It's Raw! Audio Generation with State-Space Models

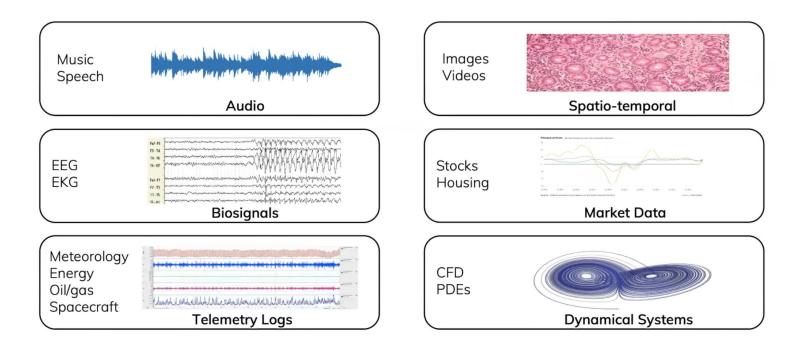
Докладчик: Спирин Иван

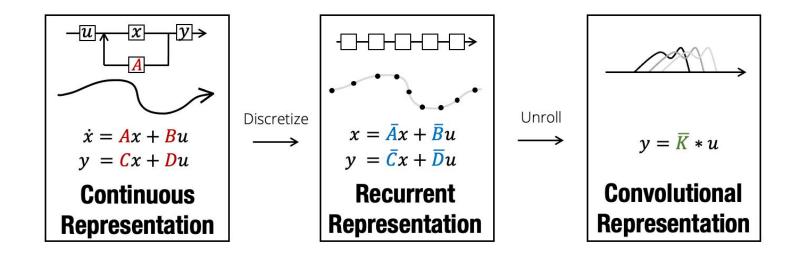
Рецензент-исследователь: Ткаченко Егор

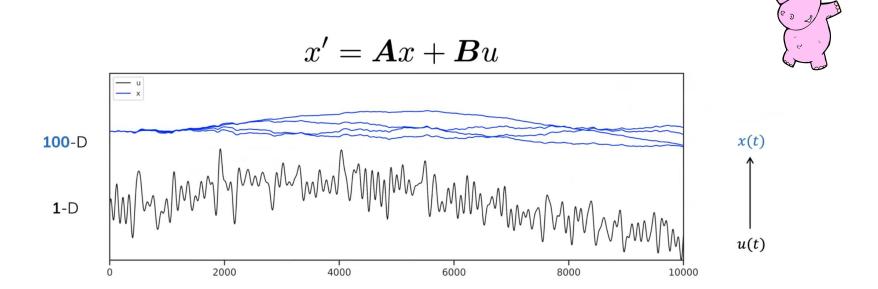
Хакер: Виноградова Ульяна

НИС МОП, 09.11.2022

Задача: моделирование длинных последовательностей

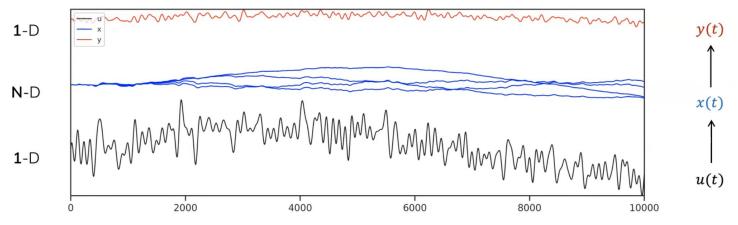






Оператор HiPPO отображает входной сигнал \boldsymbol{u} в сигнал более высокой размерности \boldsymbol{x}

$$x' = \mathbf{A}x + \mathbf{B}u$$
$$y = \mathbf{C}x + \mathbf{D}u$$



SSM отображает x в одномерный выход y

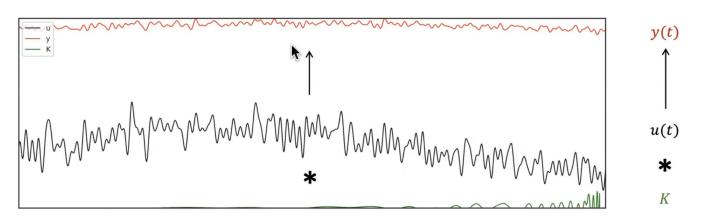
Дискретизация:

$$egin{align} \overline{m{A}} &= (m{I} - \Delta/2 \cdot m{A})^{-1} (m{I} + \Delta/2 \cdot m{A}) \ \overline{m{B}} &= (m{I} - \Delta/2 \cdot m{A})^{-1} \Delta m{B} \ \overline{m{C}} &= m{C} \ \end{aligned}$$

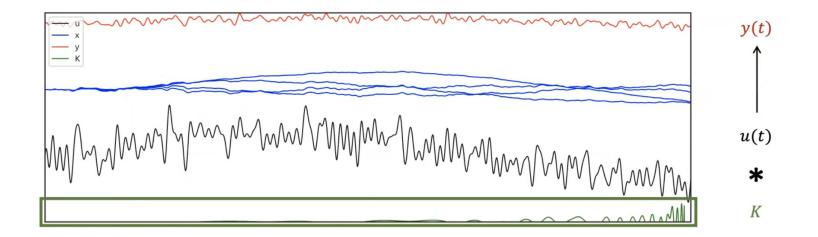
Итоговые формулы:

$$egin{aligned} x_k &= \overline{oldsymbol{A}} x_{k-1} + \overline{oldsymbol{B}} u_k \ y_k &= \overline{oldsymbol{C}} x_k \end{aligned}$$

$$y = u * K$$

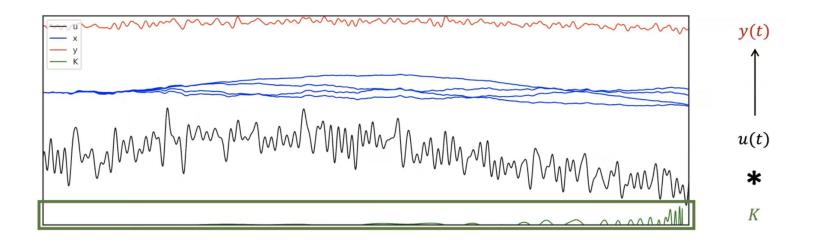


Можем смотреть на модель как на CNN с ядром К



Сложность: $O(N^2L)$

Память: O(NL)



Сложность: $O(N^2L) o O(N+L)$

Память: O(NL) o O(N+L)

Параметризация:

$$A = \Lambda - PQ^*$$

Algorithm 1 S4 Convolution Kernel (Sketch)

Input: S4 parameters $\Lambda, P, Q, B, C \in \mathbb{C}^N$ and step size Δ Output: SSM convolution kernel $\overline{K} = \mathcal{K}_L(\overline{A}, \overline{B}, \overline{C})$ for $A = \Lambda - PQ^*$ (equation (5))

1: $\widetilde{C} \leftarrow \left(I - \overline{A}^L\right)^* \overline{C}$ \triangleright Truncate SSM generating function (SSMGF) to length L2: $\begin{bmatrix} k_{00}(\omega) & k_{01}(\omega) \\ k_{10}(\omega) & k_{11}(\omega) \end{bmatrix} \leftarrow \left[\widetilde{C} \ Q\right]^* \left(\frac{2}{\Delta} \frac{1-\omega}{1+\omega} - \Lambda\right)^{-1} [B \ P]$ \triangleright Black-box Cauchy kernel

3: $\hat{K}(\omega) \leftarrow \frac{2}{1+\omega} \left[k_{00}(\omega) - k_{01}(\omega)(1+k_{11}(\omega))^{-1}k_{10}(\omega)\right]$ \triangleright Woodbury Identity

4: $\hat{K} = \{\hat{K}(\omega) : \omega = \exp(2\pi i \frac{k}{L})\}$ \triangleright Evaluate SSMGF at all roots of unity $\omega \in \Omega_L$ 5: $\overline{K} \leftarrow \mathsf{iFFT}(\hat{K})$ \triangleright Inverse Fourier Transform

Результаты:

	TRAIN	NING ST	EP (MS)	MEMORY ALLOC. (MB)				
Dim.	128	256	512	128	256	512		
LSSL S4	9.32 4.77	20.6 3.07	140.7 4.75	222.1 5.3	1685 12.6	13140 33.5		
Ratio	1.9×	6.7×	29.6 ×	42.0×	133×	392×		

Параметризация улучшает эффективность: до 30 раз по времени, до 400 раз по памяти.

Результаты:

	Convolution ³	Recurrence	Attention	S4
Parameters	LH	H^2	H^2	H^2
Training	$ ilde{L}H(B+H)$	BLH^2	$B(L^2H + LH^2)$	$BH(ilde{H}+ ilde{L})+B ilde{L}H$
Space	BLH	BLH	$B(L^2 + HL)$	BLH
Parallel	Yes	No	Yes	Yes
Inference	LH^2	H^2	$L^2H + H^2L$	H^2

Вычислительная сложность различных моделей.

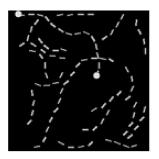
L - длина последовательности, B - размер батча, H - размер скрытого пространства.

Результаты:

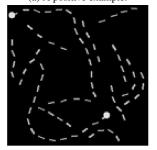
Model	LISTOPS	Техт	Retrieval	IMAGE	Pathfinder	Ратн-Х	Avg
Transformer	36.37	64.27	57.46	42.44	71.40	X	53.66
Reformer	37.27	56.10	53.40	38.07	68.50	X	50.56
BigBird	36.05	64.02	59.29	40.83	74.87	X	54.17
Linear Trans.	16.13	65.90	53.09	42.34	75.30	×	50.46
Performer	18.01	65.40	53.82	42.77	77.05	X	51.18
FNet	35.33	65.11	59.61	38.67	77.80	Х	54.42
Nyströmformer	37.15	65.52	79.56	41.58	70.94	X	57.46
Luna-256	37.25	64.57	79.29	47.38	77.72	×	59.37
S4	59.60	$\bf 86.82$	90.90	88.65	94.20	96.35	86.09

Long Range Arena: A Benchmark For Efficient Transformers

Примеры из Path-X

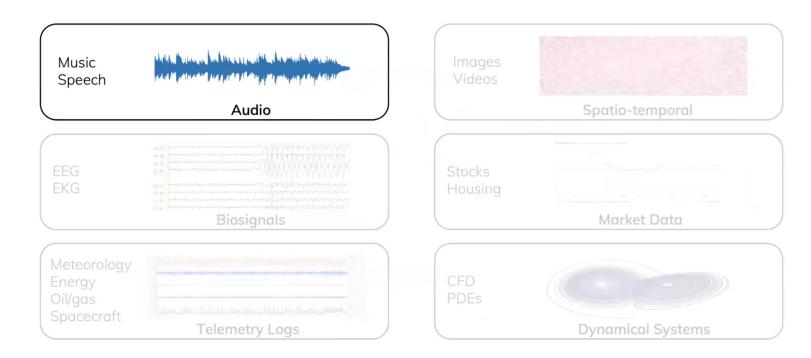


(a) A positive example.

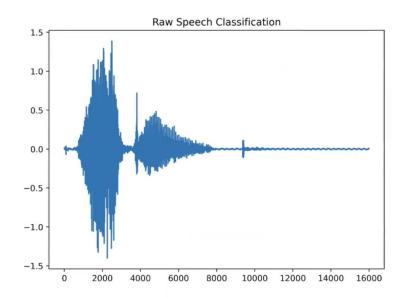


(b) A negative example.

Задачи аудио-моделирования



Задачи аудио-моделирования



1 секунда = 16000 значений

Длина:	161	16000	8000(test
	MFCC	Raw	$0.5 \times$
Transformer	90.75	x	X 30.68
Performer	80.85	30.77	
ODE-RNN	65.9	x	X
NRDE	89.8	16.49	15.12
ExpRNN	82.13	11.6	10.8
LipschitzRNN	88.38	x	X
CKConv	95.3	71.66 96.25	65.96
WaveGAN-D	X		X
LSSL S4	93.58	X	х
	<u>93.96</u>	98.32	96.30

Классификация SC10 с препроцессингом и без него

Задача аудио-генерации

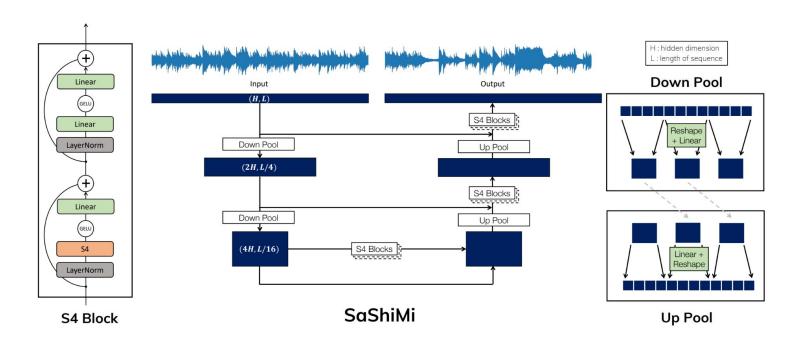
Два вида моделей: авторегрессионные (AR) - генерируют выходные значения последовательно - и неавторегрессионные (non-AR), которые выдают весь выход за раз.

non-AR:

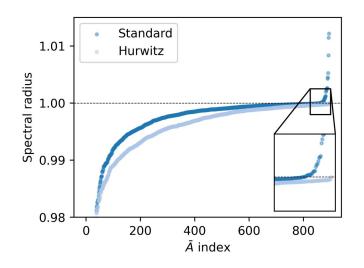
- Эффективнее на инференсе, лучше качество.
- DiffWave (Kong et al. 2021), WaveGAN (Donehue et al. 2019).

AR:

- Генерируют последовательность неограниченной длины.
- WaveNet (van den Oord et al. 2016), SampleRNN (Mehri et al. 2017).



- ullet Эксперименты показали, что параметризация $A = \Lambda PQ^*$ может вести к численной неустойчивости.
- Новая параметризация
- Теперь собственные значения строго ограничены сверху единицей по модулю.



LEARNED	FROZEN	NLL	STABLE GENERATION
$ \begin{array}{c} -\\ \Lambda + pq^*\\ \Lambda - pp^* \end{array} $	$\Lambda + pq^*$ $-$	1.445 1.420 1.419	

• non-AR модификация - двунаправленный слой S4:

$$y = Linear(Concat(S4(x), rev(S4(rev(x)))))$$

- Эксперименты проводим в двух задачах:
 - а. Генерация фортепианной музыки (Beethoven, YouTubeMix).
 - b. Генерация 1-секундных звуковых дорожек с цифрами 0-9 (SC09).

CATEGORY	Dataset	TOTAL DURATION	CHUNK LENGTH	Sampling Rate	QUANTIZATION	SPLITS (TRAIN-VAL-TEST)
Music	BEETHOVEN	10 hours	8s	16ĸHz	8-bit linear	Mehri et al. [28]
Music	YouTubeMix	4 Hours	8s	$16 \mathrm{KHz}$	8-BIT MU-LAW	88% - 6% - 6%
SPEECH	SC09	5.3 Hours	1s	16ĸHz	8-BIT MU-LAW	Warden [42]

Результаты:

Beethoven

Model	Context	NLL	$@200\mathrm{K}$ steps	@10 Hours
SAMPLERNN* WAVENET*	1024 4092	1.076 1.464	_	
SampleRNN [†] WaveNet [†]	1024 4092	1.125 1.032	1.125 1.088	1.125 1.352
SaShiMi	128000	0.946	1.007	1.095

^{*}Reported in Mehri et al. [28]

YouTubeMix

Model	Test NLL	MOS (FIDELITY)	MOS (MUSICALITY)
SAMPLERNN WAVENET SASHIMI	1.723 1.449 1.294	2.98 ± 0.08 2.91 ± 0.08 2.84 ± 0.09	1.82 ± 0.08 2.71 ± 0.08 3.11 ± 0.09
DATASET	-	3.76 ± 0.08	4.59 ± 0.07

[†]Our replication

Результаты:

Двухслойная версия SaShiMi побеждает ведущие модели на датасете YouTubeMix, имея меньше параметров.

Model	NLL	Тіме/еросн	Тнгоиднрит	Params
SampleRNN-2 TIER SampleRNN-3 TIER	1.762 1.723	800s 850s	112K SAMPLES/S 116K SAMPLES/S	51.85M $35.03M$
WAVENET-512	1.467 1.449	1000s	185K SAMPLES/S	2.67M
WAVENET-1024		1435s	182K SAMPLES/S	4.24M
SaShiMi-2 layers	1.446	205s	596K SAMPLES/S	1.29M
SaShiMi-4 layers	1.341	340s	316K SAMPLES/S	2.21M
SaShiMi-6 layers	1.315	675s	218K SAMPLES/S	3.13M
SaShiMi-8 layers	1.294	875s	129K SAMPLES/S	4.05M
ISOTROPIC S4-4 LAYERS ISOTROPIC S4-8 LAYERS	1.429	1900s	144K SAMPLES/S	2.83M
	1.524	3700s	72K SAMPLES/S	5.53M

Результаты:

SaShiMi - первая AR модель, которая может генерировать сэмплы высокого качества на датасете SC09. В сочетании с DiffWave получаем SOTA архитектуру.

Model	Params	NLL	FID ↓	IS ↑	мIS ↑	$\mathrm{AM}\downarrow$	Human (κ)	MOS			
WODEL	1 1110111110	TUBB	112 4	10	MILO	11111 4	AGREEMENT	QUALITY	Intelligibility	DIVERSITY	
SAMPLERNN	35.0M	2.042	8.96	1.71	3.02	1.76	0.321	1.18 ± 0.04	1.37 ± 0.02	2.26 ± 0.10	
WAVENET	4.2M	1.925	5.08	2.27	5.80	1.47	0.408	1.59 ± 0.06	1.72 ± 0.03	2.70 ± 0.11	
SaShiMi	4.1M	1.891	1.99	4.12	24.57	0.90	0.832	3.29 ± 0.07	3.53 ± 0.04	3.26 ± 0.09	
WAVEGAN	19.1M	-	2.03	4.90	36.10	0.80	0.840	2.98 ± 0.07	3.27 ± 0.04	3.25 ± 0.09	
DIFFWAVE	24.1M	_	1.92	5.26	51.21	0.68	0.917	4.03 ± 0.06	4.15 ± 0.03	3.45 ± 0.09	
w/ SaShiMi	23.0M	-	1.42	5.94	69.17	0.59	0.953	$\boldsymbol{4.20 \pm 0.06}$	4.33 ± 0.03	3.28 ± 0.11	
TRAIN	_	-	0.00	8.56	292.5	0.16	0.921	4.04 ± 0.06	4.27 ± 0.03	3.59 ± 0.09	
Test	-	-	0.02	8.33	257.6	0.19	0.921	4.04 ± 0.00	4.27 ± 0.03	3.59 ± 0.09	

Результаты:

SaShiMi по всем метрикам лучше WaveNet на разных этапах обучения (в качестве слоя DiffWave). Потенциально можно улучшить все модели, использующие в себе WaveNet.

Architecture 1	Params	Training Steps	FID ↓	IS ↑	MIS ↑	AM ↓	NDB ↓	Human (κ)	MOS		
THICH TECTORE	AGREEMENT QUALITY INT	Intelligibility	DIVERSITY								
SaShiMi	23.0M	800к	1.42	5.94	69.17	0.59	0.88	0.953	4.20 ± 0.06	4.33 ± 0.03	3.28 ± 0.11
WaveNet	24.1M	1000к	1.92	5.26	51.21	0.68	0.88	0.917	4.03 ± 0.06	4.15 ± 0.03	3.45 ± 0.09
SaShiMi	23.0M	500к	2.08	5.68	51.10	0.66	0.76	0.923	3.99 ± 0.06	4.13 ± 0.03	3.38 ± 0.10
WaveNet	24.1M	500к	2.25	4.68	34.55	0.80	0.90	0.848	3.53 ± 0.07	3.69 ± 0.03	3.30 ± 0.08
SaShiMi (uni.)	7.1M	500к	2.70	3.62	17.96	1.03	0.90	0.829	3.08 ± 0.07	3.29 ± 0.04	3.26 ± 0.08
SaShiMi	7.5M	500к	1.70	5.00	40.27	0.72	0.90	0.934	3.83 ± 0.07	4.00 ± 0.03	3.34 ± 0.09
WaveNet	6.8M	500к	4.53	2.80	9.02	1.30	0.94	0.446	1.85 ± 0.08	1.90 ± 0.03	3.03 ± 0.10

Источники

- https://arxiv.org/pdf/2110.13985.pdf
- https://arxiv.org/pdf/2111.00396.pdf
- https://arxiv.org/pdf/2011.04006.pdf
- https://arxiv.org/pdf/2202.09729v1.pdf
- https://srush.github.io/annotated-s4/
- https://www.youtube.com/watch?v=Za84WQ7uM68
- https://www.youtube.com/watch?v=luCBXCErkCs