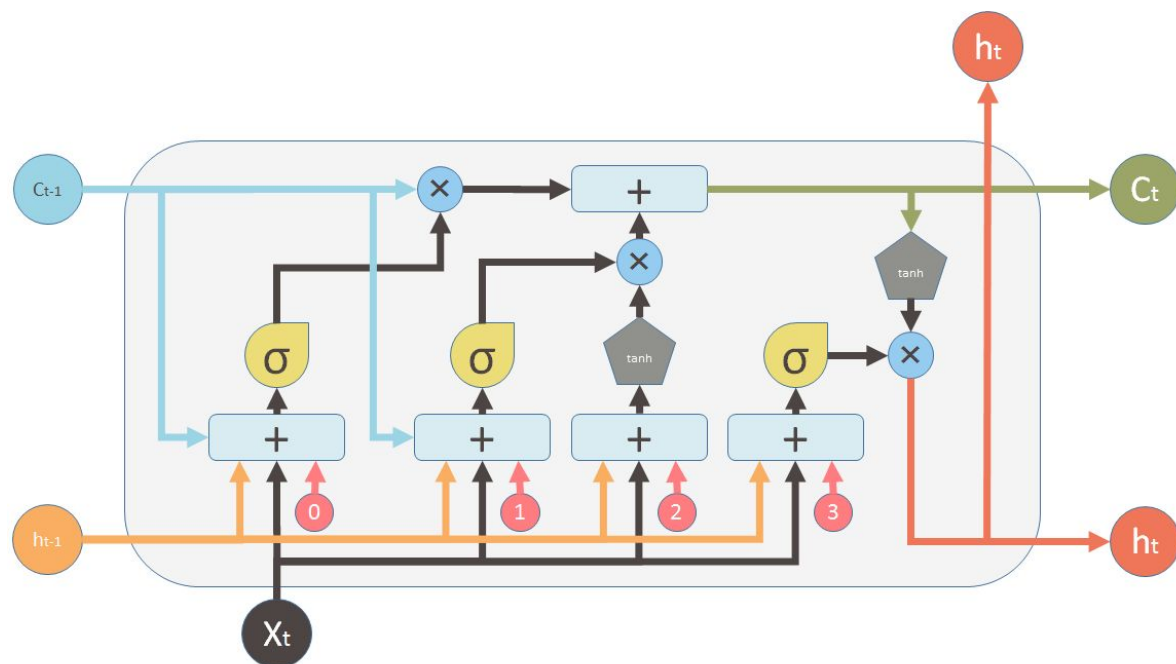
➡ Нет Gradient vanishing
(Если вентиль “забвения” открыт.)



Long Short-term Memory (LSTM)



Long Short-term Memory (LSTM)



$$\begin{aligned}
 i_t &= \sigma(x_t U^i + h_{t-1} W^i) \\
 f_t &= \sigma(x_t U^f + h_{t-1} W^f) \\
 o_t &= \sigma(x_t U^o + h_{t-1} W^o) \\
 \tilde{C}_t &= \tanh(x_t U^g + h_{t-1} W^g) \\
 C_t &= \sigma(f_t * C_{t-1} + i_t * \tilde{C}_t) \\
 h_t &= \tanh(C_t) * o_t
 \end{aligned}$$

Inputs:

- x_t Input vector
- C_{t-1} Memory from previous block
- h_{t-1} Output of previous block

outputs:

- C_t Memory from current block
- h_t Output of current block

Nonlinearities:

- σ Sigmoid
- \tanh Hyperbolic tangent

Bias:

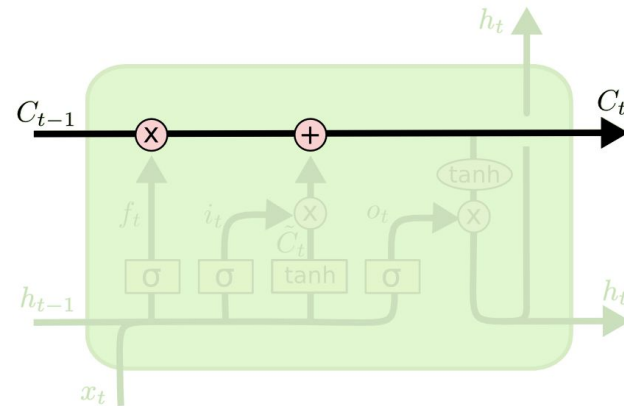
0

Vector operations:

- \times Element-wise multiplication
- $+$ Element-wise Summation / Concatenation

Cell State

- Обработывает вектор C_t который имеет ту же размерность, что и скрытый слой, h_t
- Информация может быть добавлена или удалена из этого вектора состояния через шлюзы забвения и ввода.

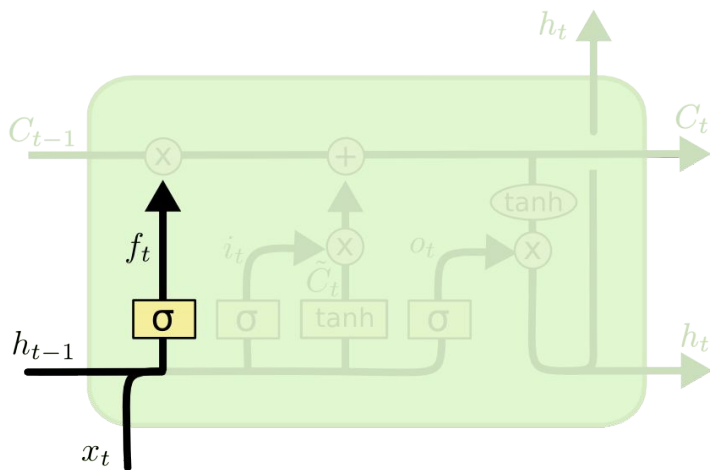


Cell State Пример

- Хотите запомнить лицо и порядок подлежащего существительного, чтобы его можно было проверить на соответствие объекту и глаголу, когда оно в конечном итоге встретится.
- Шлюз забвения удалит существующую информацию о предыдущем предмете при обнаружении нового.
- Входной шлюз «добавляет» информацию о новом предмете.

Forget Gate

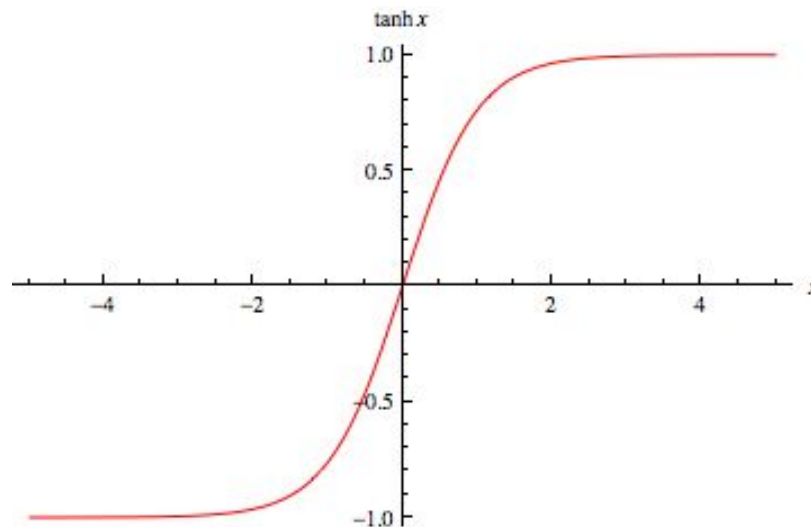
- Forget gate вычисляет значение 0-1, используя логистическую функцию вывода сигмоида из входных данных, x_t , и текущее скрытое состояние, h_t :
- Мультипликативно сочетается с состоянием ячейки, «забывая» информацию, когда клапан выводит что-то близкое к 0.



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

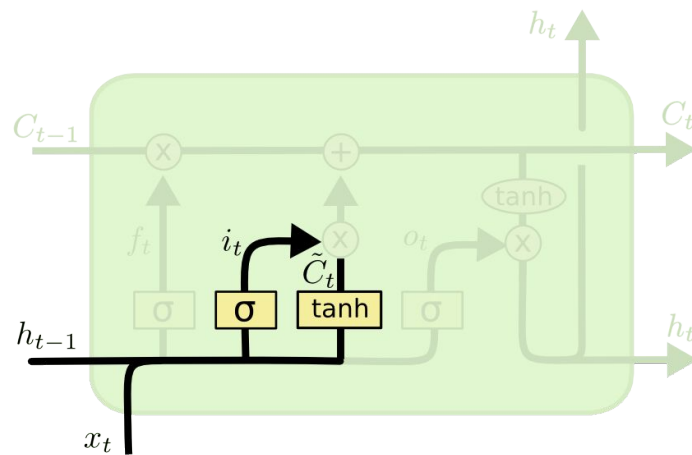
Hyperbolic Tangent Активация

- Tanh можно использовать в качестве альтернативной нелинейной функции сигмовидной логистической выходной функции (0-1).
- Используется для вывода пороговых значений от -1 до 1 .



Input Gate

- Во-первых, определяет, какие записи в состоянии ячейки обновлять, вычисляя сигмоидальный выход 0-1.
- Затем определяет, какую сумму добавить / вычесть из этих записей, вычислив функцию выхода \tanh (со значением от -1 до 1) входного и скрытого состояния.
-

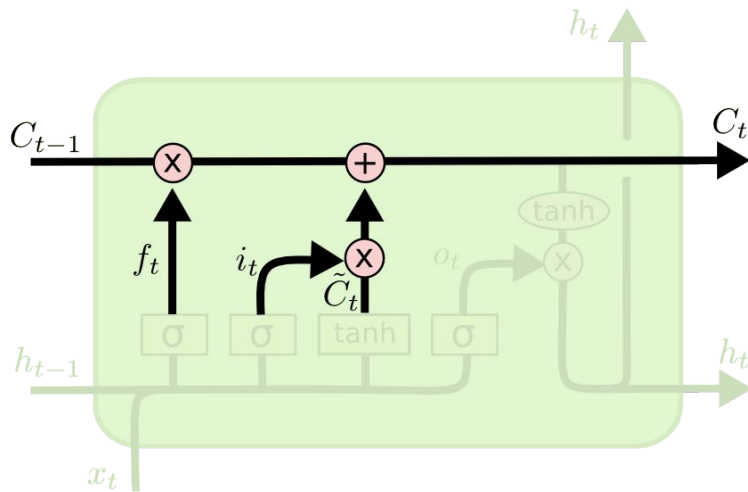


$$i_t = \sigma(W_i \cdot [h_{t-1}, x_t] + b_i)$$

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

Обновление Cell State

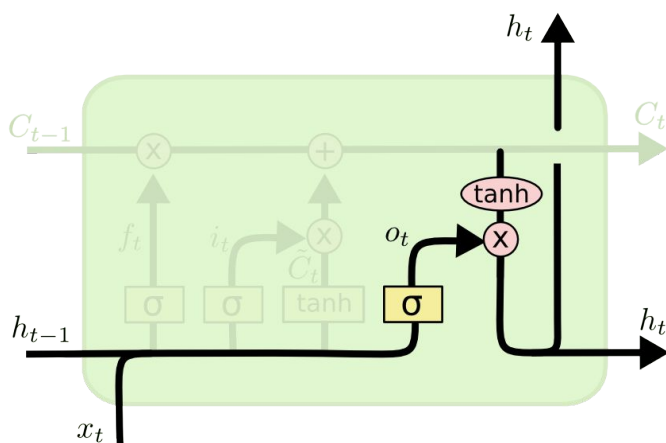
- Состояние ячейки обновляется с помощью покомпонентного векторного умножения, чтобы «забыть», и векторного сложения для «ввода» новой информации.



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

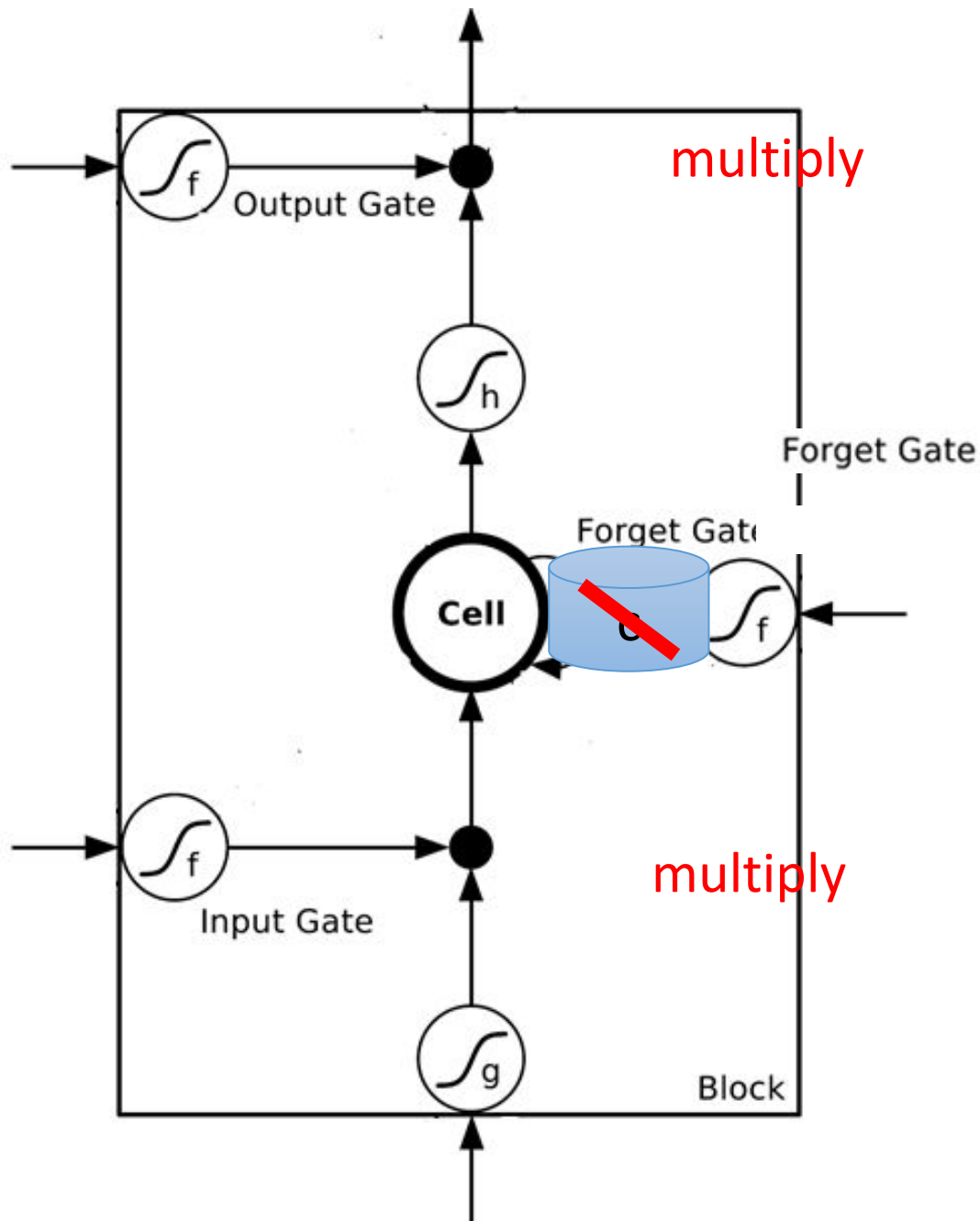
Output Gate

- Скрытое состояние обновляется на основе «отфильтрованной» версии состояния ячейки, масштабируемой от -1 до 1 с использованием \tanh .
- Выходной вентиль вычисляет сигмовидную функцию входа и текущего скрытого состояния, чтобы определить, какие элементы состояния ячейки нужно «выводить».



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

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



Функция активации - sigmoid function

Между 0 и 1

За счет этого определяется на сколько вентиль "забвения" пропускает информацию

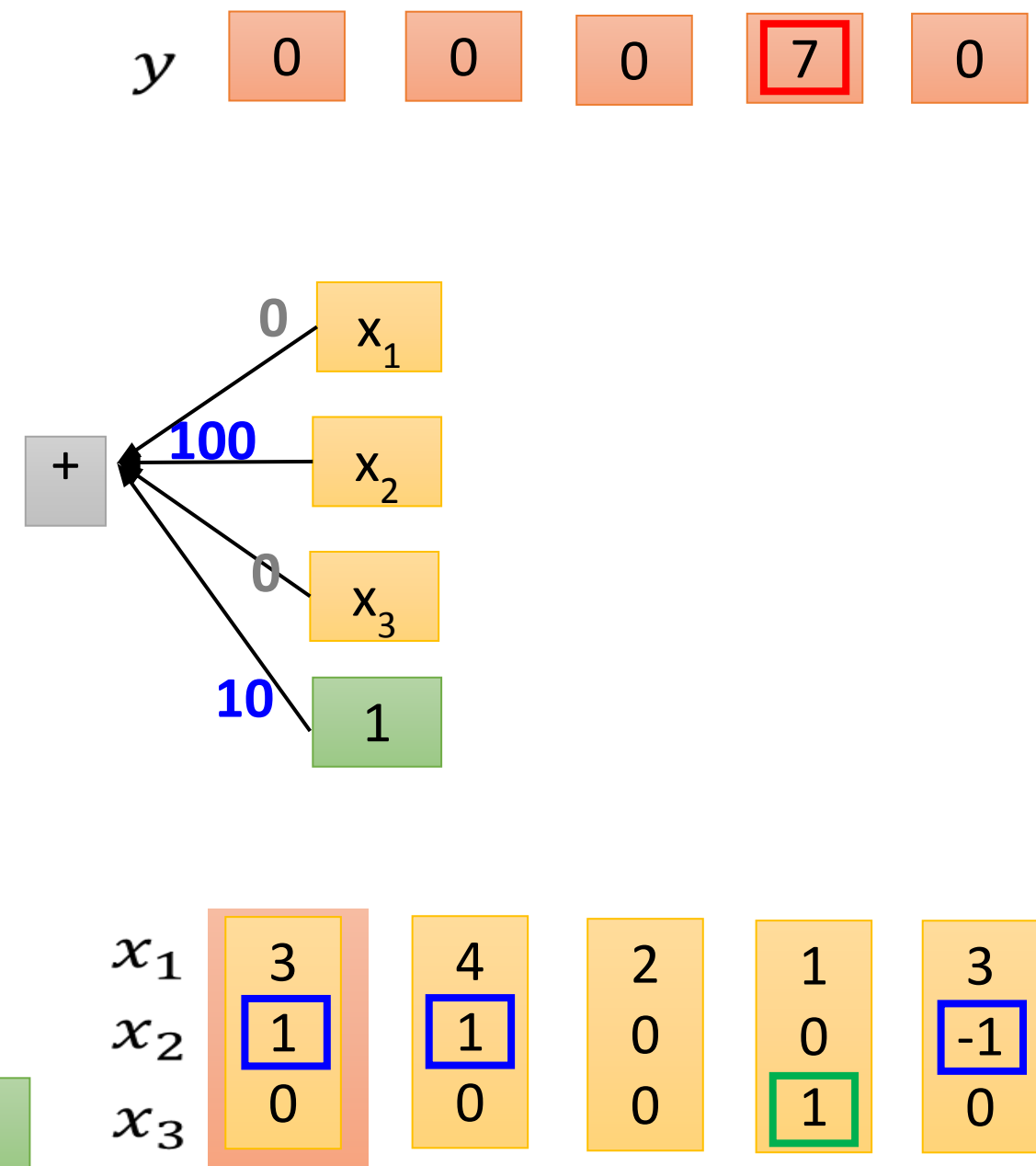
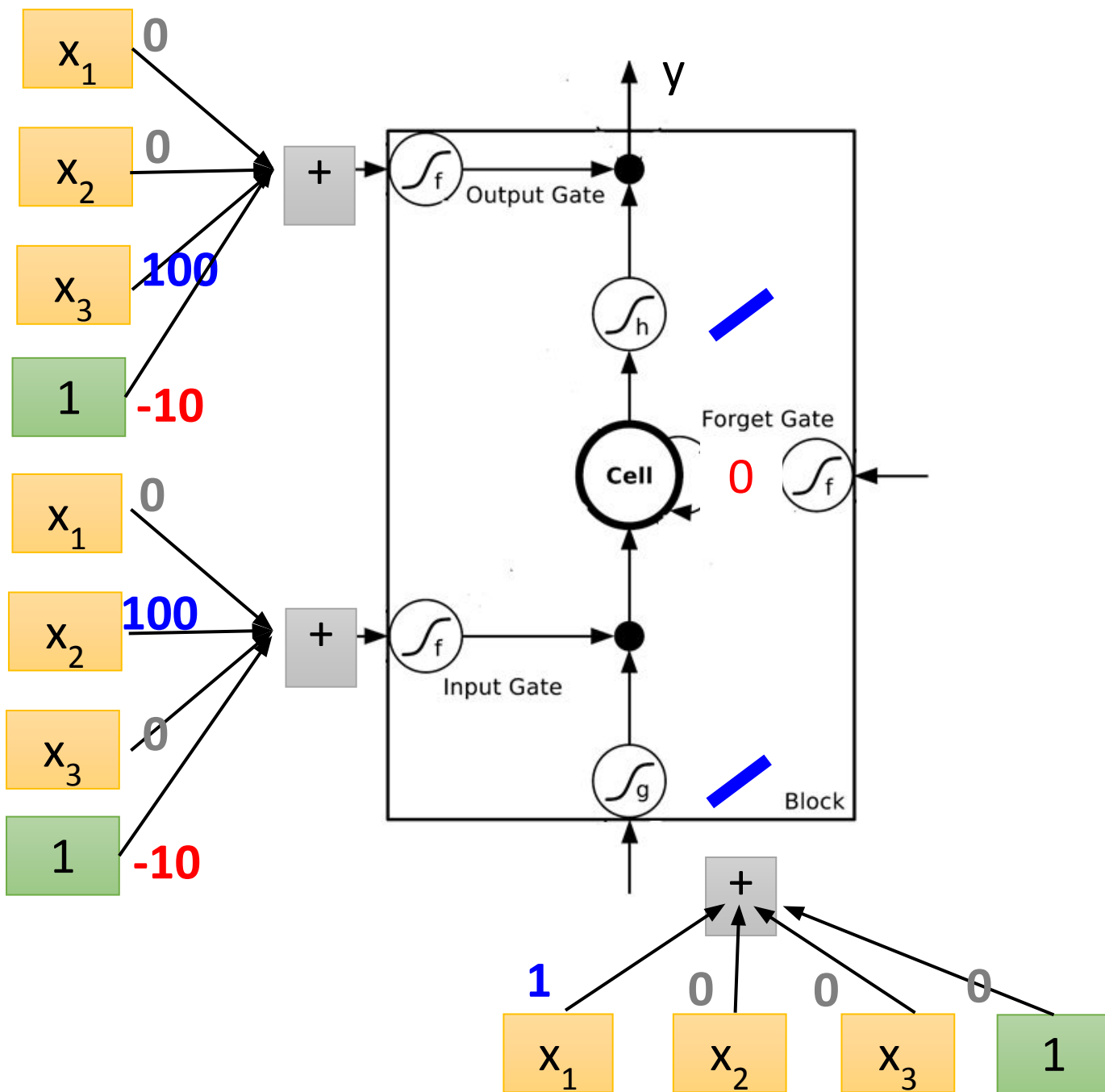
LSTM - Разбор

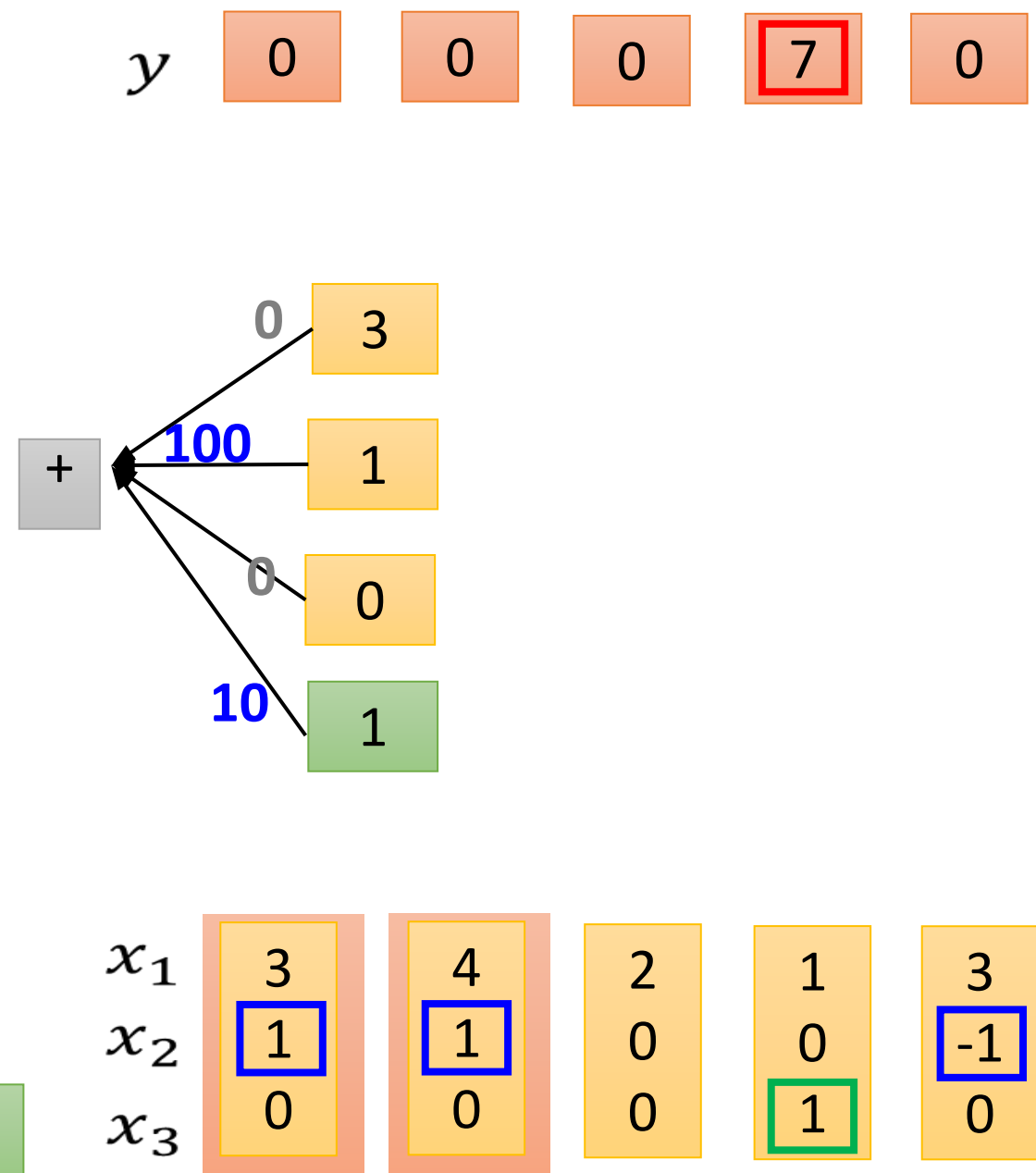
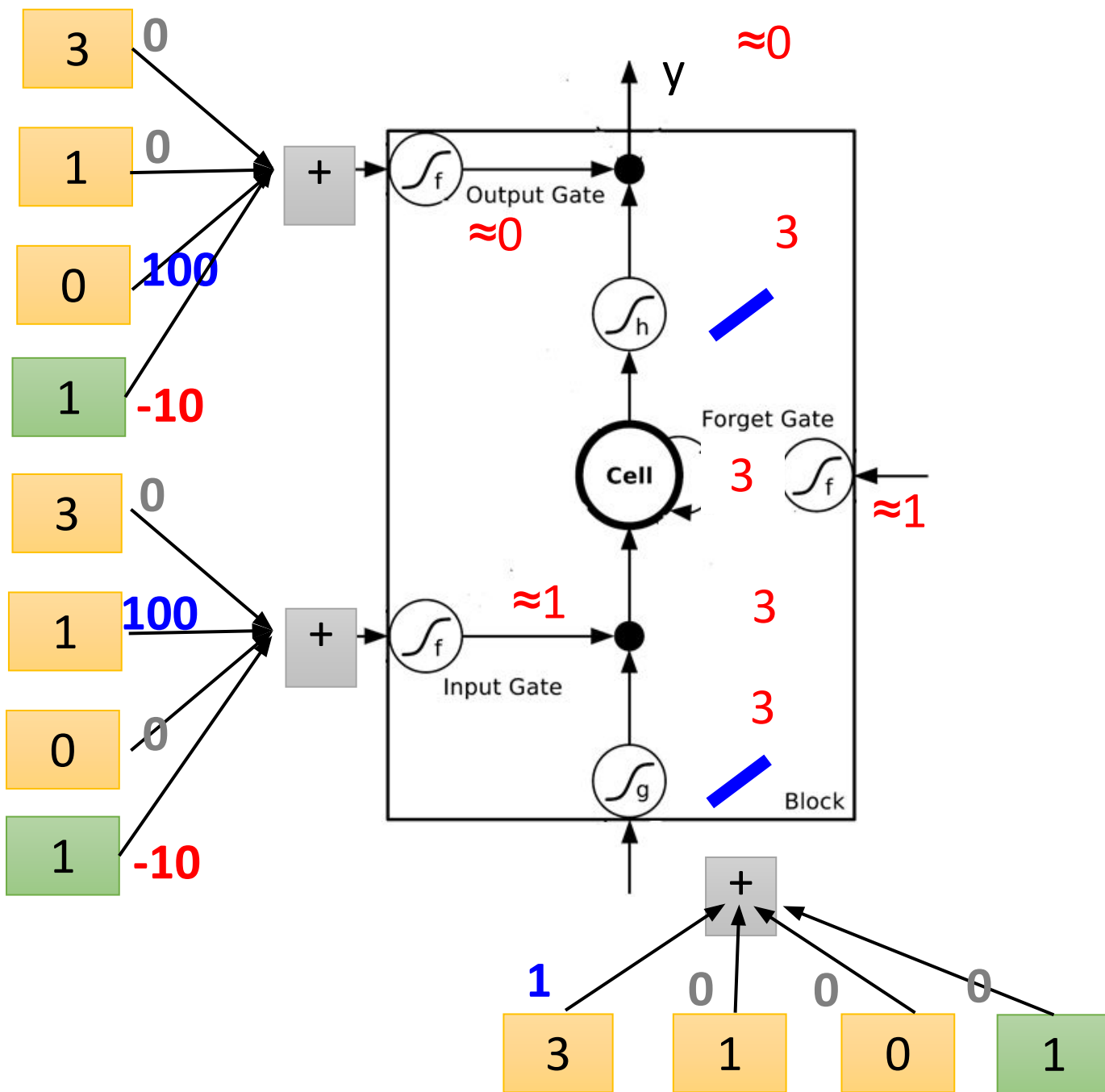
	0	0	3	3	7	7	7	0	6
x_1	1	3	2	4	2	1	3	6	1
x_2	0	1	0	1	0	0	-1	1	0
x_3	0	0	0	0	0	1	0	0	1
y	0	0	0	0	0	7	0	0	6

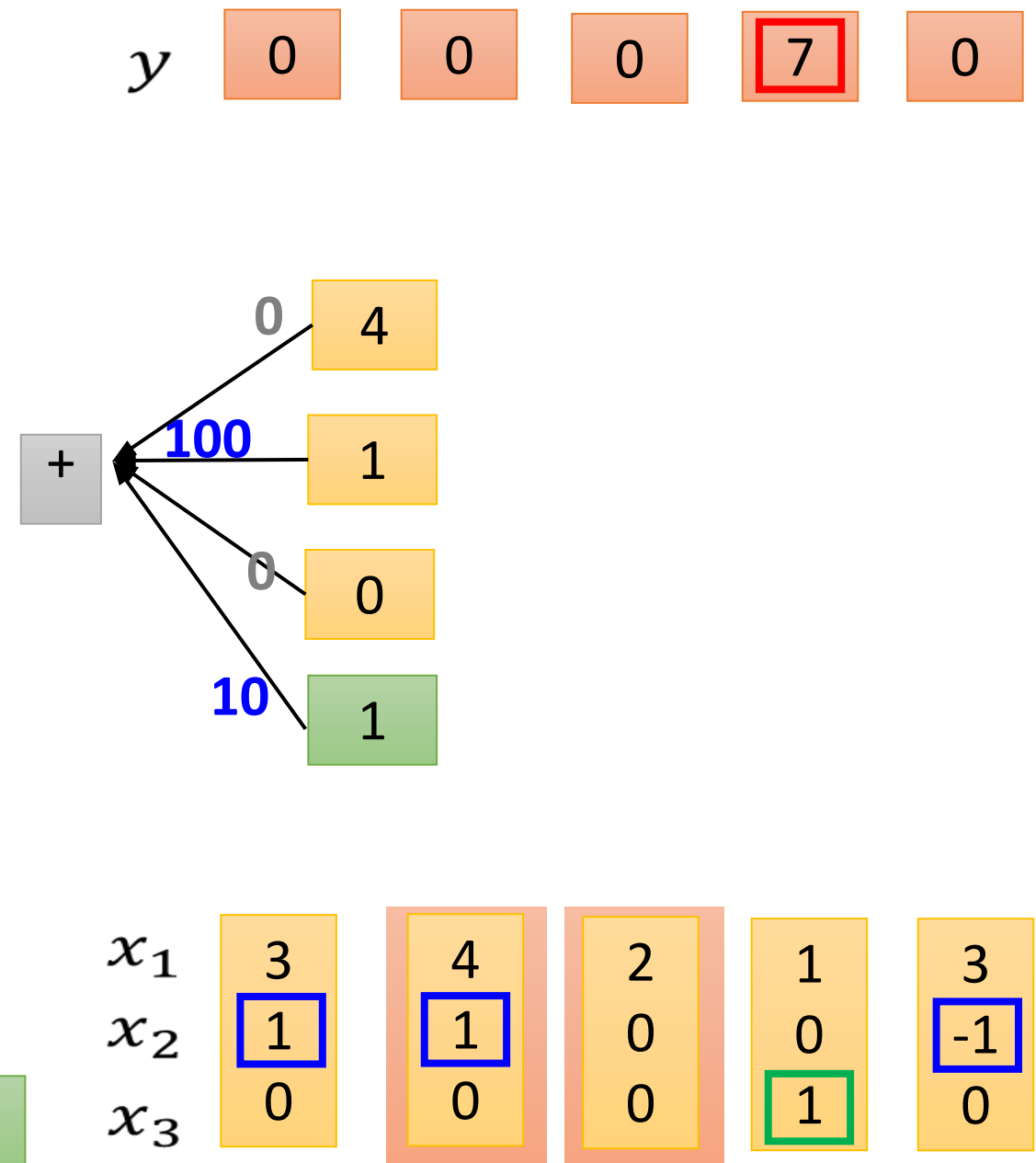
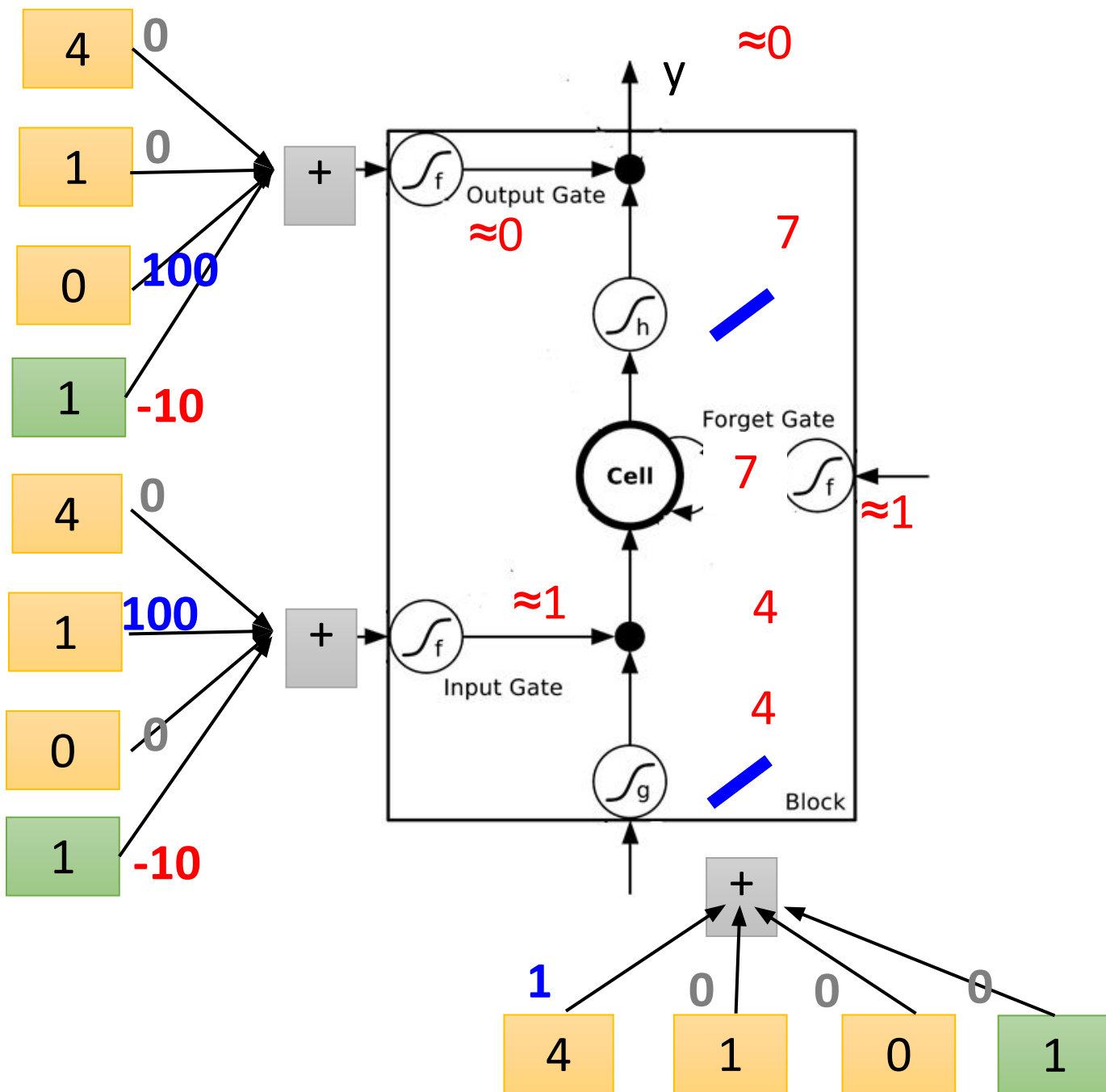
Когда $x_2 = 1$, добавляем x_1 в память

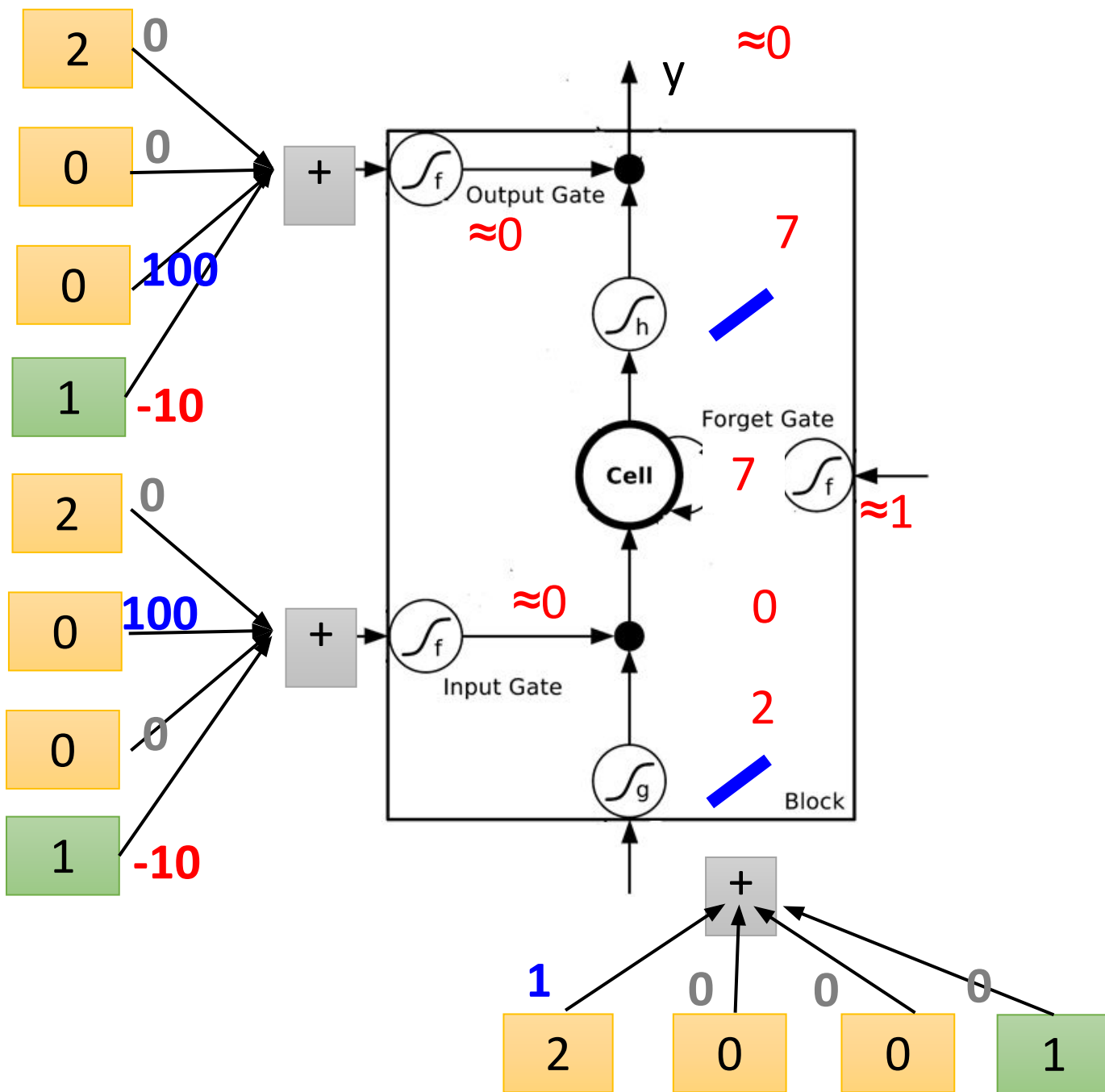
Когда $x_2 = -1$, очищаем память

Когда $x_3 = 1$, выводим данные в память

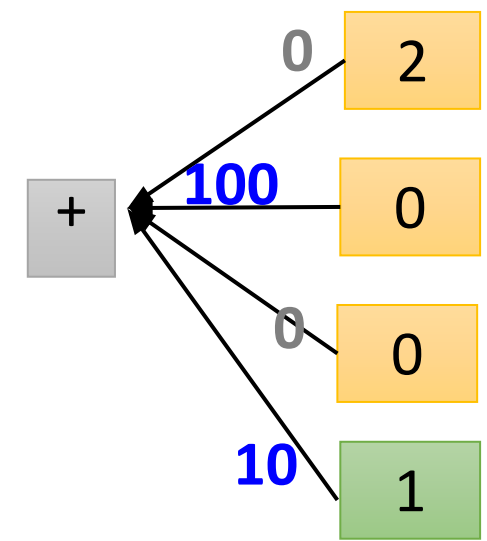




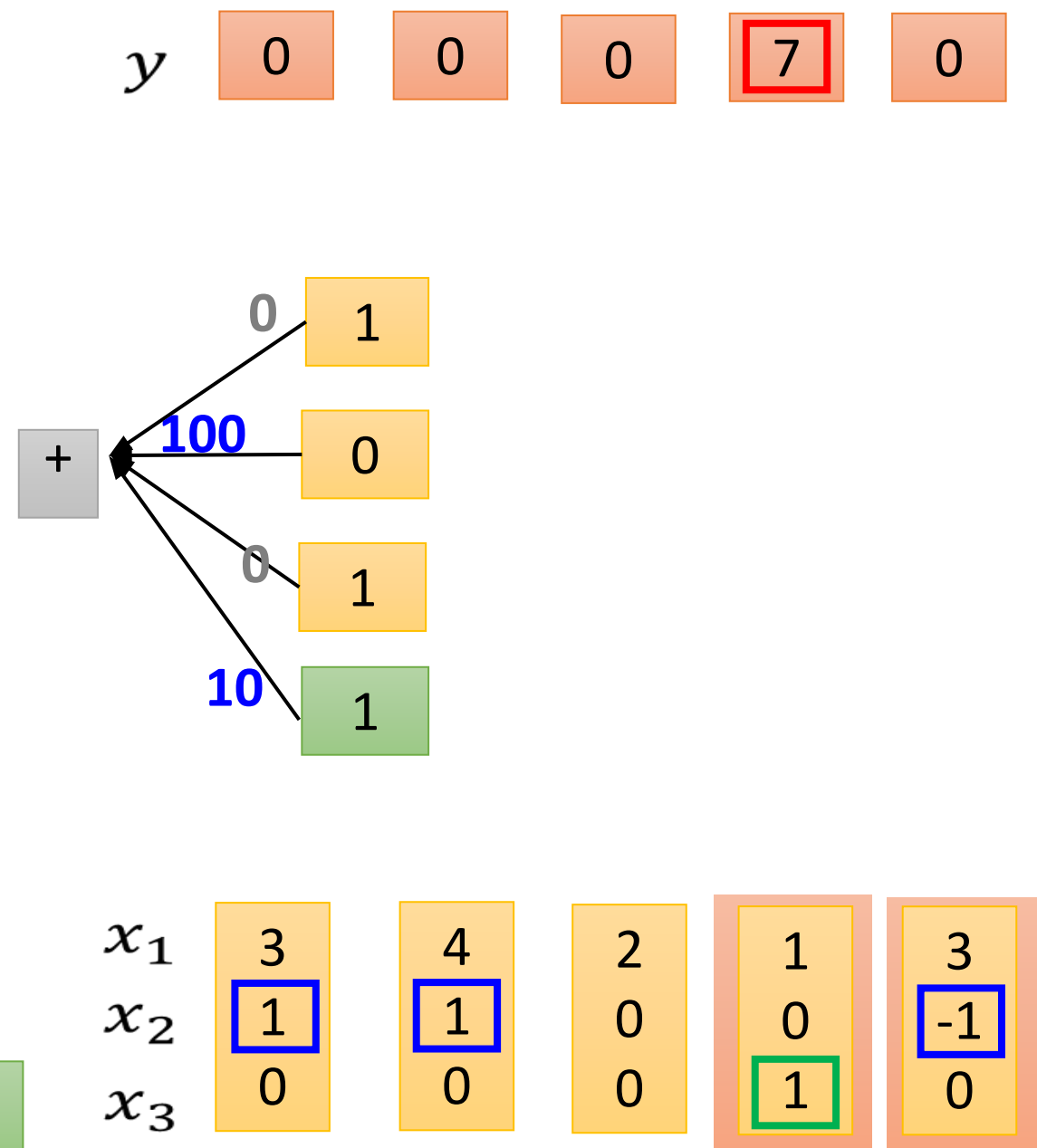
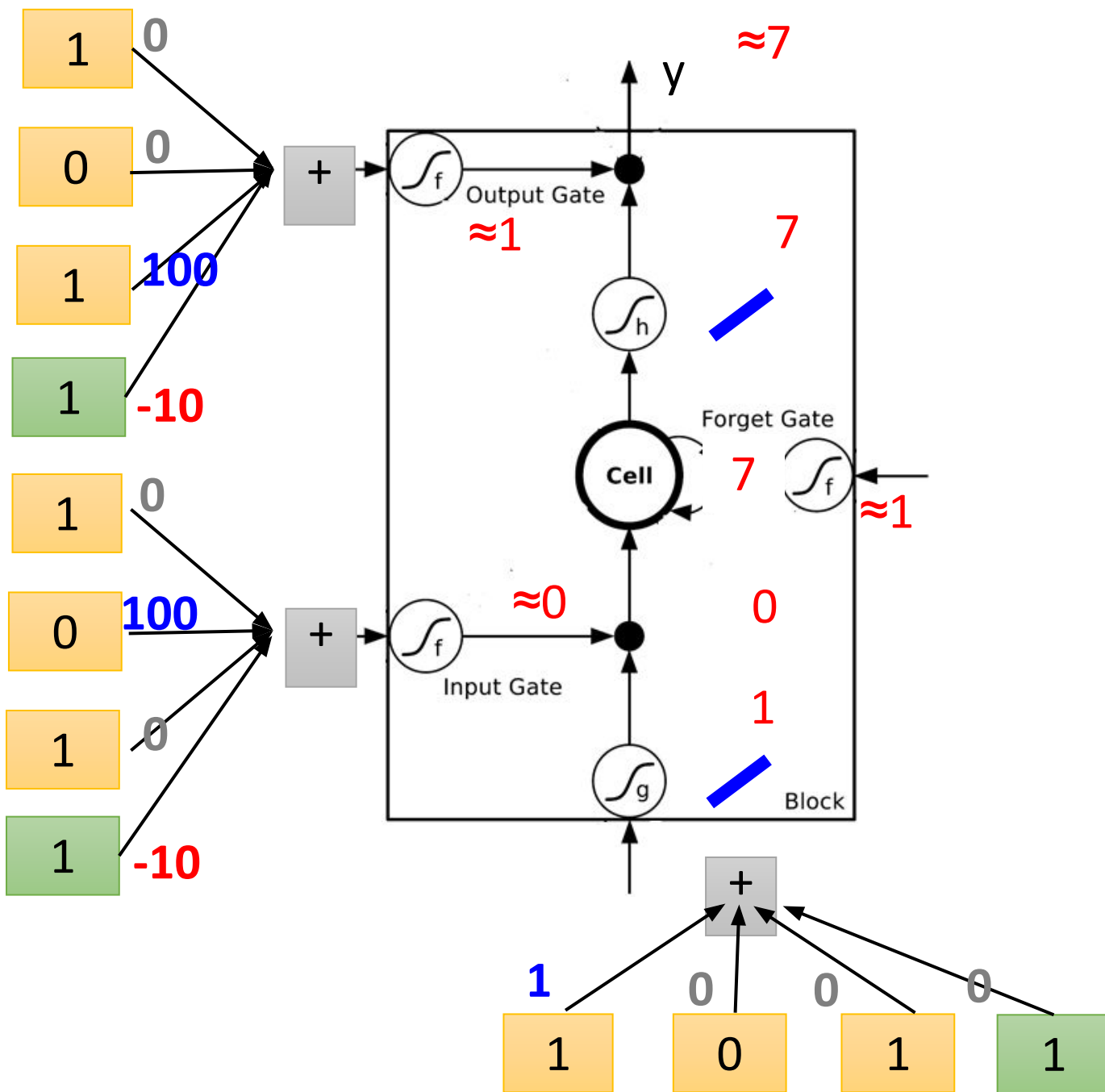


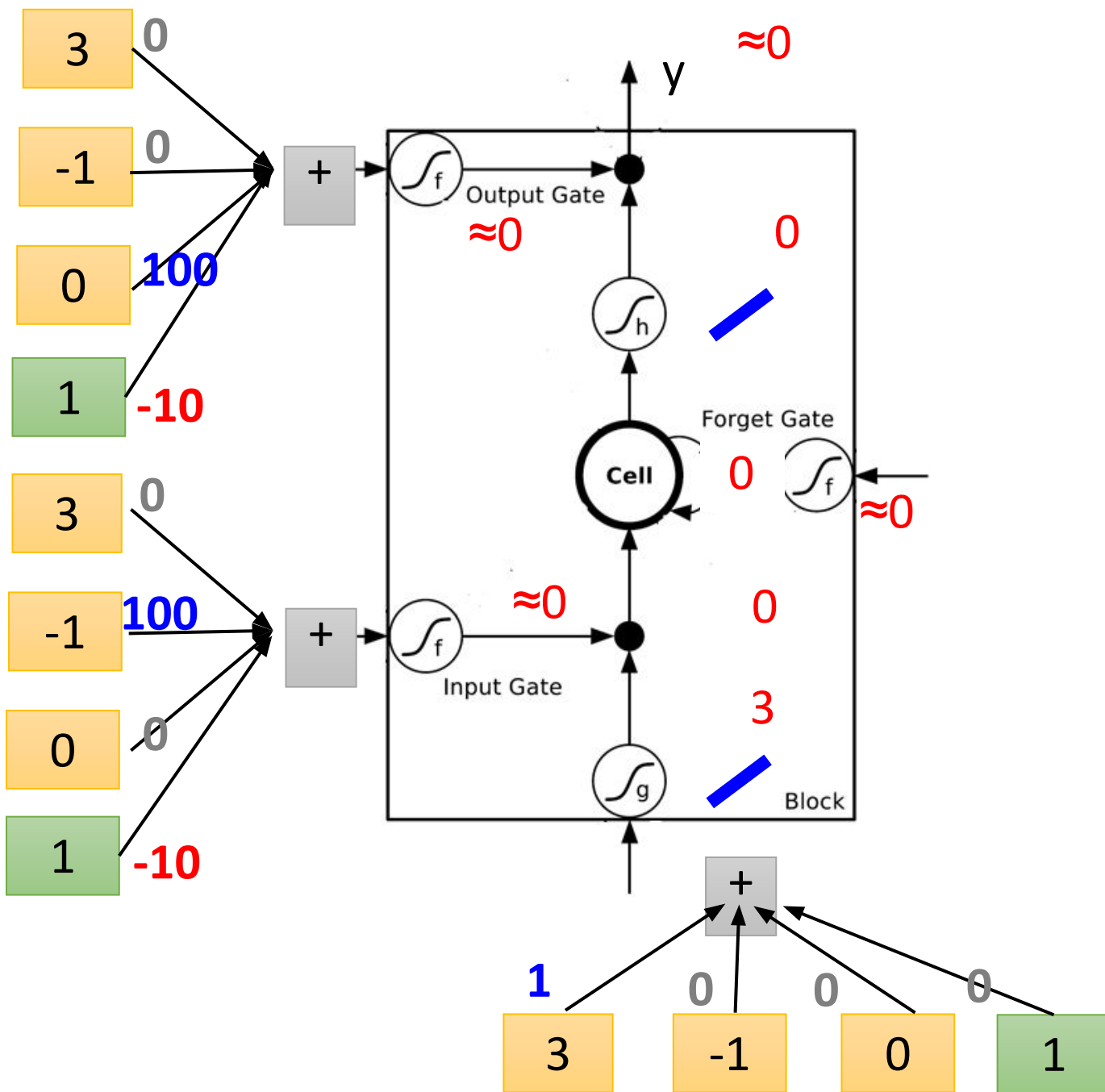


y 0 0 0 7 0

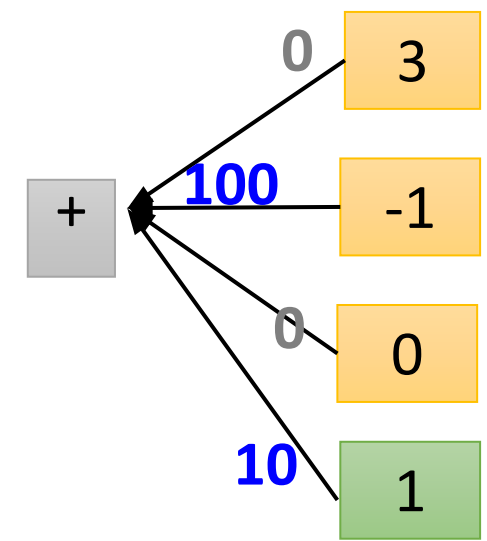


	x_1	x_2	x_3
	3	4	2
	1	1	0
	0	0	0





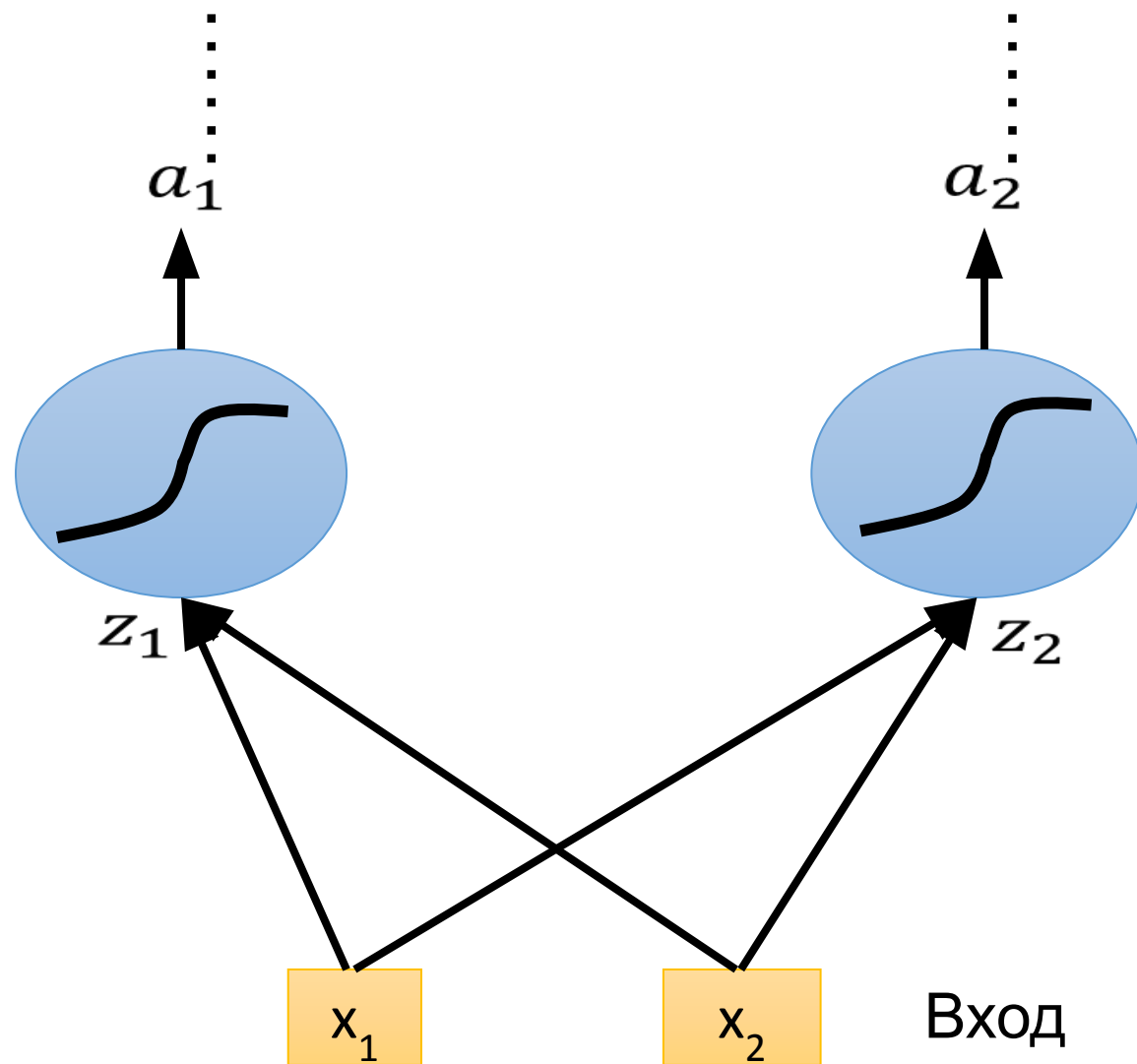
y 0 0 0 7 0

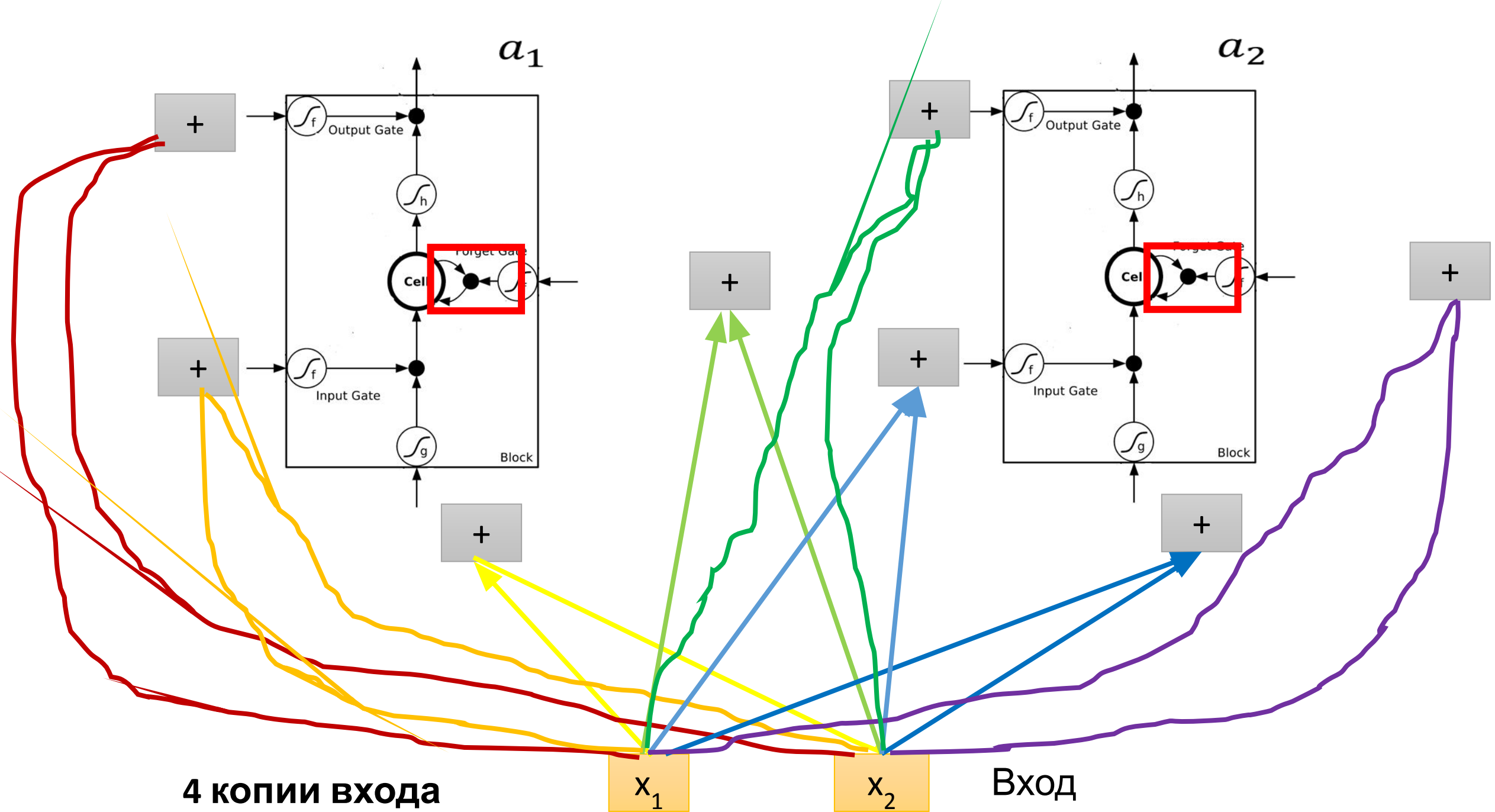


x_1	3	4	2	1	3
x_2	1	1	0	0	-1
x_3	0	0	0	1	0

Исходная сеть:

□ Заменяем обычный нейрон на LSTM

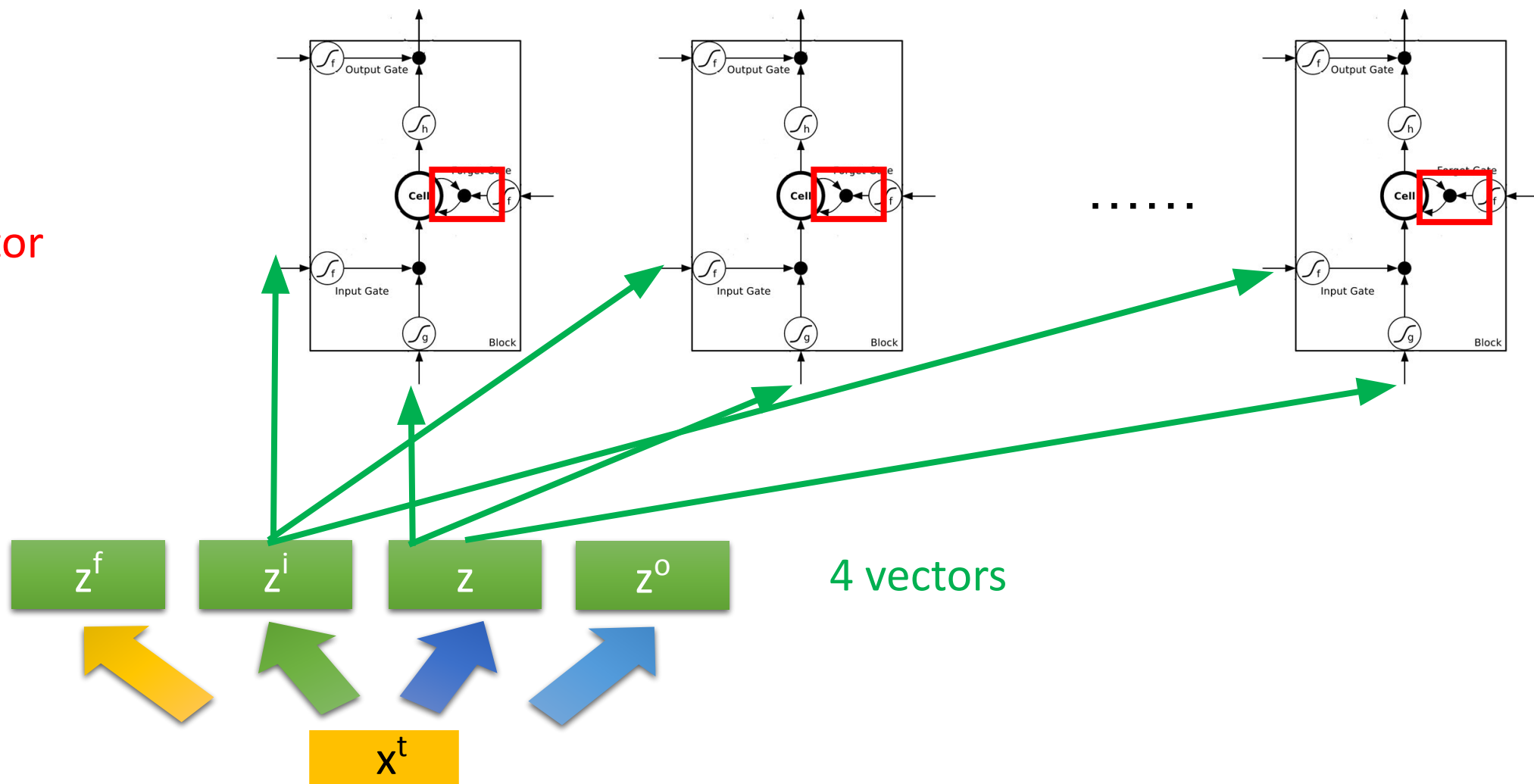




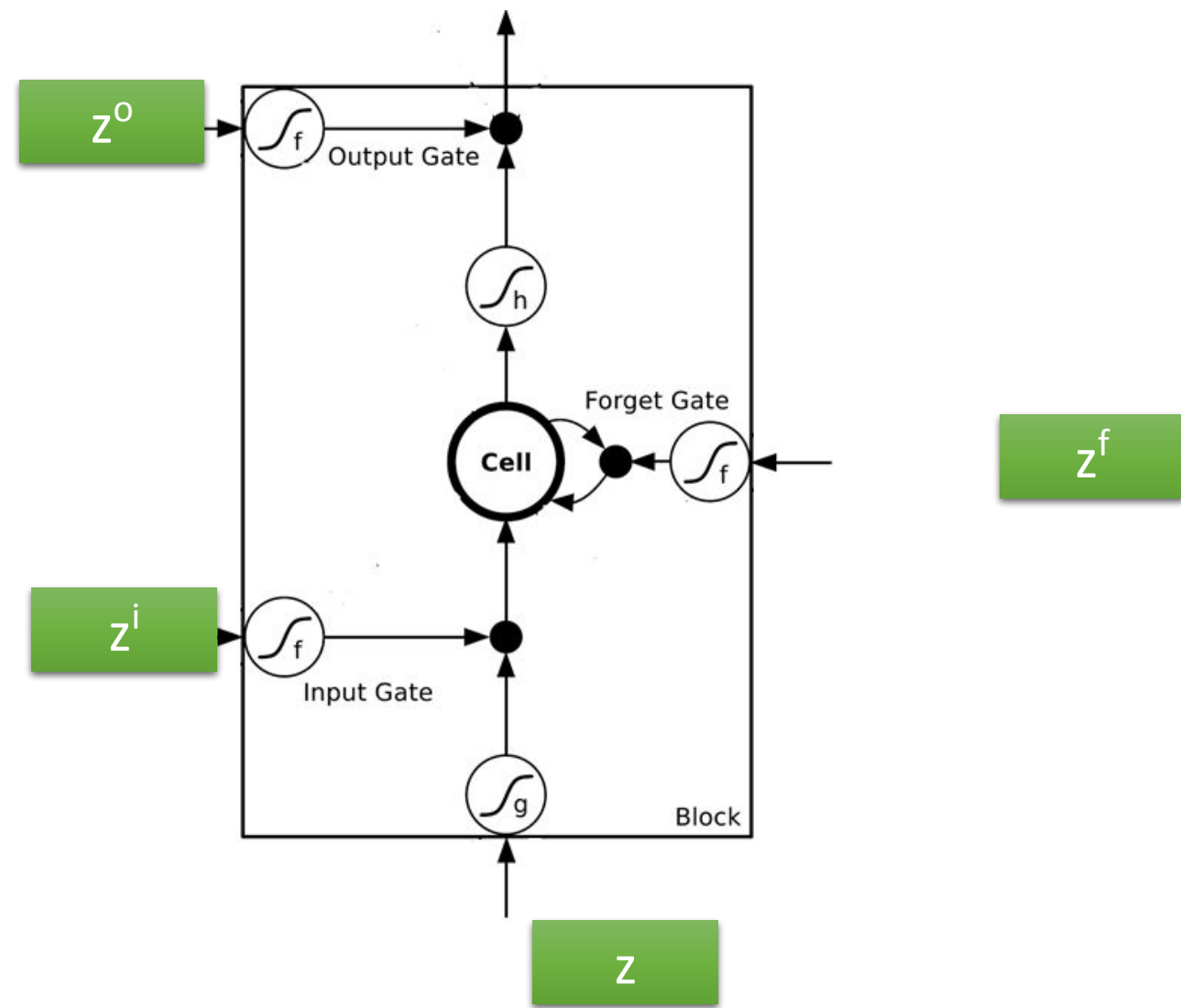
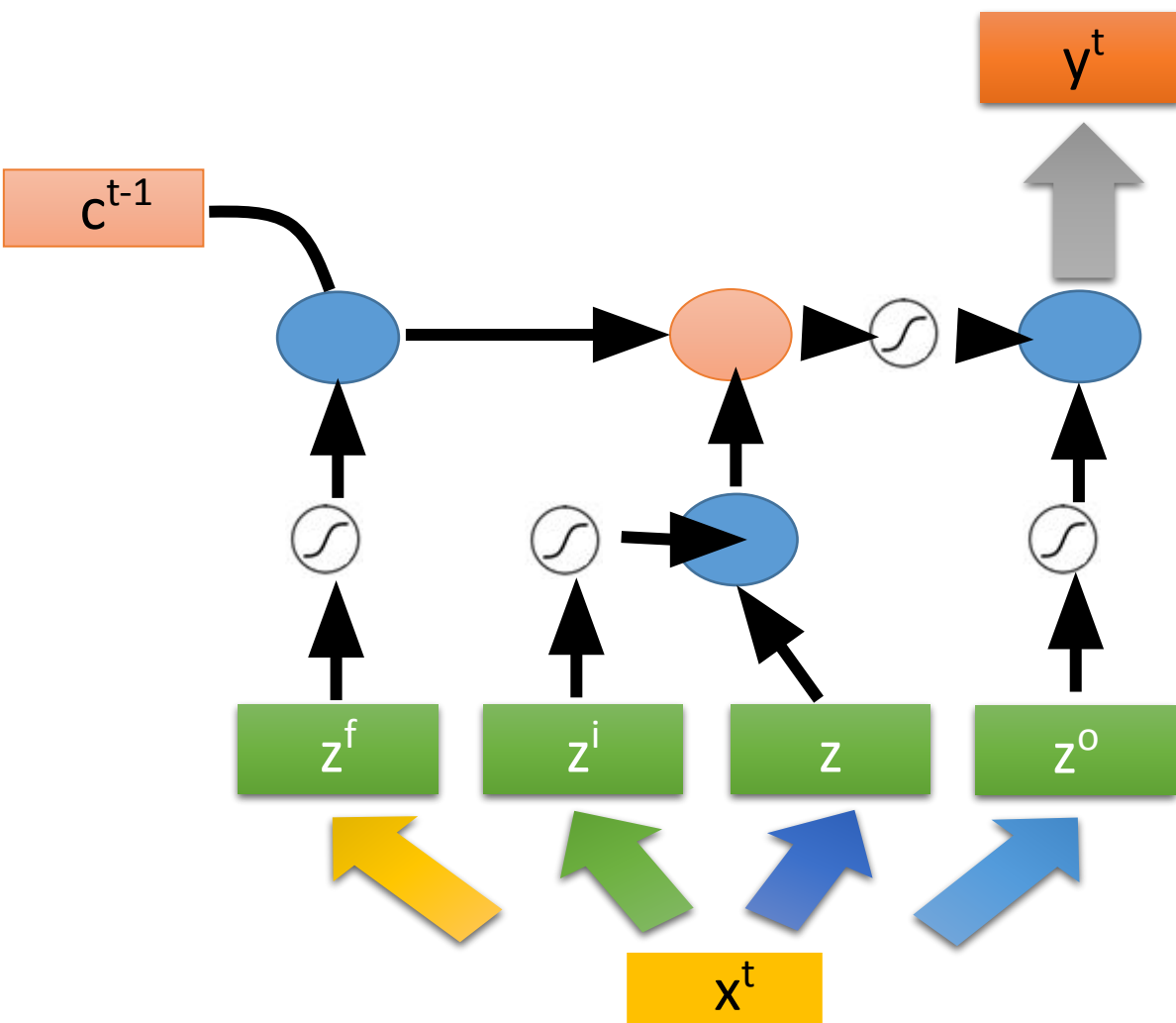
LSTM

 C^{t-1}

vector

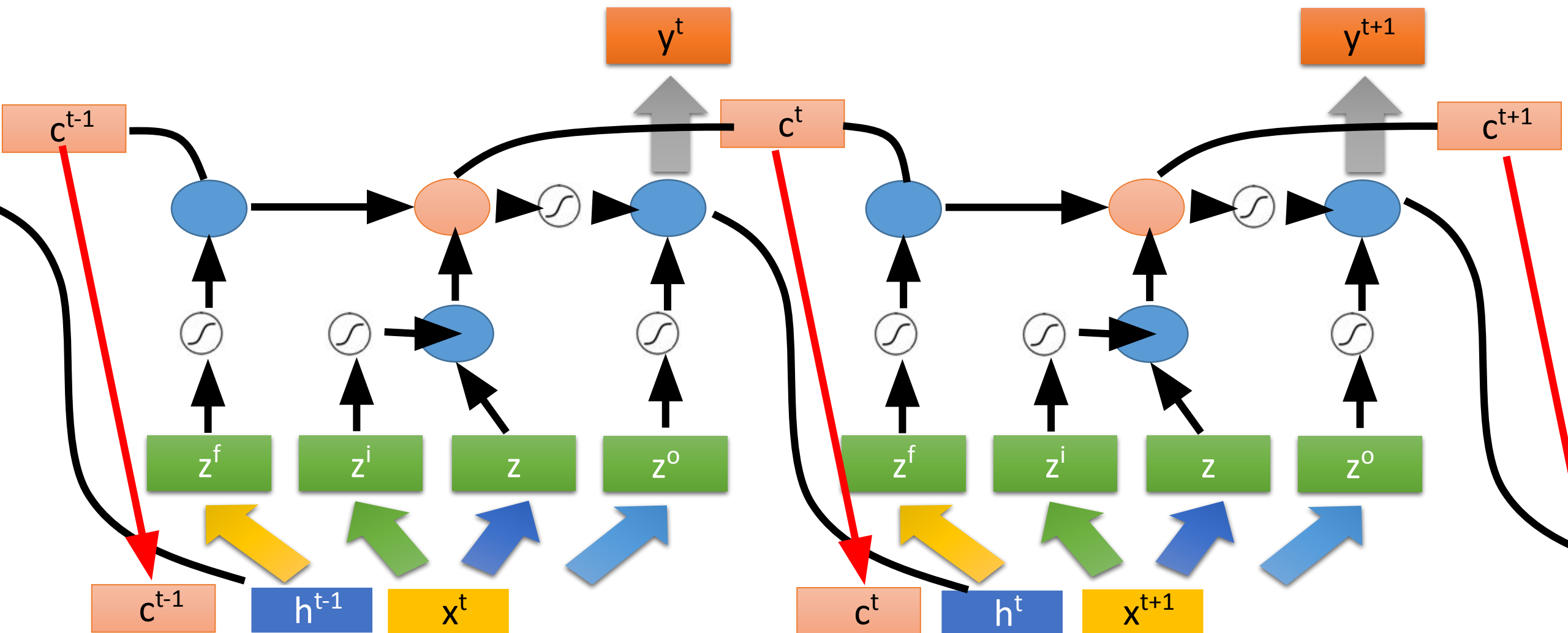


LSTM

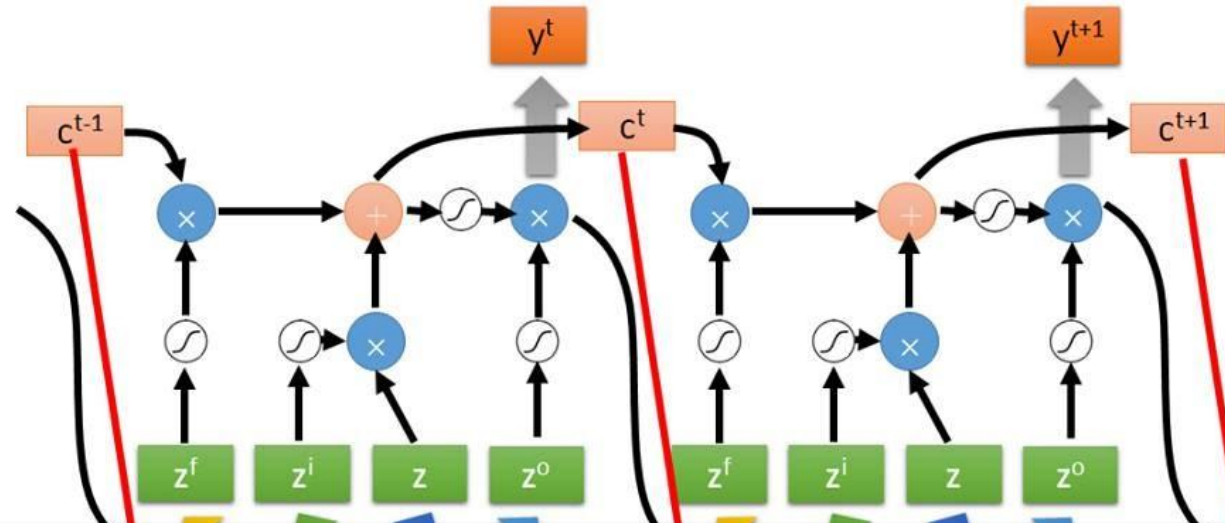


LSTM

Extension: "peephole"

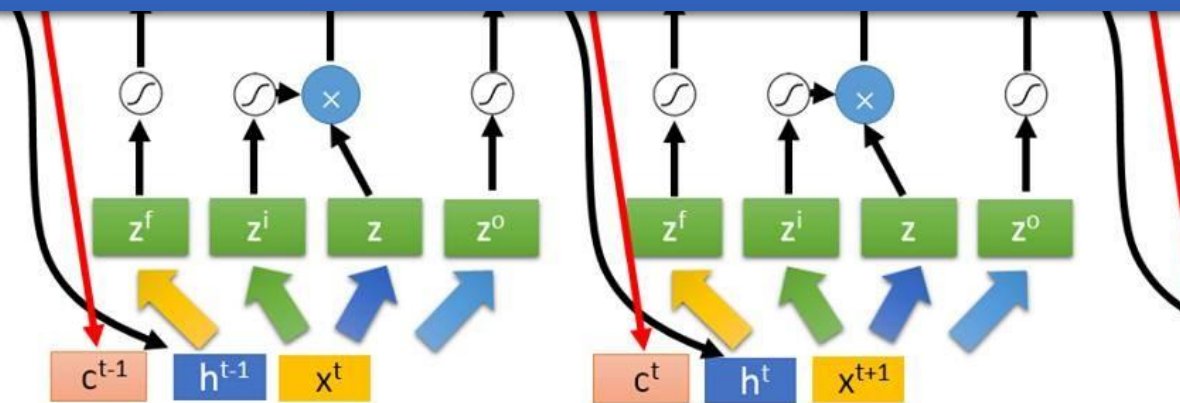


Multiple-layer LSTM



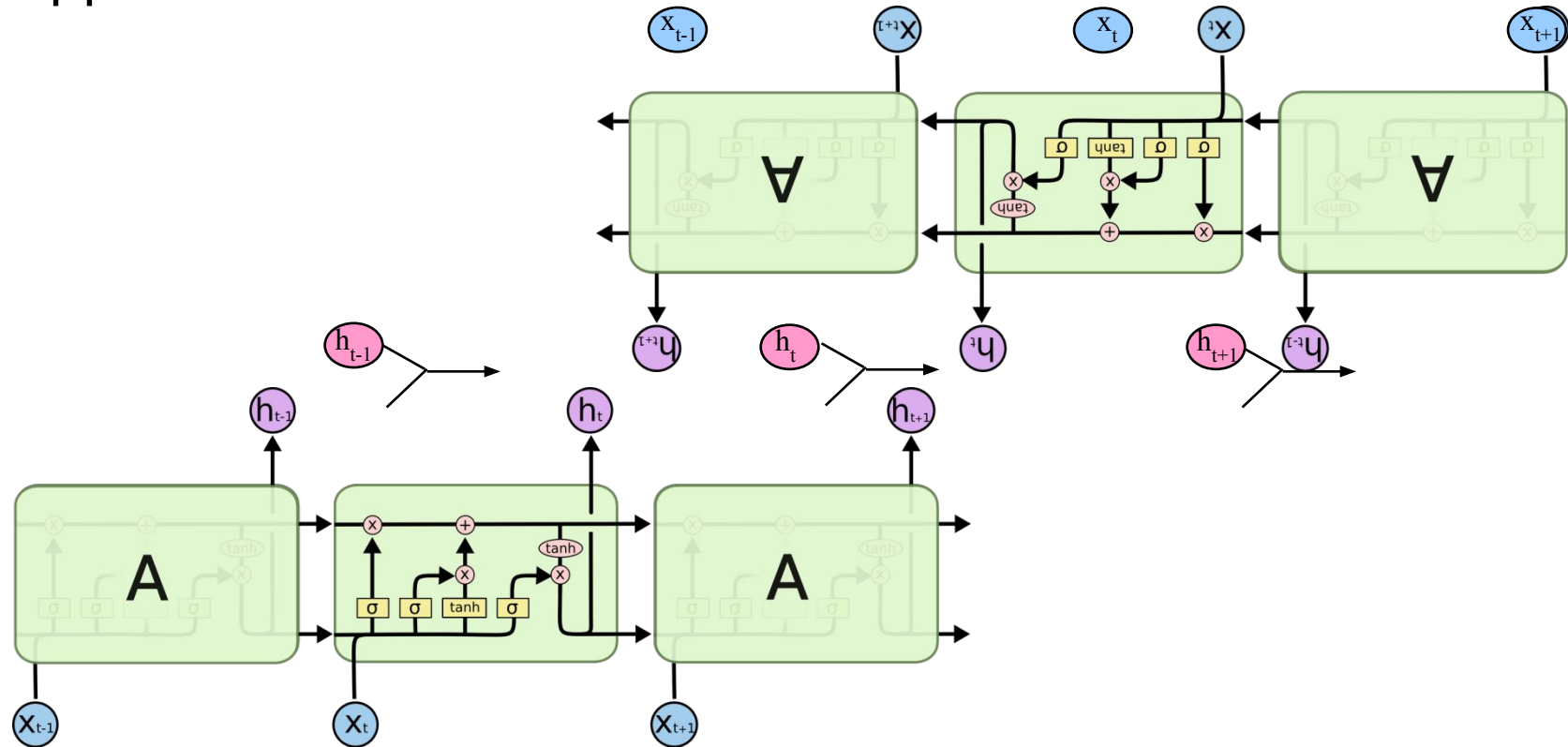
Не волнуйтесь, Keras позаботится об этом

Keras поддерживает
“LSTM”, “GRU”, “SimpleRNN” слои

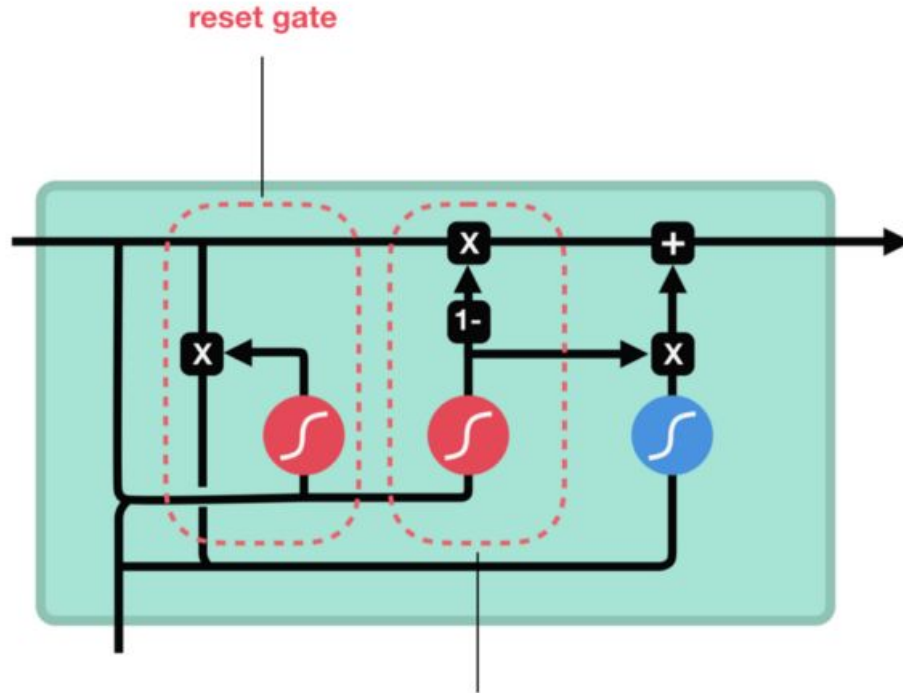


Bi-directional LSTM (Bi-LSTM)

- Отдельные LSTM обрабатывают последовательность вперед и назад, а также скрытые слои на каждом временном шаге объединяются для формирования выходных данных ячейки ».



GRU



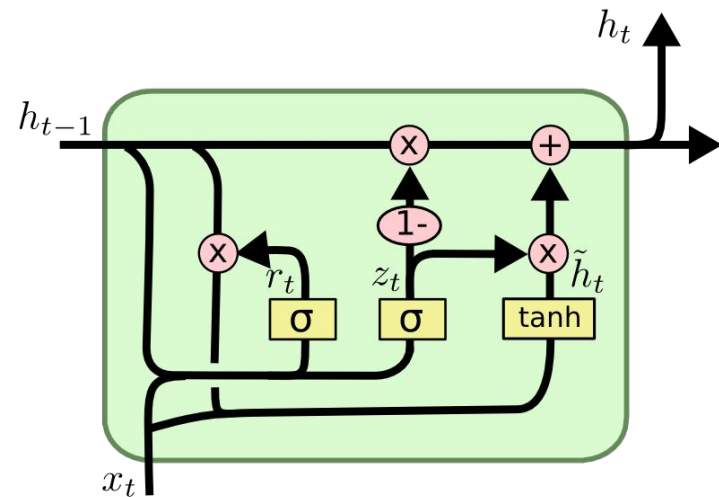
Update Gate

Шлюз обновления действует аналогично шлюзу забвения и входа LSTM. Он решает, какую информацию выбросить и какую новую добавить.

Reset Gate

Шлюз сброса - это еще один шлюз, который используется, чтобы решить, сколько прошлой информации нужно забыть.

В GRU меньше тензорных операций; поэтому они обучаются немного быстрее, чем LSTM.



update gate

$$z_t = \sigma(W_z \cdot [h_{t-1}, x_t])$$

$$r_t = \sigma(W_r \cdot [h_{t-1}, x_t])$$

$$\tilde{h}_t = \tanh(W \cdot [r_t * h_{t-1}, x_t])$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

Неизвестно, какая из архитектур лучше. Исследователи и инженеры обычно пытаются определить, этот момент экспериментальным путем..

Варианты применения LSTM

one to many

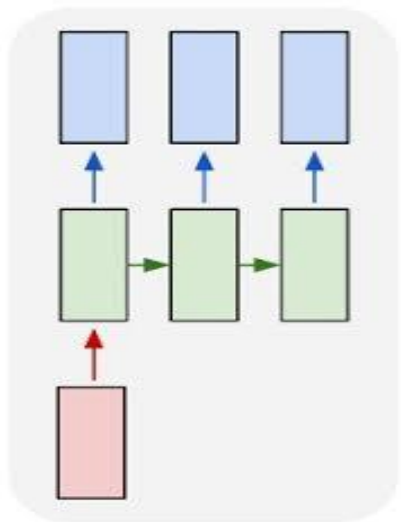
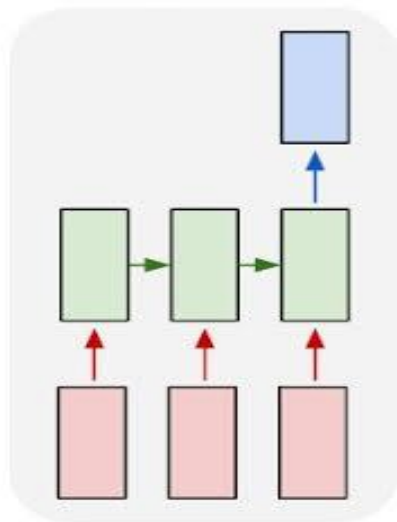


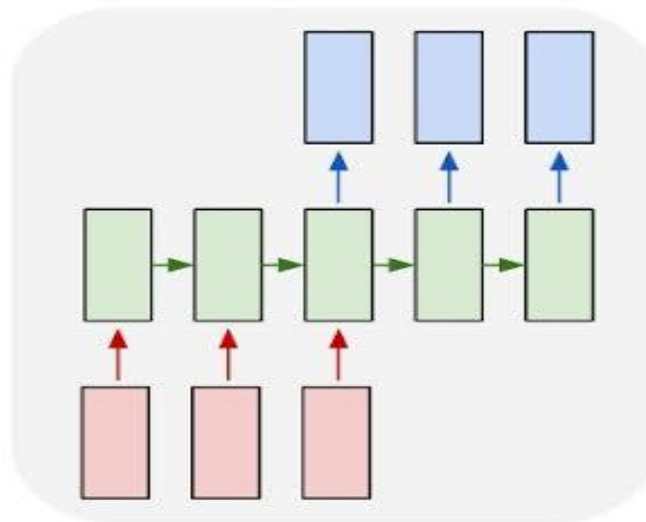
Image Captioning

many to one



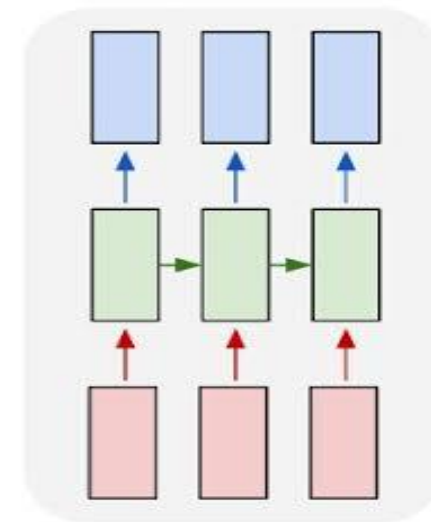
Video Activity Recog
Text Classification

many to many



Video Captioning
Machine Translation

many to many



POS Tagging
Language Modeling