



Faculty  
of  
**Computer**  
science  
Higher School of Economics

# Mamba: Linear-Time Sequence Modeling with Selective State Spaces

Подготовил:  
Казадаев Максим, БПМИ202

# План



- **SSM (State Space Model)** – повторение
- **Selective State Space Model** – улучшение модели
- **Mamba** – архитектура и интерпретация
- Эксперименты
- Вывод: **когда применять?**

# SSM (State Space Model)

Рассмотрим процесс, заданный диффузом:

$$\begin{aligned}h'(t) &= Ah(t) + Bx(t) \\ y(t) &= Ch(t)\end{aligned}$$

$x(t)$  — вход  
 $h(t)$  — скрытое состояние  
 $y(t)$  — выход  
 $A, B, C$  — параметры

Дискретизация диффура с шагом delta:  $(\Delta, A, B, C) \mapsto (\bar{A}, \bar{B}, C)$

$$\begin{aligned}h_t &= \bar{A}h_{t-1} + \bar{B}x_t \\ y_t &= Ch_t\end{aligned}$$

$\bar{A} = \exp(\Delta A)$   
 $\bar{B} = (\Delta A)^{-1}(\exp(\Delta A) - I) \cdot \Delta B$

По сути получили RNN. Но его можно вычислять с помощью свертки на GPU.

# Selective State Space Model

Увеличиваем число параметров. Добавим зависимость матриц **B**, **C** и шага **delta** от **x** (на картинке выделено **красным**). Авторы используют 1 линейный слой в качестве такой функции:

$$s_B(x) = \text{Linear}_N(x), s_C(x) = \text{Linear}_N(x), s_\Delta(x) = \text{Broadcast}_D(\text{Linear}_1(x)), \text{ and } \tau_\Delta = \text{softplus}$$

---

## Algorithm 1 SSM (S4)

---

**Input:**  $x : (B, L, D)$

**Output:**  $y : (B, L, D)$

1:  $\underline{A} : (D, N) \leftarrow \text{Parameter}$

▷ Represents structured  $N \times N$  matrix

2:  $\underline{B} : (D, N) \leftarrow \text{Parameter}$

3:  $\underline{C} : (D, N) \leftarrow \text{Parameter}$

4:  $\underline{\Delta} : (D) \leftarrow \tau_\Delta(\text{Parameter})$

5:  $\underline{\bar{A}}, \underline{\bar{B}} : (D, N) \leftarrow \text{discretize}(\underline{\Delta}, \underline{A}, \underline{B})$

6:  $y \leftarrow \text{SSM}(\underline{\bar{A}}, \underline{\bar{B}}, \underline{C})(x)$

▷ Time-invariant: recurrence or convolution

7: **return**  $y$

---

---

## Algorithm 2 SSM + Selection (S6)

---

**Input:**  $x : (B, L, D)$

**Output:**  $y : (B, L, D)$

1:  $\underline{A} : (D, N) \leftarrow \text{Parameter}$

▷ Represents structured  $N \times N$  matrix

2:  $\underline{B} : (B, L, N) \leftarrow s_B(x)$

3:  $\underline{C} : (B, L, N) \leftarrow s_C(x)$

4:  $\underline{\Delta} : (B, L, D) \leftarrow \tau_\Delta(\text{Parameter} + s_\Delta(x))$

5:  $\underline{\bar{A}}, \underline{\bar{B}} : (B, L, D, N) \leftarrow \text{discretize}(\underline{\Delta}, \underline{A}, \underline{B})$

6:  $y \leftarrow \text{SSM}(\underline{\bar{A}}, \underline{\bar{B}}, \underline{C})(x)$

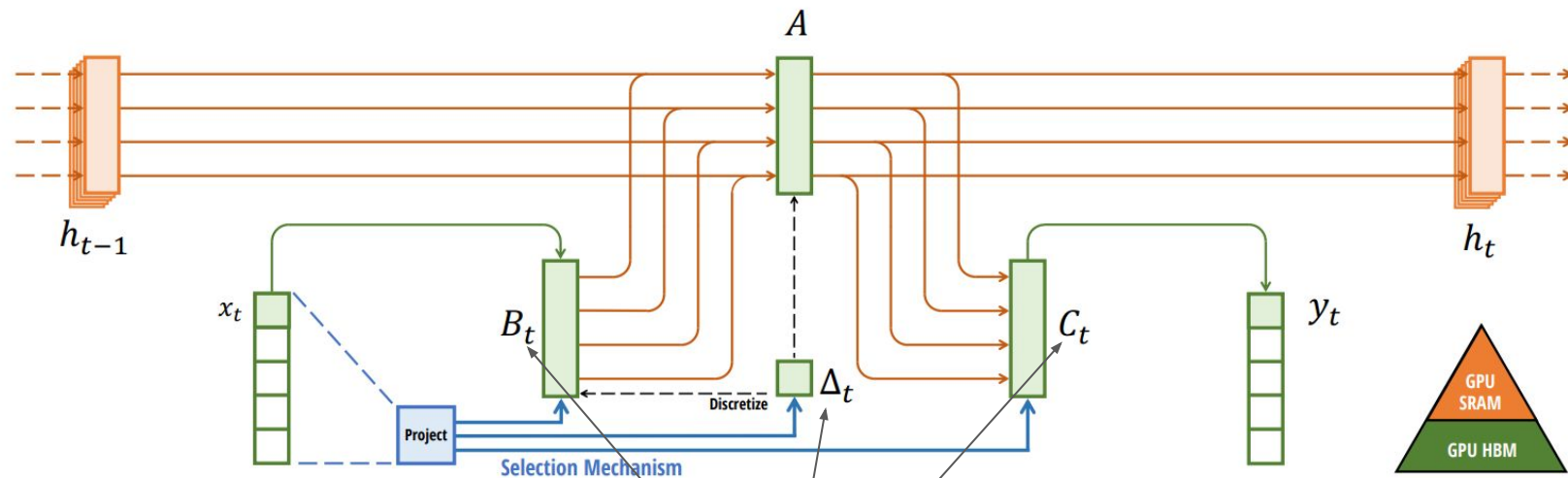
▷ Time-varying: recurrence (*scan*) only

7: **return**  $y$

---

# Selective State Space Model

## Selective State Space Model with Hardware-aware State Expansion



Зависимость параметров от входа  
(Selection Mechanism)

Параметры SSM

# Selective State Space Model

Процесс, заданный диффузом:

$$\begin{aligned}h'(t) &= Ah(t) + Bx(t) \\ y(t) &= Ch(t)\end{aligned}$$

$x(t)$  — ВХОД  
 $h(t)$  — скрытое состояние  
 $y(t)$  — ВЫХОД

Дискретизация диффура с шагом delta:  $(\Delta, A, B, C) \mapsto (\bar{A}, \bar{B}, C)$

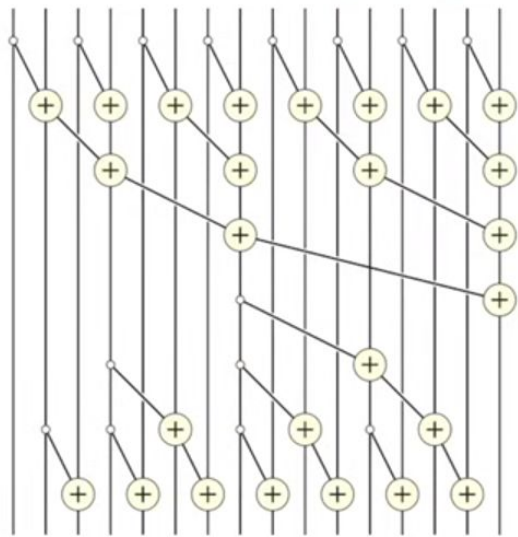
**B, C, delta** — параметры, зависящие от **x**

$$\begin{aligned}h_t &= \bar{A}h_{t-1} + \bar{B}x_t \\ y_t &= Ch_t\end{aligned}$$
$$\begin{aligned}\bar{A} &= \exp(\Delta A) \\ \bar{B} &= (\Delta A)^{-1}(\exp(\Delta A) - I) \cdot \Delta B\end{aligned}$$

# Selective State Space Model

Теперь на каждом объекте матрицы A, B, C разные, поэтому нельзя использовать свертку. Авторы придумали другой алгоритм, который называли **scan**.

9	6	7	10	8	7
---	---	---	----	---	---



9	15	22	32	40	47
---	----	----	----	----	----

$$y_k = C\bar{A}^k\bar{B}x_0 + C\bar{A}^{k-1}\bar{B}x_1 + \dots + C\bar{A}\bar{B}x_{k-1} + C\bar{B}x_k$$

Ассоциативность операции умножения:

$$(X * Y) * Z = X * (Y * Z)$$

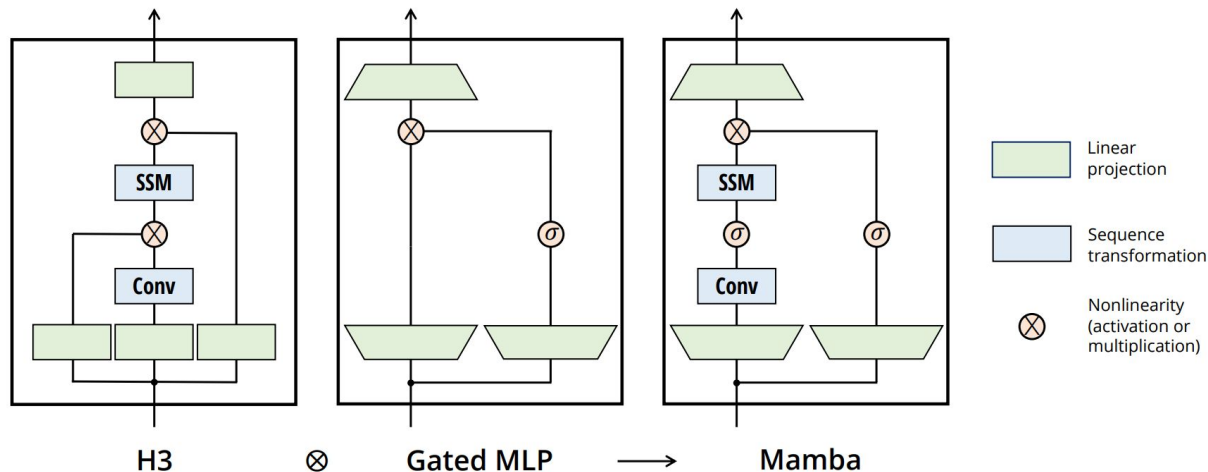
$$O(N) \rightarrow O(N / (\text{число тредов}))$$

# Mamba

Mamba (3 картинка) – это комбинация SSM и свертки

В Mamba вход проходит в следующем порядке через слои:

- Линейные слои (зеленые).
- **Свертка (Conv) и SSM.**
- **Skip-connections**
- Выход SSM и **skip-connections** перемножаются



Авторы скомбинировали 2 существующие модели (H3 и Gated MLP) и получили Mamba



# Эксперименты

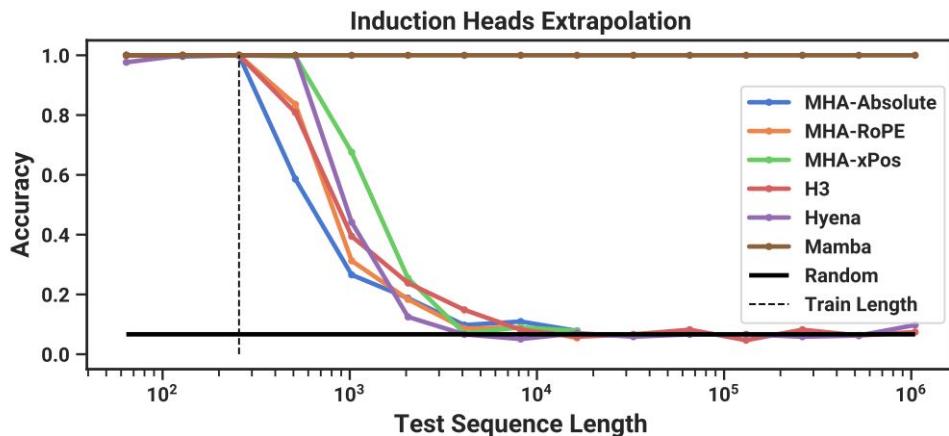


Table 2: (**Induction Heads.**) Models are trained on sequence length  $2^8 = 256$ , and tested on increasing sequence lengths of  $2^6 = 64$  up to  $2^{20} = 1048576$ . Full numbers in Table 11.

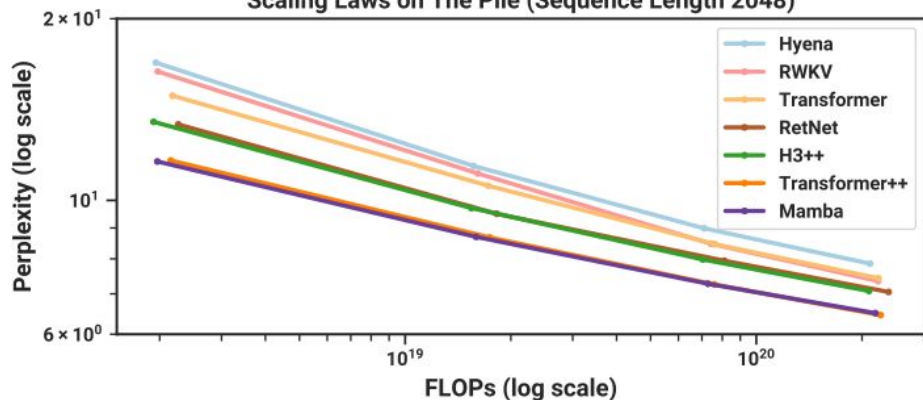
С ростом длины входа качество Mamba не падает.  
Mamba может работать с очень длинными текстами без потери качества и скорости.

**Задача:** Если модель ранее видела биграмму в последовательности, например, “Harry Potter”, то в следующий раз, когда “Harry” появится в этой же последовательности, модель должна продолжить “Potter”

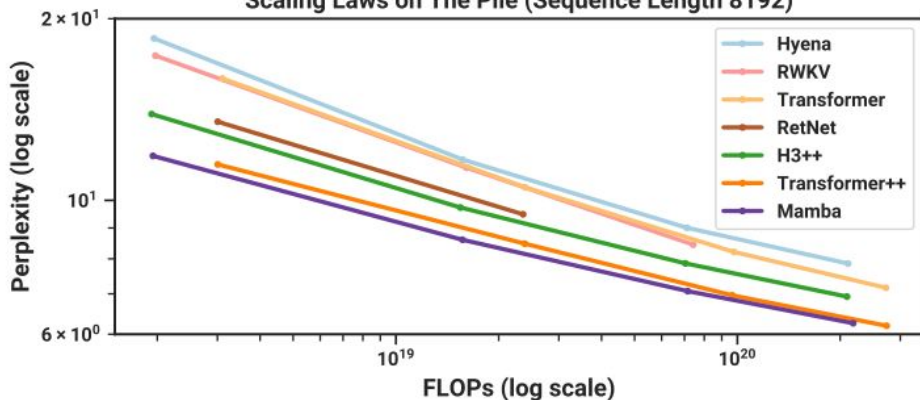
# Число операций vs Качество



Scaling Laws on The Pile (Sequence Length 2048)



Scaling Laws on The Pile (Sequence Length 8192)

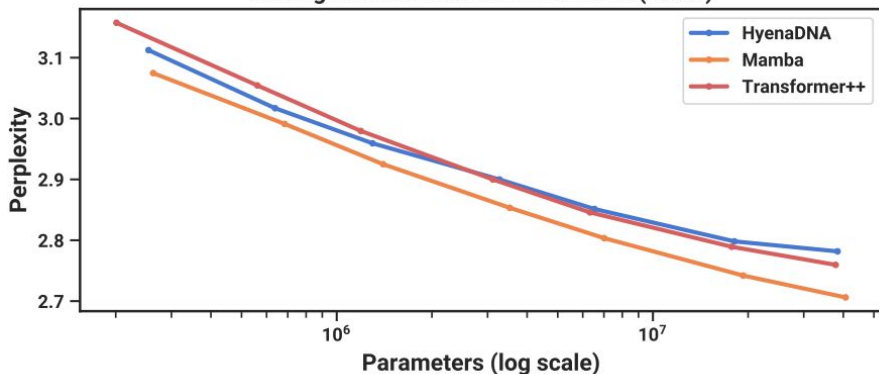


*Задача предсказания следующего слова.*

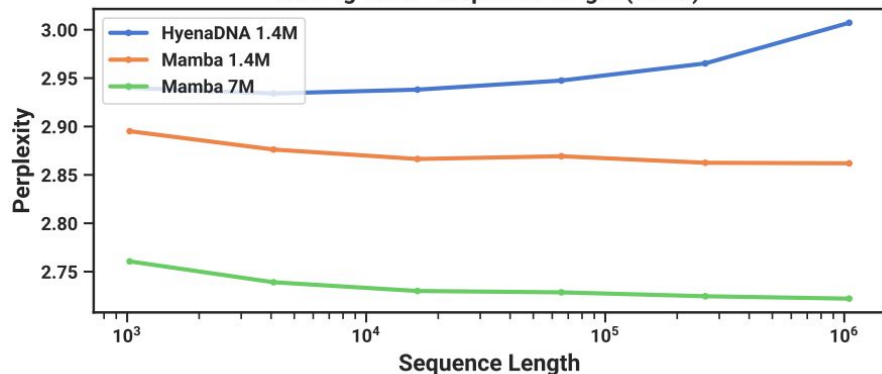
- Эксперимент сравнивает **качество одинаково сложных** моделей
- Mamba при том же числе операций FLOPs-ов добивается более хорошего качества, чем другие архитектуры.
- Качество Transformer++ и Mamba **почти одинаковое**

# DNA

Scaling Laws on the Human Genome (HG38)



Scaling Laws - Sequence Length (HG38)



**Задача:** предсказание следующего элемента в последовательности генома человека. DNA – это **длинные последовательности**, для них Mamba справляется лучше трансформеров.

Практическое применение Mamba – **работа с последовательностями, которые длиннее 10 000 слов**. Здесь трансформеры начинают уступать. Для всех остальных задач стоит использовать трансформеры.

# Плюсы и минусы Mamba

Плюсы:

- Оптимизирована для более **быстрого вычисления на GPU**
- Эффективно на **длинных последовательностях:  $O(N)$**
- Первая **модель без attention**, которая приближается к трансформерам

Минусы:

- Авторы не экспериментировали с очень **большими моделями** трансформеров
- На задачах с небольшими последовательностями Mamba **не показывает прироста в качестве**

## Практический вывод

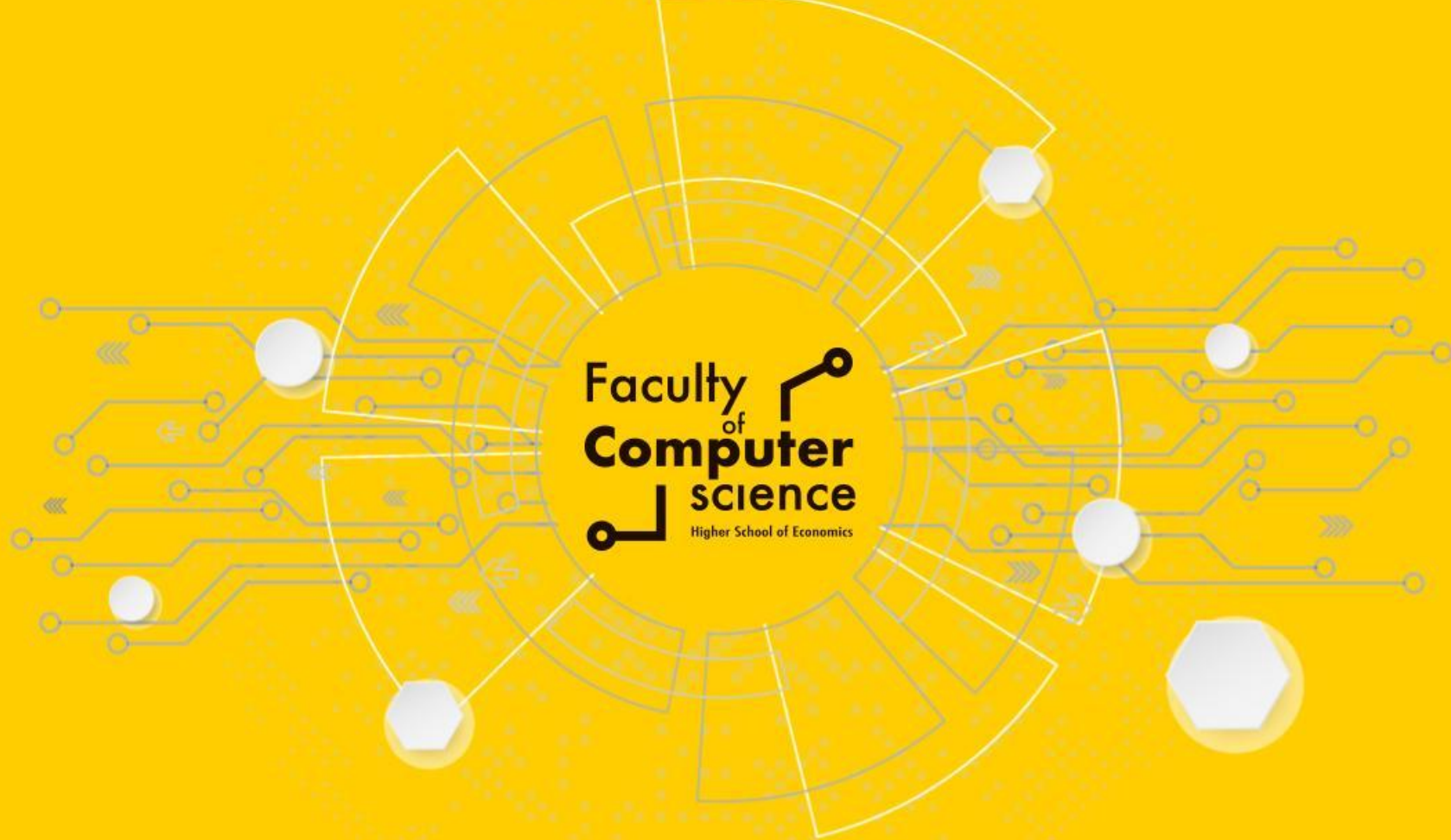
Последовательность длиннее 10 000 токенов?

- Да => Mamba
- Нет => Transformer

# Ссылки



- Оригинальная статья (сложна для понимания):  
<https://arxiv.org/ftp/arxiv/papers/2312/2312.00752.pdf>
- Гайд по Mamba: [https://www.youtube.com/watch?v=8Q\\_tqwpTpVU](https://www.youtube.com/watch?v=8Q_tqwpTpVU)
- Более простой гайд по Mamba: <https://www.youtube.com/watch?v=9dSkvxS2EB0>



[mskazadaev@edu.hse.ru](mailto:mskazadaev@edu.hse.ru)

# SSM (State Space Model)

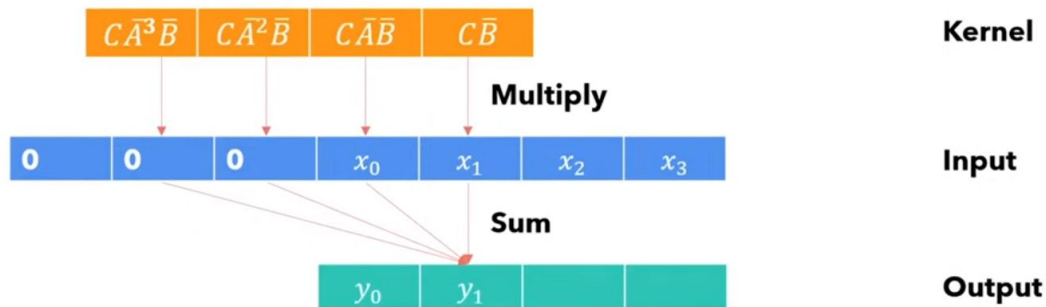
По сути получили RNN. Но его можно вычислять с помощью свертки на GPU.

$$\begin{aligned} h_t &= \bar{A}h_{t-1} + \bar{B}x_t \quad \rightarrow \quad y_k = C\bar{A}^k\bar{B}x_0 + C\bar{A}^{k-1}\bar{B}x_1 + \dots + C\bar{A}\bar{B}x_{k-1} + C\bar{B}x_k \\ y_t &= Ch_t \end{aligned}$$

Одинаковые матрицы для каждого объекта в выборке, поэтому их считаем один раз для батча

$$\bar{K} = (C\bar{B}, C\bar{A}\bar{B}, \dots, C\bar{A}^k\bar{B}, \dots)$$

$$y = x * \bar{K}$$



Интерпретация архитектуры:

1. **Свертка.** Нужна для того, чтобы перемешать входы. **Решает проблему RNN (слабый эффект далеких входов)**, т. к. теперь каждый input будет состоять не только из  $n$ -го входа, но из остальных. Каждый вход может поучаствовать на каждом шаге RNN.
2. **Selective SSM** – по сути улучшенная RNN для учета временных зависимостей. В RNN иногда при увеличении длины контекста качество наоборот падает. В SSM матрицы зависят от входа, что позволяет модели эффективно фильтровать шум, обнуляя похожие на шум входы, например, устремляя  $\delta$  к infinity. Именно поэтому модель называется **Selective**.