

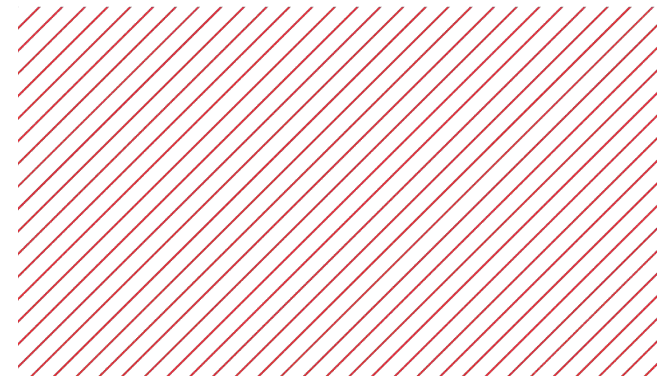
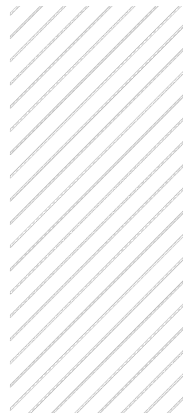
академия
больших
данных

mail.ru
group



Нормализующие потоки в TTS

Свищев Алексей
ЦРТ, ведущий научный сотрудник



Нормализующие потоки

Вывод

Пусть:

$$X = x_i, i = 1, 2, \dots, N \quad (1)$$

наблюдаемые данные
многомерного случайный вектора с неизвестной плотностью $p(X)$

Хотим найти максимально правдоподобные параметры модели для аппроксимации данных, порожденных неизвестным распределением

$$L(X|\theta) = \log(p(X|\theta)) = \sum_{i=1}^N \log p(x_i|\theta) \quad (2)$$

Максимизация логарифмической функции правдоподобия позволяет найти решение задачи

Равенство (2) возможно в силу предположения о независимости наблюдаемых данных.

$$p(X|\theta) = \prod_{i=1}^N p(x_i|\theta)$$

Нормализующие потоки

Вывод

Исходное истинное распределение вероятностей $p(X)$ а аппроксимируем

некоторым параметризованным распределением $p_X(X|\theta)$ которое задается нейросетью

Один из способов аппроксимации - применить к нормальному распределению

$$p_Z(z) = \mathcal{N}(z; 0, 1)$$

несколько нелинейных преобразований, сохраняя полную

$$\int p_Z(z) dz = 1 \quad (3)$$

$$p_X(X|\theta)$$

$$x = x(z)$$

хотим упростить, и $\frac{dx(z)}{dz} \neq 0 \forall z$; тему координат с помощью преобразования

$$\int p_X(x) dx = \int p_X(x(z)) \frac{dx(z)}{dz} dz = \int p_Z(z) \frac{dx(z)}{dz} dz \quad (4)$$

$$\int_a^b f(x) dx = \int_a^b f[\varphi(t)] \varphi'(t) dt.$$

Нормализующие потоки

Вывод

Пусть

$$z = z(x) = f(x) \quad (5)$$

и существует

$$g(z) = f^{-1}(z(x)) = x \quad (6)$$

и

$$\frac{dg(z)}{dz} \neq 0, \forall z \quad \frac{df(x)}{dx} \neq 0, \forall x \quad (7)$$

тогда (4) можно переписать

$$\int p_X(x) dx = \int p_Z(z) \left| \frac{dg(z)}{dz} \right| dz = \int p_Z(z) \left| \left(\frac{dz(x)}{dx} \right)^{-1} \right| dz = \int p_Z(z) \left| \left(\frac{df(x)}{dx} \right)^{-1} \right| dz \quad (8)$$

Нормализующие потоки

Вывод

Трансформация одного распределения с сохранением нормировки и есть простейший поток (*flow*).

Совокупность потоков строится как последовательное применение функций преобразования (композиции функций):

$$f : f = f_1 \circ f_2 \circ \dots \circ f_M$$

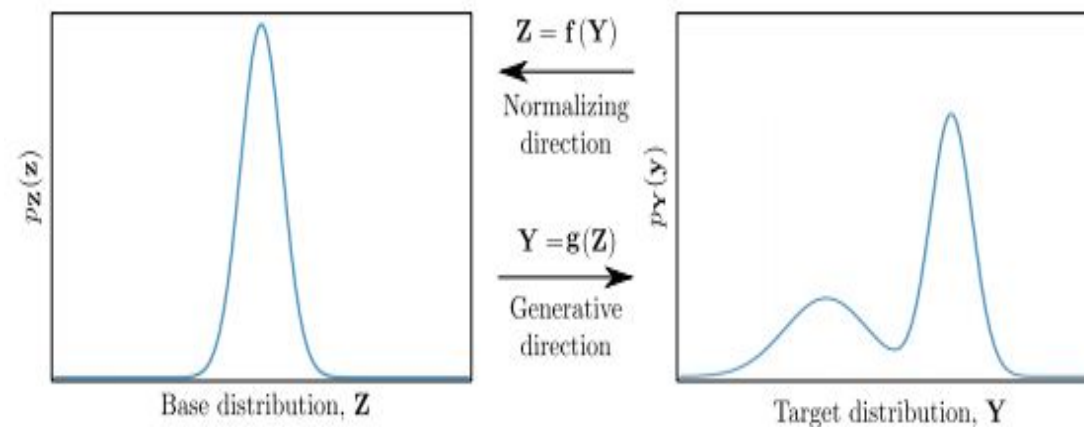
(9)

$f(x)$ нормализующий поток (normalizing flow)

$g(z)$ генеративный поток (generative flow)

$$g = g_M \circ g_{M-1} \circ \dots \circ g_1$$

(10)



Идея метода

на обучении нормализующий поток

используя обратимость $f(x)$

на инференсе генеративный

Нормализующие потоки функция потерь

Потребуем

$$\frac{\partial f_{ij}(x)}{\partial x} = \frac{\partial z_i}{\partial x_j} \neq 0, \forall i, j = 1, \dots, n \quad (11)$$

также пусть существует
и ее производные не равны нулю

$$g_{ij}(z) = f_{ij}^{-1}(z(x)) \frac{\partial g_{ij}}{\partial z} \neq 0 \quad (12)$$

Тогда можно составить матрицу
Якоби
которая будет обратимой

$$J_{f(x)} = \begin{pmatrix} \frac{\partial z_1}{\partial x_1} & \dots & \frac{\partial z_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial z_n}{\partial x_1} & \dots & \frac{\partial z_n}{\partial x_n} \end{pmatrix} = \begin{pmatrix} \frac{\partial f_1(x)}{\partial x_1} & \dots & \frac{\partial f_1(x)}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_n(x)}{\partial x_1} & \dots & \frac{\partial f_n(x)}{\partial x_n} \end{pmatrix} \quad (13)$$

Замена переменных в многомерном
случае

$$\int_{-\infty}^{\infty} p_X(x_1, \dots, x_n) dx_1 \dots dx_n = \int_{-\infty}^{\infty} p_Z(z_1 \dots z_n) |\det J_{f(z)}|^{-1} dz_1 \dots dz_n \quad (14)$$

Нормализующие потоки функция потерь

Рассматривая M последовательно
применяемых потоков

$$f : f = f_1 \circ f_2 \circ \dots \circ f_M \quad (15)$$

Согласно правилу
дифференцирования
сложной функции детерминант
Якобиана такой цепочки

$$\det J_{f(z)} = \prod_{m=1}^M \det J_{f_m(z)} \quad (16)$$

$$\det J_{g(x)} = \prod_{m=1}^M \det J_{g_m(x)} \quad (17)$$

Результат применения совокупности
нормализующих потоков к исходному
простому

$$p_X(x|\theta) = p_Z(f(x)) |\det J_g(x)(x)| = p_Z(f(x)) \prod_{m=1}^M |\det J_{g_m(x)}| \quad (18)$$

Нормализующие потоки функция потерь

В левой части (18) аппроксимирует $p(x_i|\theta)$
из формулы (2)

Подставим (16) и (18) в (2)

$$\begin{aligned} L(X|\theta) &= \sum_{i=1}^N \log p_X(x_i|\theta) = \sum_{i=1}^N \log (p_Z(z) | \prod_{m=1}^M \det J_{g_m}(x, \theta) |) \\ \sum_{i=1}^N \log (p_Z(z) | \prod_{m=1}^M \det J_{g_m}(x, \theta) |) &= \sum_{i=1}^N \log (p_Z(z)) - \sum_{i=1}^N \sum_{m=1}^M \log |\det J_{f_m}(z, \theta)| \end{aligned} \quad (19)$$

$$\det J_{f_m}(x, \theta) = |\det J_{g_m}(z, \theta)|^{-1}$$

Нормализующие потоки Flowtron

В качестве исходного распределения в модели предлагается использовать нормальное

$$z \sim N(z; 0, 1) \quad (20)$$

Преобразование исходного распределения помощью потоков выглядит так

$$\begin{aligned} p_X(x) &\rightarrow p_Z(z) : z = f_1 \circ f_2 \circ \dots \circ f_M(x) \\ p_Z(z) &\rightarrow p_X(x) : x = f_M^{-1} \circ f_{M-1}^{-1} \circ \dots \circ f_1^{-1}(z) = g_M \circ g_{M-1} \circ \dots \circ g_1(z) \end{aligned} \quad (21)$$

В качестве функции потока был выбран аффинный связывающий слой

$$x_t^{m-1} = f_m(x_t^m) = s_t^k \circ x_t^m + b_t^m \quad (22)$$

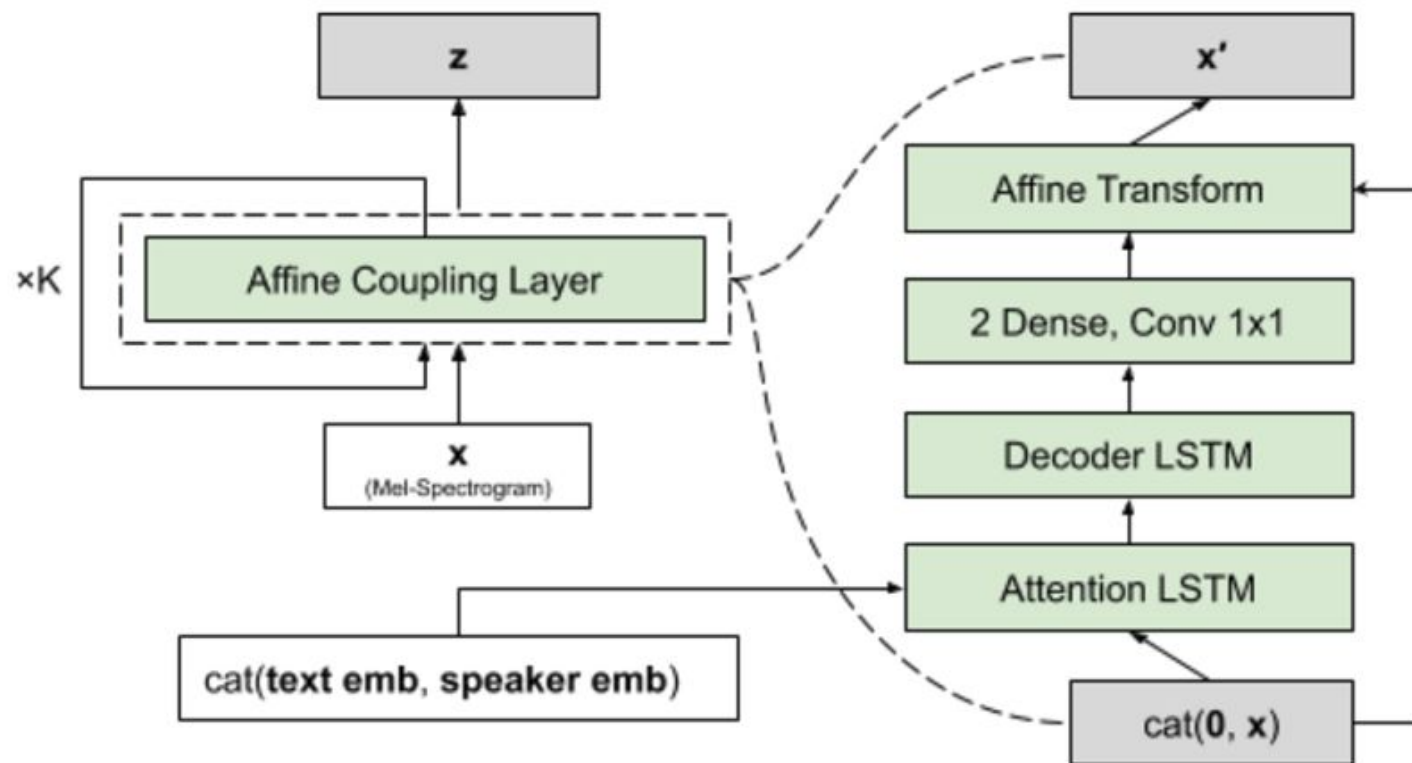
$$(\log s_t^m, b_t^m) = NN_m(x_{1:t-1}^m, speaker, text) \quad (23)$$

Нормализующие потоки Flowtron

В обратном направлении первым входом в поток служит Mel-спектрограмма известной фразы (*text*), сказанной диктором (*speaker*).

Нейронная сеть:

- выстраивает attention
- вычисляет scale и bias



Нормализующие потоки Flowtron

Если хотим генерировать речь на основе текста и диктора то:

$$x_t^m = g_k(x_t^{m-1}) = \frac{x_t^{m-1} - b_t^m}{s_t^m} \quad (24)$$

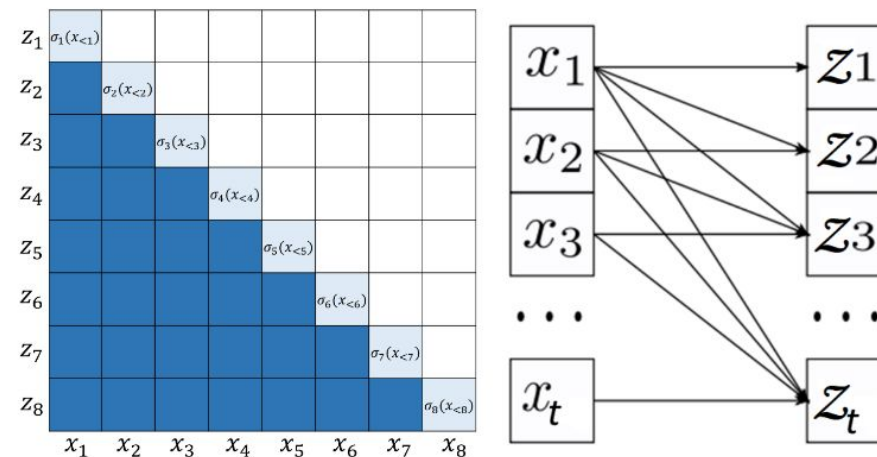
$$(\log s_t^m, b_t^m) = NN_m(x_{1:t-1}^{m-1}, speaker, text) \quad (25)$$

Коэффициенты предсказания потока зависят от результатов предыдущего предсказания

$$x_t^{m-1} = f_m(x_{1:t-1}^m) \quad (26)$$

Из этого следует что матрица зависимостей является ниже-треугольной, а значит

$$\det J_{f_m(x)} = \prod_{t=1}^D \frac{\partial f_m(x_t)}{\partial x_t} \quad (27)$$



Нормализующие потоки Flowtron

Принимая во внимание (22)
производная

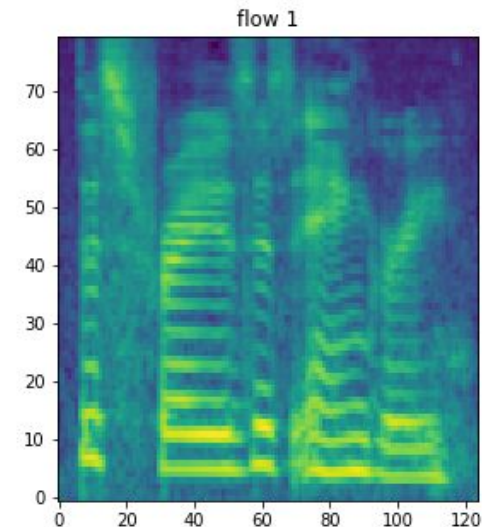
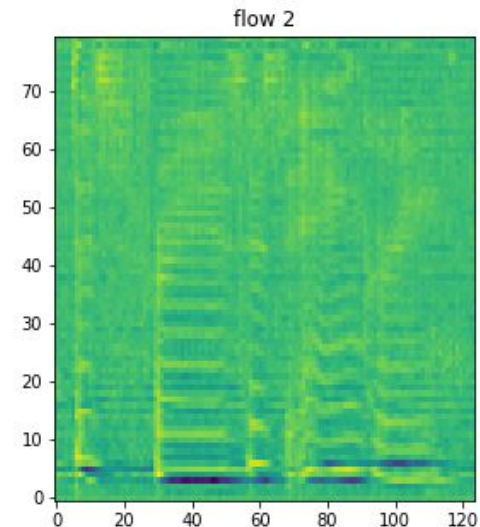
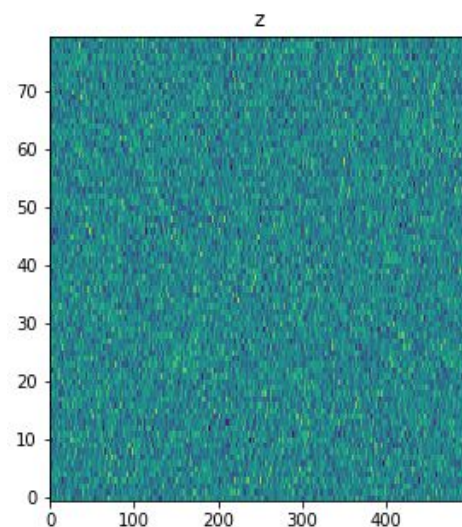
$$f_m(x_t) = s_t^m \circ x_t^m + b_t^m \Rightarrow \frac{\partial f_m(x_t)}{\partial x_t} = s_t^m \quad (28)$$

Подставляя в (19) выражения (27) и (28) получим

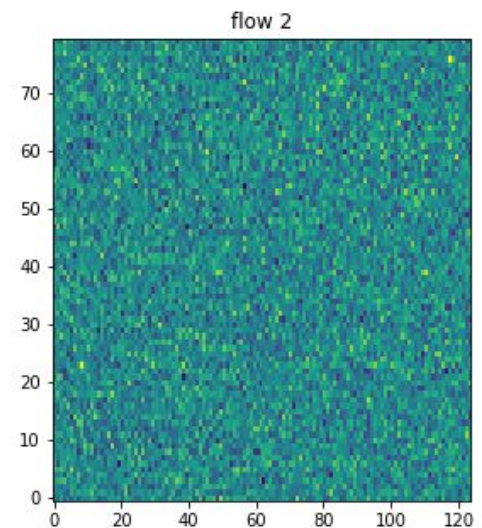
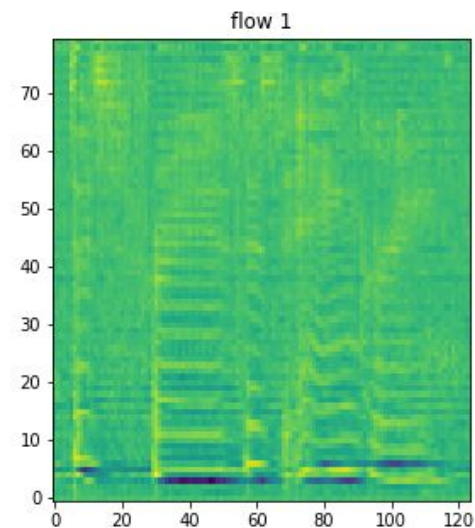
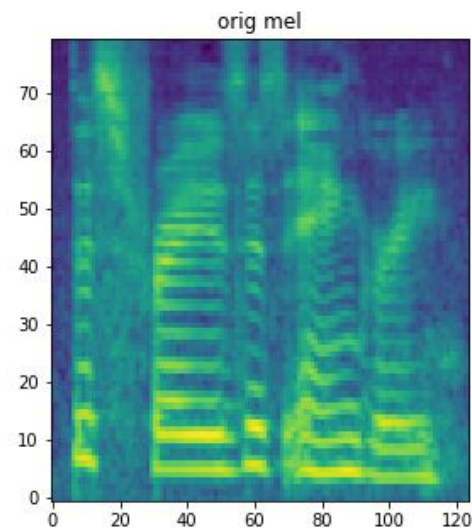
$$L(X|\theta) = \sum_{i=1}^N \log p_Z(z) - \sum_{i=1}^N \sum_{m=0}^M \log |\det J_{f_m(x)}| = \sum_{i=1}^N \log p_Z(z) - \sum_{i=1}^N \sum_{m=0}^M \sum_{t=0}^D \log |\det s_t^m| \quad (29)$$

Нормализующие потоки Flowtron

Направление генерации



Направление нормализации



Нормализующие потоки некоторые функции потоков

Linear Flows (A - обратимая

$$\mathbf{g}(\mathbf{x}) = \mathbf{A}\mathbf{x} + \mathbf{b}$$

Функции потоков определяют выразительность потока

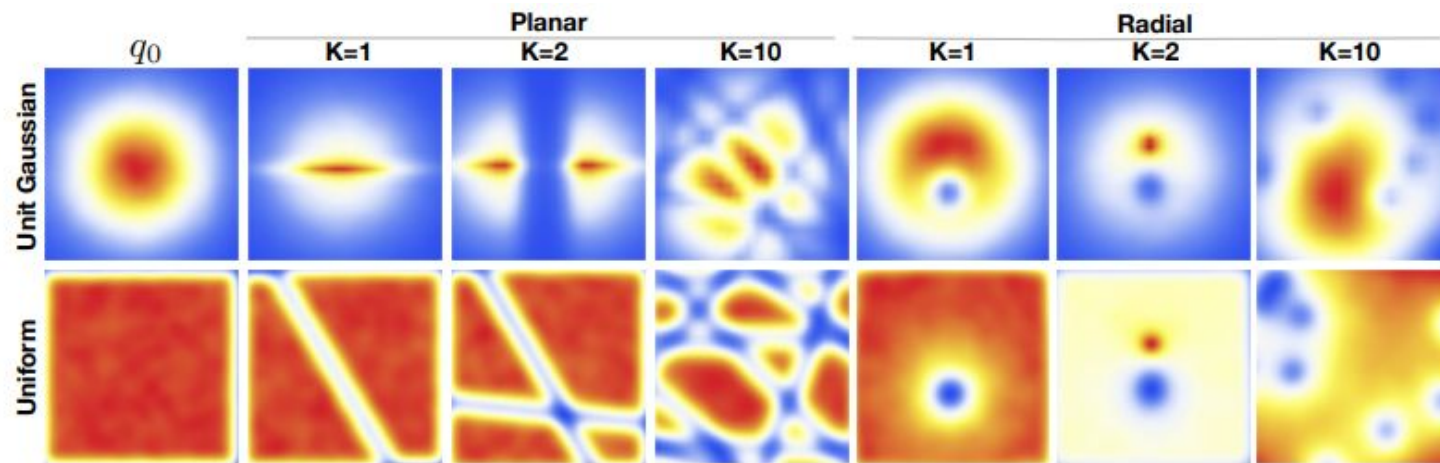
От функции зависит вычисление Якобиана

Planar and Radial Flows

$$\mathbf{g}(\mathbf{x}) = \mathbf{x} + \mathbf{u}h(\mathbf{w}^T\mathbf{x} + b)$$

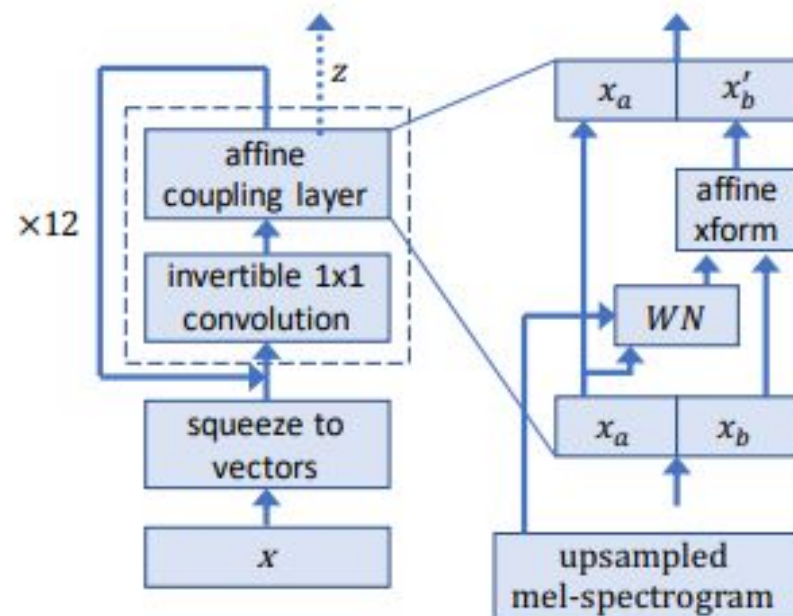
$$\begin{aligned}\det\left(\frac{\partial \mathbf{g}}{\partial \mathbf{x}}\right) &= \det(\mathbf{1}_D + \mathbf{u}h'(\mathbf{w}^T\mathbf{x} + b)\mathbf{w}^T) \\ &= 1 + h'(\mathbf{w}^T\mathbf{x} + b)\mathbf{u}^T\mathbf{w},\end{aligned}$$

$$\mathbf{g}(\mathbf{x}) = \mathbf{x} + \frac{\beta}{\alpha + \|\mathbf{x} - \mathbf{x}_0\|}(\mathbf{x} - \mathbf{x}_0)$$



Нормализующие потоки WaveGlow

- Архитектура схожа с Wavenet
- Параллельный инференс (быстро)
- Последовательное обучение (долго)
- Требуется шумоподавление



$$\begin{aligned} \log p_{\theta}(\mathbf{x}) = & -\frac{\mathbf{z}(\mathbf{x})^T \mathbf{z}(\mathbf{x})}{2\sigma^2} \\ & + \sum_{j=0}^{\#coupling} \log s_j(\mathbf{x}, mel-spectrogram) \\ & + \sum_{k=0}^{\#conv} \log \det |\mathbf{W}_k| \end{aligned}$$

$$\begin{aligned} \mathbf{x}_a, \mathbf{x}_b &= split(\mathbf{x}) \\ (\log \mathbf{s}, \mathbf{t}) &= WN(\mathbf{x}_a, mel-spectrogram) \\ \mathbf{x}_b' &= \mathbf{s} \odot \mathbf{x}_b + \mathbf{t} \\ \mathbf{f}_{coupling}^{-1}(\mathbf{x}) &= concat(\mathbf{x}_a, \mathbf{x}_b') \end{aligned}$$