КАК РАБОТАЮТ ДИФФУЗИОННЫЕ МОДЕЛИ

Сергей Николенко СПбГУ— Санкт-Петербург 20 ноября 2025 г.

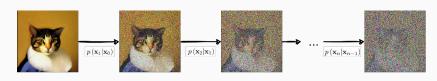




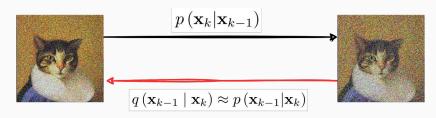
Random facts:

- 20 ноября в ООН Всемирный день ребёнка; в этот день в 1959 г. была принята «Декларация прав ребёнка», а в 1989 — «Конвенция прав ребёнка»
- 20 ноября 1805 г. венский мясник Иоганн Ланер изобрёл сосиски
- 20 ноября 1910 г. Франсиско Мадеро опубликовал «план Сан-Луис-Потоси», в котором объявлял результаты президентских выборов недействительными и призывал к борьбе с режимом Порфирио Диаса, правившего к тому времени 34 года, чем запустил Мексиканскую революцию; а в доме Ивана Озолина, начальника станции Астапово, в тот же день умер Лев Толстой
- 20 ноября 1947 г. в Вестминстерском аббатстве принцесса Елизавета вышла замуж за лейтенанта Филипа Маунтбэттена, который при этом стал герцогом Эдинбургским
- 20 ноября 1979 г. у стен Каабы появился человек и объявил, что настало время для исполнения древнего пророчества, во имя которого в пределах Запретной Мечети была пролита кровь; захват заложников длился до 4 декабря, погибло 255 человек
- 20 ноября 1985 г. вышла первая графическая операционная система от Microsoft, Windows 1.0

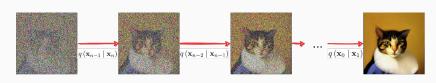
- (Sohl-Dickstein et al., 2015): Deep Unsupervised Learning using Nonequilibrium Thermodynamics
- Идея из статистической физики: давайте обучать марковскую цепь, которая постепенно сделает из случайного шума нужное распределение
- Для этого начнём с цепи, которая сделает из нужного распределения случайный шум:



• А потом попробуем обратить каждый шаг:



• Если получится, то потом мы можем начать со случайного шума и двигаться обратно:

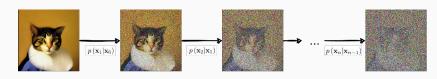


• Forward diffusion: начнём с $p\left(\mathbf{x}_{0}\right)$ и будем добавлять гауссовский шум:

$$q\left(\mathbf{x}_{t} \mid \mathbf{x}_{t-1}\right) = N\left(\mathbf{x}_{t} \middle| \sqrt{1-\beta_{t}} \mathbf{x}_{t-1}, \beta_{t} \mathbf{I}\right).$$

• Тогда постепенно \mathbf{x}_0 превращается в шум, скажем, за T шагов:

$$q\left(\mathbf{x}_{1:T} \mid \mathbf{x}_{0}\right) = \prod_{t=1}^{T} q\left(\mathbf{x}_{t} \mid \mathbf{x}_{t-1}\right).$$



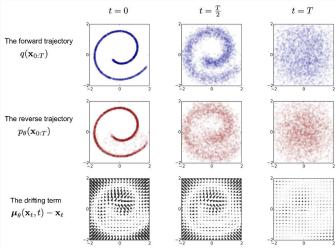
• И для каждого t мы можем просэмплировать \mathbf{x}_t в замкнутой форме через reparametrization trick: пусть $\epsilon \sim N\left(0,\mathbf{I}\right)$ и обозначим $\alpha_t=1-\beta_t,\,A_t=\prod_{i=1}^T\alpha_i$; тогда

$$\begin{split} \mathbf{x}_t &= \sqrt{\alpha_t} \mathbf{x}_{t-1} + \sqrt{1 - \alpha_t} \epsilon = \sqrt{\alpha_t \alpha_{t-1}} \mathbf{x}_{t-2} + \sqrt{1 - \alpha_t \alpha_{t-1}} \epsilon = \dots \\ \dots &= \sqrt{A_t} \mathbf{x}_0 + \sqrt{1 - A_t} \epsilon, \text{ TO еСТЬ } q\left(\mathbf{x}_t \mid \mathbf{x}_0\right) = N\left(\mathbf{x}_t \middle| \sqrt{A_t} \mathbf{x}_0, (1 - A_t) \mathbf{I}\right) \end{split}$$

потому что когда мы складываем два гауссиана $N\left(0,\sigma_1^2\mathbf{I}\right)$ и $N\left(0,\sigma_1^2\mathbf{I}\right)$, получается гауссиан $N\left(0,\left(\sigma_1^2+\sigma_2^2\right)\mathbf{I}\right)$, и у нас это

$$(1-\alpha_t)+\alpha_t\left(1-\alpha_t\right)=1-\alpha_t\alpha_{t-1}.$$

. А теперь задача, казалось бы, в том, чтобы обратить этот процесс, научиться сэмплировать из $q\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\right)$



- Для маленьких eta_t обратные распределения тоже будут гауссовскими, но параметры их теперь уже зависят от точек данных.
- · Reverse diffusion: теперь уже надо обучать распределения

$$p_{\theta}\left(\mathbf{x}_{0:T}\right) = p_{\theta}\left(\mathbf{x}_{T}\right) \prod_{t=1}^{T} p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\right), \quad \text{где}$$

$$p_{\boldsymbol{\theta}}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\right) = N\left(\mathbf{x}_{t-1} \middle| \boldsymbol{\mu}_{\boldsymbol{\theta}}(\mathbf{x}_{t}, t), \boldsymbol{\Sigma}_{\boldsymbol{\theta}}(\mathbf{x}_{t}, t)\right).$$

• Если бы мы знали \mathbf{x}_0 , то обратить было бы можно аналитически, кстати:

$$\begin{split} q\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}, \mathbf{x}_{0}\right) &= q\left(\mathbf{x}_{t} \mid \mathbf{x}_{t-1}, \mathbf{x}_{0}\right) \frac{q\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{0}\right)}{q\left(\mathbf{x}_{t} \mid \mathbf{x}_{0}\right)} \propto \\ &\propto e^{-\frac{1}{2}\left(\frac{\left(\mathbf{x}_{t} - \sqrt{\alpha_{t}}\mathbf{x}_{t-1}\right)^{2}}{\beta_{t}} + \frac{\left(\mathbf{x}_{t-1} - \sqrt{A_{t-1}}\mathbf{x}_{0}\right)^{2}}{1 - A_{t-1}} + \frac{\left(\mathbf{x}_{t} - \sqrt{A_{t}}\mathbf{x}_{0}\right)^{2}}{1 - A_{t}}\right)} \end{split}$$

и в этой формуле теперь надо выделить полный квадрат по \mathbf{x}_{t-1} ...

• Получится

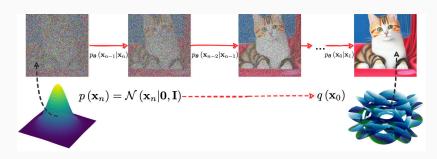
$$q\left(\mathbf{x}_{t-1}\mid\mathbf{x}_{t},\mathbf{x}_{0}
ight)=N\left(\mathbf{x}_{t-1}\middle|\tilde{\mu}\left(\mathbf{x}_{t},\mathbf{x}_{0}
ight),\tilde{eta}_{t}\mathbf{I}
ight),$$
 где

$$\begin{split} \tilde{\beta}_t &= \frac{1 - A_{t-1}}{1 - A_t} \cdot \beta_t, \\ \tilde{\mu}\left(\mathbf{x}_t, \mathbf{x}_0\right) &= \frac{\sqrt{\alpha_t} \left(1 - A_{t-1}\right)}{1 - A_t} \mathbf{x}_t + \frac{\sqrt{A_{t-1}} \beta_t}{1 - A_t} \mathbf{x}_0, \end{split}$$

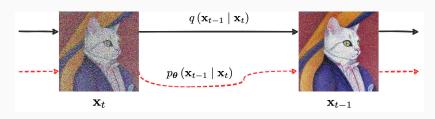
и поскольку мы знаем, что $\mathbf{x}_t = \sqrt{A_t}\mathbf{x}_0 + \sqrt{1-A_t}\epsilon$,

$$\tilde{\mu}\left(\mathbf{x}_{t},\mathbf{x}_{0}\right)=\frac{1}{\sqrt{A_{t}}}\left(\mathbf{x}_{t}-\frac{1-\alpha_{t}}{\sqrt{1-A_{t}}}\epsilon_{t}\right).$$

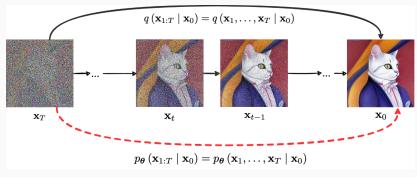
• В конечном счёте хочется, чтобы обратная диффузия реконструировала $q\left(\mathbf{x}_{0}\right)$ из стандартного входа в \mathbf{x}_{n} ; примерно так:



- Что делать? Как всегда в байесовском выводе, приближать!
- На каждом шаге модель должна стать хорошим приближением к $q\left(\mathbf{x}_{t} \mid \mathbf{x}_{t-1}\right)$, без условия на неизвестное \mathbf{x}_{0} :



- Чтобы получить эту аппроксимацию, нужна вариационная нижняя оценка, похожая на ту, которая была в VAE и DALL-E
- Начинаем с оценки на глобальное распределение $q(\mathbf{x}_{1:T} \mid \mathbf{x}_0) = q(\mathbf{x}_1, ..., \mathbf{x}_T \mid \mathbf{x}_0)$:



• А потом она разложится на оценки для отдельных шагов диффузионного процесса

• Вспомним идею вариационных приближений:

$$\begin{split} p\left(\mathbf{x}, \mathbf{z}\right) &= p\left(\mathbf{x}\right) p\left(\mathbf{z} | \mathbf{x}\right), \\ \log p\left(\mathbf{x}\right) &= \log p\left(\mathbf{x}, \mathbf{z}\right) - \log p\left(\mathbf{z} | \mathbf{x}\right), \end{split}$$

возьмём ожидание по $q\left(\mathbf{z}\right)$:

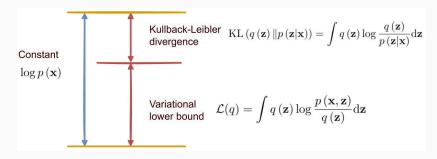
$$\mathbb{E}_{q(\mathbf{z})}\left[\log p\left(\mathbf{x}\right)\right] = \mathbb{E}_{q(\mathbf{z})}\left[\log p\left(\mathbf{x},\mathbf{z}\right)\right] - \mathbb{E}_{q(\mathbf{z})}\left[\log p\left(\mathbf{z}|\mathbf{x}\right)\right],$$

а затем сделаем стандартные преобразования:

$$\begin{split} \log p\left(\mathbf{x}\right) = & \mathbb{E}_{q(\mathbf{z})} \left[\log p\left(\mathbf{x}, \mathbf{z}\right) \right] - \mathbb{E}_{q(\mathbf{z})} \left[\log q\left(\mathbf{z}\right) \right] + \\ & + \mathbb{E}_{q(\mathbf{z})} \left[\log q\left(\mathbf{z}\right) \right] - \mathbb{E}_{q(\mathbf{z})} \left[\log p\left(\mathbf{z}|\mathbf{x}\right) \right], \\ \log p\left(\mathbf{x}\right) = & \int q\left(\mathbf{z}\right) \log \frac{p\left(\mathbf{x}, \mathbf{z}\right)}{q\left(\mathbf{z}\right)} \mathrm{d}\mathbf{z} + \int q\left(\mathbf{z}\right) \log \frac{q\left(\mathbf{z}\right)}{p\left(\mathbf{z}|\mathbf{x}\right)} \mathrm{d}\mathbf{z}. \end{split}$$

• Итого получается, что

$$\log p\left(\mathbf{x}
ight) = L(q) + \mathrm{KL}\left(q\left(\mathbf{z}
ight) \| p\left(\mathbf{z}|\mathbf{x}
ight)
ight),$$
 где $L(q) = \int q\left(\mathbf{z}
ight) \log rac{p\left(\mathbf{x},\mathbf{z}
ight)}{q\left(\mathbf{z}
ight)} \mathrm{d}\mathbf{z}.$



• Давайте для диффузионной модели выведем вариационную оценку снова из первых принципов:

$$\log p_{\theta}\left(\mathbf{x}_{0}\right) = \log p_{\theta}\left(\mathbf{x}_{0}, \ldots, \mathbf{x}_{T}\right) - \log p_{\theta}\left(\mathbf{x}_{1}, \ldots, \mathbf{x}_{T} \mid \mathbf{x}_{0}\right),$$

$$\begin{split} & \mathbb{E}_{q(\mathbf{x}_0)}\left[\log p_{\theta}\left(\mathbf{x}_0\right)\right] = \mathbb{E}_{q(\mathbf{x}_{0:T})}\left[\log p_{\theta}\left(\mathbf{x}_{0:T}\right)\right] - \mathbb{E}_{q(\mathbf{x}_{0:T})}\left[\log p_{\theta}\left(\mathbf{x}_{1:T} \mid \mathbf{x}_0\right)\right] = \\ & = \mathbb{E}_{q(\mathbf{x}_{0:T})}\left[\log \frac{p_{\theta}\left(\mathbf{x}_{0:T}\right)}{q\left(\mathbf{x}_{1:T} \mid \mathbf{x}_0\right)}\right] + \mathbb{E}_{q(\mathbf{x}_{0:T})}\left[\log \frac{q\left(\mathbf{x}_{1:T} \mid \mathbf{x}_0\right)}{p_{\theta}\left(\mathbf{x}_{1:T} \mid \mathbf{x}_0\right)}\right]. \end{split}$$

• Здесь второе слагаемое — это КL-дивергенция, и мы получаем функцию ошибки как минус вариационную нижнюю оценку:

$$L = \mathbb{E}_{q(\mathbf{x}_{0:T})} \left[\log \frac{q\left(\mathbf{x}_{1:T} \mid \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{0:T}\right)} \right] \geq -\mathbb{E}_{q(\mathbf{x}_{0})} \left[\log p_{\theta}\left(\mathbf{x}_{0}\right) \right].$$

• Эту оценку можно подсчитать как сумму:

$$\begin{split} L = & \mathbb{E}_{q} \left[\log \frac{q\left(\mathbf{x}_{1:T} \mid \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{0:T}\right)} \right] = \mathbb{E}_{q} \left[\log \frac{\prod_{t=1}^{T} q\left(\mathbf{x}_{t} \mid \mathbf{x}_{t-1}\right)}{p_{\theta}\left(\mathbf{x}_{T}\right) \prod_{t=1}^{T} p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\right)} \right] \\ = & \mathbb{E}_{q} \left[-\log p_{\theta}\left(\mathbf{x}_{T}\right) + \sum_{t=1}^{T} \log \frac{q\left(\mathbf{x}_{t} \mid \mathbf{x}_{t-1}\right)}{p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\right)} \right] \\ = & \mathbb{E}_{q} \left[-\log p_{\theta}\left(\mathbf{x}_{T}\right) + \sum_{t=2}^{T} \log \frac{q\left(\mathbf{x}_{t} \mid \mathbf{x}_{t-1}\right)}{p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\right)} + \log \frac{q\left(\mathbf{x}_{1} \mid \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{0} \mid \mathbf{x}_{1}\right)} \right] \\ = & \mathbb{E}_{q} \left[-\log p_{\theta}\left(\mathbf{x}_{T}\right) + \sum_{t=2}^{T} \log \left(\frac{q\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}, \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{0}\right)} + \log \frac{q\left(\mathbf{x}_{1} \mid \mathbf{x}_{0}\right)}{q\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{0}\right)} + \log \frac{q\left(\mathbf{x}_{1} \mid \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{0} \mid \mathbf{x}_{1}\right)} \right] \\ = & \mathbb{E}_{q} \left[\log p_{\theta}\left(\mathbf{x}_{T}\right) + \sum_{t=2}^{T} \log \frac{q\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}, \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{0}\right)} + \log \frac{q\left(\mathbf{x}_{1} \mid \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{0} \mid \mathbf{x}_{1}\right)} \right] \\ = & \mathbb{E}_{q} \left[\log \frac{q\left(\mathbf{x}_{T} \mid \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{T}\right)} + \sum_{t=2}^{T} \log \frac{q\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}, \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{0}\right)} - \log p_{\theta}\left(\mathbf{x}_{0} \mid \mathbf{x}_{1}\right) \right]. \end{split}$$

• Итого получается

$$\begin{split} L = & L_T + L_{T-1} + \ldots + L_0, \quad \text{where} \\ L_T = & \text{KL}\left(q\left(\mathbf{x}_T \mid \mathbf{x}_0\right) \| p_{\theta}\left(\mathbf{x}_T\right)\right), \\ L_t = & \text{KL}\left(q\left(\mathbf{x}_t \mid \mathbf{x}_{t+1}, \mathbf{x}_0\right) \| p_{\theta}\left(\mathbf{x}_t \mid \mathbf{x}_{t+1}\right)\right), \quad t = 1, \ldots, T-1, \\ L_0 = & -\log p_{\theta}\left(\mathbf{x}_0 \mid \mathbf{x}_1\right). \end{split}$$

- Всё это KL-дивергенции между гауссианами, их можно подсчитать и подставить в функцию ошибки
- \cdot Например, в L_t мы используем гауссовскую параметризацию

$$p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\right) = N\left(\mathbf{x}_{t-1} \middle| \mu_{\theta}(\mathbf{x}_{t}, t), \Sigma_{\theta}(\mathbf{x}_{t}, t)\right)$$

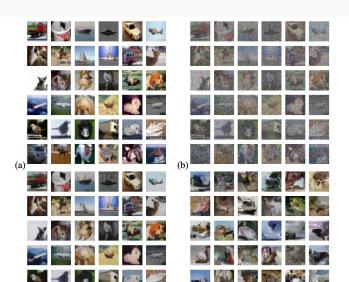
и пытаемся её параметры совместить с $q\left(\mathbf{x}_{t}\mid\mathbf{x}_{t-1},\mathbf{x}_{0}\right)$. Для среднего, например,

$$\mu_{\boldsymbol{\theta}}(\mathbf{x}_t,t) \approx \tilde{\mu}_t\left(\mathbf{x}_t,\mathbf{x}_0\right) = \frac{1}{\sqrt{A_t}}\left(\mathbf{x}_t - \frac{1-\alpha_t}{\sqrt{1-A_t}}\epsilon_t\right),$$

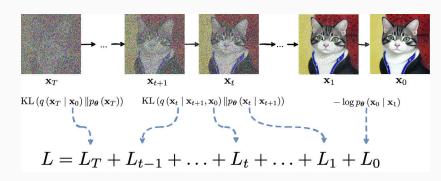
и поскольку мы знаем \mathbf{x}_t при обучении, мы можем параметризовать шум напрямую:

$$p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\right) = N\left(\mathbf{x}_{t-1} \middle| \frac{1}{\sqrt{\alpha_{t}}} \left(\mathbf{x}_{t} - \frac{1-\alpha_{t}}{\sqrt{1-A_{t}}} \epsilon_{\theta}(\mathbf{x}_{t}, t)\right), \Sigma_{\theta}(\mathbf{x}_{t}, t)\right).$$

• (Sohl-Dickstein et al., 2015) само по себе не слишком хорошо работало, разве что на CIFAR:



- Следующий шаг: "Denoising Diffusion Probabilistic Models" (DDPM; Ho et al., 2020)
- Та же основная идея, та же структура функции ошибки:



DDPM

- · Но дальше DDPM делает следующие замечания:
 - дисперсии прямой диффузии β_t постоянные гиперпараметры, их не обучают, так что все q не обучаются; т.к. $p_{\theta}\left(\mathbf{x}_T\right)$ фиксированное распределение, из которого мы будем сэмплировать, L_T является константой;
 - · для промежуточных шагов дисперсии в $p_{\theta}\left(\mathbf{x}_{t}\mid\mathbf{x}_{t+1}\right)$ тоже не обучаются, и можно найти простую замкнутую форму для L_{t} ;
 - \cdot главное давайте введём отдельный дискретный декодер для L_0 ; если данные состоят из целых чисел от 0 до 255, отмасштабированных в [-1,1] (это естественное представление для картинок), DDPM моделирует

$$p_{\boldsymbol{\theta}}\left(\mathbf{x}_{0} \mid \mathbf{x}_{1}\right) = \prod_{i=1}^{D} \int_{\delta_{-}\left(x_{0}, i\right)}^{\delta_{+}\left(x_{0}, i\right)} N\left(x \middle| \mu_{\boldsymbol{\theta}, i}\left(\mathbf{x}_{1}\right), \sigma_{1}^{2}\right) \mathrm{d}x,$$

где i идёт по пикселям, $\mu_{\theta}(\mathbf{x}_1)$ — декодер, пределы интегрирования $\frac{1}{255}$ в каждую сторону от $x_{0,i}$

• Т.е. можно подставить другую модель на последнем шаге и использовать $\mu_{\theta}(\mathbf{x}_1)$ без шума во время сэмплирования

DDPM

• DDPM уже давала отличное порождение, сравнимое с лучшими GAN'ами того времени:



- Следующий шаг пришёл очень скоро: "Denoising Diffusion Implicit Models" (DDIM; Song et al., 2020)
- Важный недостаток всех диффузионных моделей пока что в том, что они *очень медленные*
- Порождение должно пройти по всем шагам диффузии, а их буквально тысячи, и они идут последовательно, не параллелизуются
- Song et al. упоминают, что сэмплирование из обученного GAN примерно в 1000 раз быстрее, чем сэмплирование из DDPM для того же размера картинок

DDIM

- Как можно ускорить диффузионные модели?
- Song et al. обобщают диффузионные модели и в частности DDPM
- Функция ошибки L не зависит напрямую от совместного распределения $q\left(\mathbf{x}_{1:T} \mid \mathbf{x}_{0}\right)$, а только от маргиналов $q\left(\mathbf{x}_{t} \mid \mathbf{x}_{0}\right)$
- Это значит, что мы можем переиспользовать ту же функцию ошибки для другого совместного распределения, если у него те же маргиналы
- В частности, обобщить можно даже на немарковские диффузионные процессы, если для них получится найти обратную марковскую цепь:



- Зададим вероятностную модель для данных \mathbf{x}_0 (из распределения $q(\mathbf{x}_0)$) через латентные переменные $\mathbf{x}_1,\dots,\mathbf{x}_T$
- Сначала прямой проход (noising) $q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0)$, в котором к данным постепенно добавляется шум:

$$q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0) = q_{\sigma}(\mathbf{x}_T \mid \mathbf{x}_0) \, \prod_{t=2}^T q_{\sigma}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0).$$

- \cdot Здесь мы по определению задаём, как связаны $\mathbf{x}_1, \dots, \mathbf{x}_T$ с \mathbf{x}_0
- Важный момент: переходы $q_{\sigma}(\mathbf{x}_{t-1}\mid\mathbf{x}_t,\mathbf{x}_0)$ зависят и от \mathbf{x}_t , и от \mathbf{x}_0 , т.е. процесс $\mathbf{x}_0,\mathbf{x}_1,\dots,\mathbf{x}_T$ не марковский! это осознанное решение

 Мы выбираем следующие распределения (с тем же маргиналом, что и в DDPM):

$$\begin{split} q_{\sigma}(\mathbf{x}_T \mid \mathbf{x}_0) &= \mathcal{N} \big(\sqrt{A_T} \, \mathbf{x}_0, \; (1-A_T) \, \mathbf{I} \big), \\ q_{\sigma}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) &= \mathcal{N} \Big(\sqrt{A_{t-1}} \mathbf{x}_0 + \sqrt{1-A_{t-1} - \sigma_t^2} \, \frac{\mathbf{x}_t - \sqrt{A_t} \mathbf{x}_0}{\sqrt{1-A_t}}, \; \sigma_t^2 \, \mathbf{I} \Big). \end{split}$$

· σ_t^2 – дисперсия добавочного шума в обратном переходе; форма среднего такая для того, чтобы маргиналы $q_\sigma(\mathbf{x}_t\mid\mathbf{x}_0)$ оставались гауссианами и совпадали с DDPM

 \cdot Ключевое утверждение: для любого t

$$q_{\sigma}(\mathbf{x}_t \mid \mathbf{x}_0) = \mathcal{N}\big(\sqrt{A_t}\,\mathbf{x}_0,\; (1-A_t)\,\mathbf{I}\big).$$

• Эквивалентно можно записать в виде «одноступенчатой» формулы:

$$\mathbf{x}_t = \sqrt{A_t}\,\mathbf{x}_0 + \sqrt{1-A_t}\,\epsilon_t, \qquad \epsilon_t \sim \mathcal{N}(\mathbf{0},\,\mathbf{I}).$$

• Это ключевая формула DDPM/DDIM: любое зашумлённое состояние \mathbf{x}_t можно получить напрямую из \mathbf{x}_0 , не прогоняя цепочку шаг за шагом

• Доказательство по индукции. Предположим, что

$$\begin{split} q_{\sigma}(\mathbf{x}_t \mid \mathbf{x}_0) &= \mathcal{N}\big(\sqrt{A_t}\,\mathbf{x}_0,\; (1-A_t)\,\mathbf{I}\big), \quad \text{то есть} \\ \mathbf{x}_t &= \sqrt{A_t}\,\mathbf{x}_0 + \sqrt{1-A_t}\,\epsilon_t, \qquad \epsilon_t \sim \mathcal{N}(\mathbf{0},\mathbf{I}). \end{split}$$

· Подставим это в $q_{\sigma}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0)$, получим

$$q_{\sigma}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) = \mathcal{N}\Big(\sqrt{A_{t-1}}\mathbf{x}_0 + \sqrt{1 - A_{t-1} - \sigma_t^2}\,\epsilon_t, \ \sigma_t^2\,\mathbf{I}\Big).$$

 \cdot Мы просто заменили $\dfrac{\mathbf{x}_t-\sqrt{A_t}\mathbf{x}_0}{\sqrt{1-A_t}}$ на ϵ_t , используя выражение для \mathbf{x}_t через \mathbf{x}_0 и ϵ_t ; среднее стало линейной функцией \mathbf{x}_0 и ϵ_t

• Используем факт про гауссианы: если

$$\begin{split} p(\mathbf{x}) &= \mathcal{N}(\mathbf{x} \mid \boldsymbol{\mu}, \ \boldsymbol{\Lambda}^{-1}), \\ p(\mathbf{y} \mid \mathbf{x}) &= \mathcal{N}(\mathbf{y} \mid A\mathbf{x} + \mathbf{b}, \ L^{-1}), \quad \text{to} \\ p(\mathbf{y}) &= \mathcal{N}(\mathbf{y} \mid A\boldsymbol{\mu} + \mathbf{b}, \ L^{-1} + A\boldsymbol{\Lambda}^{-1}A^{\top}). \end{split}$$

• В нашем случае ${\bf x}$ – это ϵ_t , распределённый как $\mathcal{N}({\bf 0},{\bf I});$ ${\bf y}$ – это ${\bf x}_{t-1};$ матрица A и сдвиг ${\bf b}$ задаются коэффициентами при ϵ_t и ${\bf x}_0$

• Сопоставим параметры:

$$\mu \equiv \mathbf{0}, \quad \Lambda^{-1} \equiv \mathbf{I},$$

$$\mathbf{A} \equiv \sqrt{1 - A_{t-1} - \sigma_t^2} \, \mathbf{I}, \quad \mathbf{b} \equiv \sqrt{A_{t-1}} \, \mathbf{x}_0, \quad L^{-1} \equiv \sigma_t^2 \, \mathbf{I}.$$

• Подставив в формулу, получим

$$q_{\sigma}(\mathbf{x}_{t-1} \mid \mathbf{x}_0) = \mathcal{N}(A\mu + \mathbf{b}, L^{-1} + A\Lambda^{-1}A^{\top}) = \mathcal{N}(\sqrt{A_{t-1}}\mathbf{x}_0, (1 - A_{t-1})\mathbf{I}),$$

потому что

$$\begin{split} A \mu + \mathbf{b} &= \sqrt{A_{t-1}} \, \mathbf{x}_0, \\ L^{-1} + A \Lambda^{-1} A^\top &= \sigma_t^2 \, \mathbf{I} + (1 - A_{t-1} - \sigma_t^2) \, \mathbf{I} = (1 - A_{t-1}) \, \mathbf{I}. \end{split}$$

• Из полученного выражения

$$\mathbf{x}_t = \sqrt{A_t}\,\mathbf{x}_0 + \sqrt{1-A_t}\,\epsilon_t$$

можно выразить \mathbf{x}_0 через \mathbf{x}_t и ϵ_t :

$$\mathbf{x}_0 = \frac{\mathbf{x}_t - \sqrt{1 - A_t} \, \epsilon_t}{\sqrt{A_t}}.$$

- Эта формула важна для порождения, т.к. в обратном процессе мы будем оценивать ϵ_t нейросетью и тем самым получать оценку \mathbf{x}_0 из \mathbf{x}_t

• Поскольку истинный шум ϵ_t нам неизвестен, мы вводим аппроксимацию

$$\epsilon_{\theta}^{(t)}(\mathbf{x}_t) \approx \epsilon_t,$$

где $\epsilon_{ heta}^{(t)}$ – нейросеть, зависящая от \mathbf{x}_t и (неявно) от шага t

 \cdot Тогда оценка \mathbf{x}_0 на шаге t задаётся как

$$f_{\theta}^{(t)}(\mathbf{x}_t) \; \equiv \; \frac{\mathbf{x}_t - \sqrt{1 - A_t} \, \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}{\sqrt{A_t}}. \label{eq:ftheta}$$

• Обратное выражение:

$$\epsilon_{\theta}^{(t)}(\mathbf{x}_t) = \frac{\mathbf{x}_t - \sqrt{A_t} \, f_{\theta}^{(t)}(\mathbf{x}_t)}{\sqrt{1 - A_t}}.$$

- Теперь переходим к обратному (порождающему) процессу
- Совместное распределение модели задаётся как

$$p_{\theta}(\mathbf{x}_{0:T}) = p_{\theta}(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}^{(t)}(\mathbf{x}_{t-1} \mid \mathbf{x}_t),$$

где обычно берут

$$p_{\theta}(\mathbf{x}_T) = \mathcal{N}(\mathbf{0}, \mathbf{I}),$$

а переходы $p_{ heta}^{(t)}$ параметризуются через нейросеть

• Правдоподобие данных:

$$\log p(\mathbf{x}_0) = \log \int p_{\theta}(\mathbf{x}_{0:T}) d\mathbf{x}_{1:T}$$

• Введём вспомогательное распределение – наш прямой процесс $q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0)$, умножим и поделим на него внутри интеграла и применим неравенство Йенсена:

$$\begin{split} \log p(\mathbf{x}_0) &= \log \int p_{\theta}(\mathbf{x}_{0:T}) \, d\mathbf{x}_{1:T} = \log \int \frac{p_{\theta}(\mathbf{x}_{0:T})}{q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \, q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0) \, d\mathbf{x}_{1:T} \\ &= \log \mathbb{E}_{q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \left[\frac{p_{\theta}(\mathbf{x}_{0:T})}{q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \right] \\ &\geq \mathbb{E}_{q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \left[\log p_{\theta}(\mathbf{x}_{0:T}) - \log q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0) \right]. \end{split}$$

· Обозначим $J_{\sigma} = - \mathbb{E}_{q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0)} \left[\log p_{\theta}(\mathbf{x}_{0:T}) - \log q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0) \right],$ подставим факторизации p_{θ} и q_{σ} :

$$\begin{split} q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_0) &= q_{\sigma}(\mathbf{x}_T \mid \mathbf{x}_0) \prod_{t=2}^T q_{\sigma}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0), \\ p_{\theta}(\mathbf{x}_{0:T}) &= p_{\theta}(\mathbf{x}_T) \prod_{t=1}^T p_{\theta}^{(t)}(\mathbf{x}_{t-1} \mid \mathbf{x}_t). \end{split}$$

• Тогда после раскрытия логарифмов получаем

$$\begin{split} J_{\sigma} &= \mathbb{E}_{q_{\sigma}(\mathbf{x}_{1:T} \mid \mathbf{x}_{0})} \Bigg[\log q_{\sigma}(\mathbf{x}_{T} \mid \mathbf{x}_{0}) + \sum_{t=2}^{T} \log q_{\sigma}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}, \mathbf{x}_{0}) \\ &- \sum_{t=1}^{T} \log p_{\theta}^{(t)}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}) - \log p_{\theta}(\mathbf{x}_{T}) \Bigg]. \end{split}$$

• Следующий шаг – перегруппировать всё это так, чтобы получилось что-то похожее на сумму КL-дивергенций, плюс отдельное слагаемое для реконструкции \mathbf{x}_0

• После некоторых мучений получим:

$$\begin{split} J_{\sigma} &= \sum_{t=2}^{T} \mathbb{E}_{q(\mathbf{x}_{t} \mid \mathbf{x}_{0})} \Big[D_{\mathrm{KL}} \big(q_{\sigma}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}, \mathbf{x}_{0}) \, \| \, p_{\theta}^{(t)}(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}) \big) \Big] \\ &- \mathbb{E}_{q(\mathbf{x}_{1} \mid \mathbf{x}_{0})} \big[\log p_{\theta}^{(0)}(\mathbf{x}_{0} \mid \mathbf{x}_{1}) \big]. \end{split}$$

• Здесь $p_{\theta}^{(0)}(\mathbf{x}_0 \mid \mathbf{x}_1)$ – отдельный шаг, как в DDPM; в реальных реализациях его часто опускают или сливают с общей ошибкой

• Параметризация обратного шага и KL-дивергенция: для t>1 возьмём

$$p_{\theta}^{(t)}(\mathbf{x}_{t-1}\mid\mathbf{x}_{t}) = \mathcal{N}\Big(\sqrt{A_{t-1}}\,f_{\theta}^{(t)}(\mathbf{x}_{t}) + \sqrt{1-A_{t-1}-\sigma_{t}^{2}}\,\epsilon_{\theta}^{(t)}(\mathbf{x}_{t}),\;\sigma_{t}^{2}\,\mathbf{I}\Big).$$

- · Это та же структура, что и у $q_\sigma(\mathbf{x}_{t-1}\mid\mathbf{x}_t,\mathbf{x}_0)$, только вместо истинного \mathbf{x}_0 и истинного шума ϵ_t мы подставляем оценки $f_\theta^{(t)}(\mathbf{x}_t)$ и $\epsilon_\theta^{(t)}(\mathbf{x}_t)$
- КL между двумя гауссианами с одинаковой ковариацией имеет особенно простую форму – это просто квадратичная функция разности средних

6

• У нас:

$$\begin{split} q_{\sigma}(\mathbf{x}_{t-1} \mid \mathbf{x}_t, \mathbf{x}_0) &= \mathcal{N}\Big(\sqrt{A_{t-1}}\mathbf{x}_0 + \sqrt{1 - A_{t-1} - \sigma_t^2}\,\epsilon_t, \ \sigma_t^2\,\mathbf{I}\Big), \\ p_{\theta}^{(t)}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) &= \mathcal{N}\Big(\sqrt{A_{t-1}}\,f_{\theta}^{(t)}(\mathbf{x}_t) + \sqrt{1 - A_{t-1} - \sigma_t^2}\,\epsilon_{\theta}^{(t)}(\mathbf{x}_t), \ \sigma_t^2\,\mathbf{I}\Big). \end{split}$$

· Для двух гауссианов с одинаковой ковариацией $\Sigma = \sigma_t^2 \mathbf{I}$ КL-дивергенция равна

$$\mathrm{KL}\big(\mathcal{N}(\boldsymbol{\mu}_q, \boldsymbol{\Sigma}) \, \| \, \mathcal{N}(\boldsymbol{\mu}_p, \boldsymbol{\Sigma})\big) = \frac{1}{2} \, (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p)^\top \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_q - \boldsymbol{\mu}_p).$$

• Подставляя наши средние и $\Sigma^{-1} = \sigma_t^{-2} \mathbf{I}$, получаем

$$\begin{split} &D_{\mathrm{KL}}\big(q_{\sigma}(\mathbf{x}_{t-1}\mid\mathbf{x}_{t},\mathbf{x}_{0})\,\|\,p_{\theta}^{(t)}(\mathbf{x}_{t-1}\mid\mathbf{x}_{t})\big) = \\ &= \frac{1}{2\sigma_{t}^{2}}\Big\|\sqrt{A_{t-1}}(\mathbf{x}_{0} - f_{\theta}^{(t)}(\mathbf{x}_{t})) + \sqrt{1 - A_{t-1} - \sigma_{t}^{2}}(\epsilon_{t} - \epsilon_{\theta}^{(t)}(\mathbf{x}_{t}))\Big\|^{2}. \end{split}$$

• Используя связь

$$\mathbf{x}_0 = \frac{\mathbf{x}_t - \sqrt{1 - A_t} \, \epsilon_t}{\sqrt{A_t}}, \qquad f_{\theta}^{(t)}(\mathbf{x}_t) = \frac{\mathbf{x}_t - \sqrt{1 - A_t} \, \epsilon_{\theta}^{(t)}(\mathbf{x}_t)}{\sqrt{A_t}},$$

можно показать (после упрощений), что KL-член пропорционален

$$D_{\mathrm{KL}}(\cdot) = \frac{\gamma_t}{2\sigma_t^2} \big\| \epsilon_t - \epsilon_{\theta}^{(t)}(\mathbf{x}_t) \big\|^2 + \mathrm{const},$$

где γ_t – некоторая известная функция A_t и σ_t^2 , которую можно заранее посчитать по расписанию шумов, как вес для шага t

• Главный вывод: оптимизация КL между истинным и параметризованным обратным переходом эквивалентна MSE между истинным шумом ϵ_t и предсказанным шумом $\epsilon_{\theta}^{(t)}(\mathbf{x}_t)$, с некоторым весом, зависящим от t

6

• Если собрать все шаги t и учесть ожидание по $\mathbf{x}_0 \sim q(\mathbf{x}_0)$ и $\epsilon_t \sim \mathcal{N}(0,\mathbf{I})$, то окончательная функция ошибки (с весами) имеет вид

$$\mathbb{E}_{t,\mathbf{x}_0,\epsilon_t} \left[\frac{\gamma_t}{2\sigma_t^2} \big\| \epsilon_t - \epsilon_{\theta}^{(t)}(\mathbf{x}_t) \big\|^2 \right].$$

• На практике, как и в DDPM, обычно игнорируют веса $\frac{\gamma_t}{2\sigma_t^2}$, беря более простой loss

$$L_{\mathrm{simple}} = \mathbb{E}_{t,\mathbf{x}_0,\epsilon} \big[\| \epsilon - \epsilon_{\theta}^{(t)}(\mathbf{x}_t) \|^2 \big],$$

где t выборочно сэмплируется из $\{1,\dots,T\}$

- Итого обучение идёт точно как в DDPM: берём \mathbf{x}_0 , сэмплируем \mathbf{x}_t , обновляем параметры сети по ошибке $L=\|\epsilon-\epsilon_{\theta}^{(t)}(\mathbf{x}_t)\|$
- · А вот inference после обучения имеет вид

$$p_{\theta}^{(t)}(\mathbf{x}_{t-1} \mid \mathbf{x}_t) = \mathcal{N}\big(\mu_{\theta}(\mathbf{x}_t, t), \ \sigma_t^2 \, \mathbf{I}\big),$$

где среднее

$$\boldsymbol{\mu}_{\theta}(\mathbf{x}_t,t) = \sqrt{A_{t-1}} \, f_{\theta}^{(t)}(\mathbf{x}_t) + \sqrt{1 - A_{t-1} - \sigma_t^2} \, \boldsymbol{\epsilon}_{\theta}^{(t)}(\mathbf{x}_t)$$

• Если мы явно подставим $f_{\theta}^{(t)}(\mathbf{x}_t)$ через $\epsilon_{\theta}^{(t)}(\mathbf{x}_t)$, то получится более удобная форма:

$$\mathbf{x}_{t-1} = \sqrt{A_{t-1}} \left(\frac{\mathbf{x}_t - \sqrt{1 - A_t} \, \boldsymbol{\epsilon}_{\boldsymbol{\theta}}^{(t)}(\mathbf{x}_t)}{\sqrt{A_t}} \right) + \sqrt{1 - A_{t-1} - \sigma_t^2} \, \boldsymbol{\epsilon}_{\boldsymbol{\theta}}^{(t)}(\mathbf{x}_t) + \sigma_t$$

где $\epsilon_t \sim \mathcal{N}(\mathbf{0}, \mathbf{I})$ – новый случайный шум

• Часто это записывают как

$$\mathbf{x}_{t-1} = \underbrace{\sqrt{A_{t-1}} \widehat{\mathbf{x}}_0^{(t)}}_{\text{проекция предсказанного } \mathbf{x}_0} + \underbrace{\sqrt{1 - A_{t-1} - \sigma_t^2}}_{\text{направление к текущему } \mathbf{x}_t} + \underbrace{\sigma_t \, \epsilon_t}_{\text{случайный шум}},$$

где
$$\hat{\mathbf{x}}_0^{(t)} = f_{\theta}^{(t)}(\mathbf{x}_t).$$

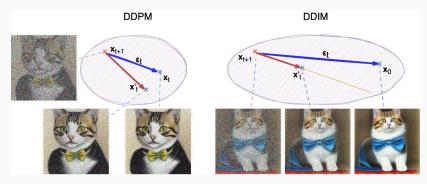
- Получается, что первый член возвращает нас к предсказанной версии \mathbf{x}_0 , второй корректирует направление в пространстве к текущему \mathbf{x}_t , третий добавляет случайность
- В DDIM можно взять $\sigma_t = 0$ и получить детерминированный обратный процесс; это тоже очень важное преимущество

6

- Но главное, конечно, в том, что теперь можно перепрыгивать через некоторые шаги!
- В DDIM можно рассмотреть подмножество шагов $au = \{ au_1, \dots, au_S\} \subset \{1, \dots, T\}$, по которым мы действительно делаем обратные переходы, и делать всё то же, но по au
- Интуитивно говоря, мы сохраняем маргинал $q(\mathbf{x}_t \mid \mathbf{x}_0)$ для всех t, но обратный процесс будем реализовывать только на подмножестве au, перепрыгивая через остальные шаги
- Таким образом можно тысячу шагов превратить в примерно пятьдесят

DDIM

• Иллюстрация:



• Сложно сказать, какое приближение лучше, но теперь мы разделили зависимости ϵ_t от \mathbf{x}_0 , и мы можем перепрыгнуть на несколько шагов сразу, двигаясь от \mathbf{x}_t к \mathbf{x}_{t+k} за дин шаг с увеличенным $\epsilon!$

7

DDIM

• Это привело к 10x–100x ускорению по сравнению с DDPM без потери качества:

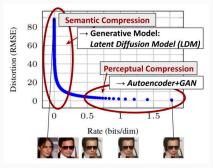


DDIM

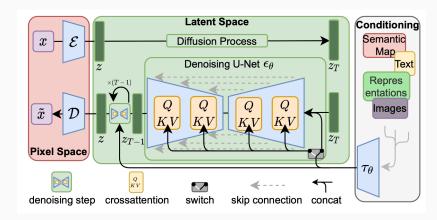
- Порождение в DDIMs тоже не обязано быть стохастическим; можно дисперсию обратной диффузии установить в ноль при порождении, $\sigma \to 0$
- И теперь латентный код в пространстве \mathbf{x}_T соответствует ровно одной картинке, а DDIM хорошая модель для латентных представлений; вот, например, интерполяции в латентном пространстве:



- Ещё один прорыв пришёл тогда, когда стали использовать этот диффузионный процесс в латентном пространстве
- · Stable Diffusion (Rombach, Blattmann et al., 2022):



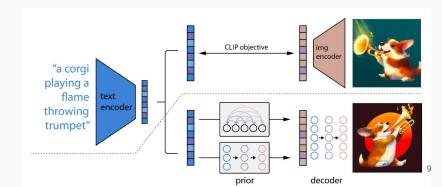
· Stable Diffusion (Rombach, Blattmann et al., 2022):



· unCLIP, он же DALL-E 2 (Ramesh et al., 2022), обучает

$$p\left(\mathbf{x}|y\right) = p\left(\mathbf{x}|\mathbf{c},y\right)p\left(\mathbf{c}|y\right),$$

где prior model $p\left(\mathbf{c}|y\right)$ строит CLIP embedding по тексту, а в $p\left(\mathbf{x}|\mathbf{c},y\right)$ сначала порождается «шум» (prior) для картинки, а потом по ней диффузионная модель рисует саму картинку:



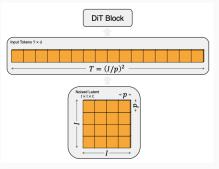
DALL-E 2



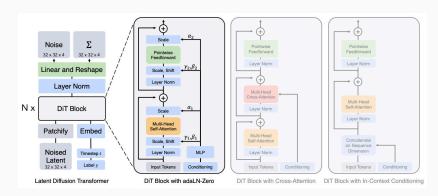
MIDJOURNEY



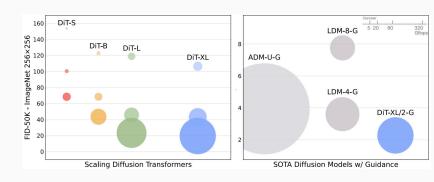
- Уже вполне современная архитектура DiT (Peebles, Xie, 2022)
- Это модель латентной диффузии, но вместо U-Net-подобной архитектуры для латентных кодов используют трансформер
- Вход трансформера тензор $I \times I \times C$, "patchified" в последовательность $p \times p$ патчей длины $T = (I/p)^2$; p гиперпараметр, квадратично влияющий на сложность



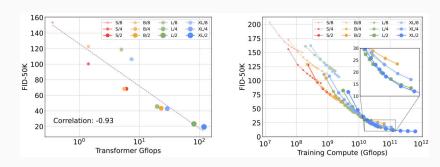
• Затем его обрабатывают блоками DiT; интересная часть здесь – adaptive layer norm с параметрами, взятыми из условия



• Получается, что DiT куда более вычислительно эффективны:

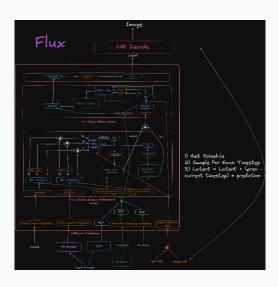


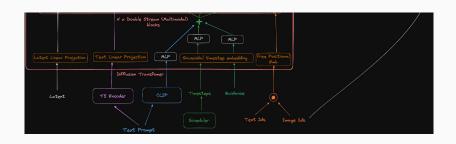
• Transformer GFlops сильно скоррелированы с качеством, а большие модели DiT используют вычисления более эффективно:

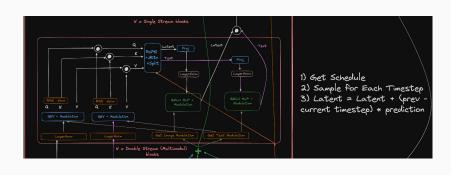


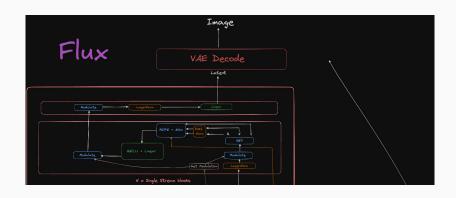
• Хорошие сэмплы:































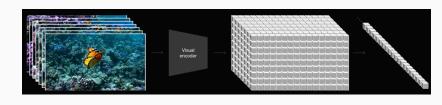
OPENAL SORA

• Текущее применение DiT: OpenAl Sora (Feb 2024, и вот ещё совсем недавно)



OPENAI SORA

· OpenAl всех деталей не раскрывает, но там точно есть DiT:





От дискретных процессов к непрерывным

ОТ DDIM к непрерывному времени

- Начнём с идеи DDIM: вместо шага $\mathbf{x}_t o \mathbf{x}_{t-1}$ мы приближаем направление сразу к \mathbf{x}_0 .
- Пусть

$$\alpha_t = 1 - \beta_t, \quad A_t = \prod_{i=1}^t \alpha_i, \quad \hat{\mathbf{x}}_0(\mathbf{x}_t, t) = \frac{\mathbf{x}_t - \sqrt{1 - A_t} \, \epsilon_\theta(\mathbf{x}_t, t)}{\sqrt{A_t}}.$$

· Один шаг DDIM с параметром $\eta=0$:

$$\mathbf{x}_{t-1} = \sqrt{A_{t-1}}\,\hat{\mathbf{x}}_0(x_t,t) + \sqrt{1-A_{t-1}}\,\hat{\epsilon}(\mathbf{x}_t,t).$$

• Интуитивно это значит, что мы движемся по «детерминированной» траектории от \mathbf{x}_t к \mathbf{x}_0 .

ОТ DDIM к непрерывному времени

- Дальше переходим к непрерывному времени: время $t \in [0,T]$, шаги $\Delta t \to 0$.
- Будем описывать прямую диффузию как стохастическое дифференциальное уравнение (SDE); обратный процесс как *обратную* SDE.
- Зашумление это стохастическое дифференциальное уравнение

$$\mathrm{d}\mathbf{x} = \mathbf{f}(\mathbf{x},t)\,\mathrm{d}t + g(t)\,\mathrm{d}\mathbf{w}_t, \quad \mathbf{f}(\mathbf{x},t) = -\tfrac{1}{2}\beta(t)\,\mathbf{x}, \ g(t) = \sqrt{\beta(t)},$$

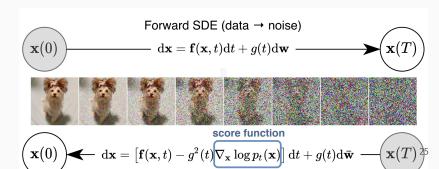
где \mathbf{w}_t – винеровский процесс (броуновское движение), \mathbf{f} – детерминированный дрейф, g – интенсивность шума; этот процесс постепенно превращает данные в шум (обычно стандартный гауссиан); пусть $p_t(\mathbf{x})$ – это плотность $\mathbf{x}(t)$

ОТ DDIM к непрерывному времени

• Самый главный здесь результат — Anderson (1982): для такого процесса обратный тоже является диффузионным процессом и задаётся уравнением

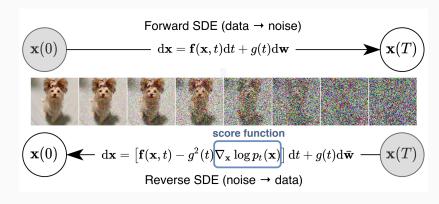
$$\mathrm{d}\mathbf{x} = \left[\mathbf{f}(\mathbf{x},t) - g(t)^2 \, \nabla_{\mathbf{x}} \log p_t(\mathbf{x})\right] \mathrm{d}t + g(t) \, \mathrm{d}\bar{\mathbf{w}},$$

где $\bar{\mathbf{w}}$ – тоже стандартный винеровский процесс, но для времени от T до 0, а $\mathrm{d}t$ теперь отрицательное



ОТ DDIM к непрерывному времени

- Функцию $s_*(\mathbf{x},t) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ в этой науке называют score; если score известен для каждого t, то процесс можно обратить
- Так что наш вопрос сводится к тому, откуда взять score

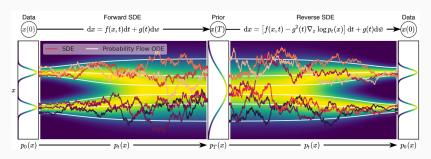


OT DDIM к непрерывному времени

• A score, конечно же, можно и нужно обучать через score matching: ищем s_{θ} с параметрами θ , минимизируя

$$\mathbb{E}_t \left[\lambda(t) \mathbb{E}_{\mathbf{x}(0)} \left[\mathbb{E}_{\mathbf{x}(t) \mid \mathbf{x}(0)} \left[\left\| s_{\theta}(\mathbf{x}(t), t) - \nabla_{\mathbf{x}(t)} \log p_t(\mathbf{x}(t) | \mathbf{x}(0)) \right\|_2^2 \right] \right] \right],$$

где λ – веса, а $p_t(\mathbf{x}(t)|\mathbf{x}(0))$ обычно можно так или иначе подсчитать (не будем уж в подробностях)



• И теперь DDPM – это дискретизация такого процесса: наша марковская цепь $\mathbf{x}_t = \sqrt{1-\beta_t}\mathbf{x}_{t-1} + \sqrt{\beta_t}\epsilon$ сходится в пределе к стохастическому дифференциальному уравнению

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x} dt + \sqrt{\beta(t)} d\mathbf{w},$$

т.е. как раз к такому процессу, который мы рассматривали

- Там есть два варианта: variance preserving (VP) как выше (дисперсия сохраняется единичной) и variance exploding (VE), в котором она всегда растёт, но это нам сейчас не важно
- Важно, что DDIM это просто более грубая дискретизация того же процесса!

• И даже более того – удивительно, но для любого диффузионного процесса есть и детерминированный процесс с теми же маргиналами $p_t(\mathbf{x})!$ Для $s(\mathbf{x},t) = \nabla_{\mathbf{x}} \log p_t(\mathbf{x})$ это

$$d\mathbf{x} = \left[\mathbf{f}(\mathbf{x}, t) - \frac{1}{2}g(t)^2 s(\mathbf{x}, t)\right] dt$$

· Например, для DDPM как выше это будет

$$d\mathbf{x} = -\frac{1}{2}\beta(t)\mathbf{x} - \frac{1}{2}\beta(t)s(\mathbf{x}, t)$$

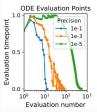
• Интуитивно говоря, в SDE вероятностная масса течёт и размазывается, а в таком процессе течёт без шума, но так, что плотность p_t по времени остаётся той же

• И если теперь заменить s_* на s_{θ} , мы получим практическую формулу для порождения:

$$\frac{\mathrm{d}\mathbf{x}}{\mathrm{d}t} \approx -\frac{1}{2}\beta(t)\Big(\mathbf{x} + s_{\theta}(\mathbf{x}, t)\Big).$$

- Это и есть непрерывная версия DDIM: движение к \mathbf{x}_0 задаётся ОДУ, а не стохастическим процессом

- · Как обучать на практике: берём \mathbf{x}_0 из данных, получаем \mathbf{x}_t из замкнутой формы распределения, минимизируем квадрат отклонения $s_{\theta}(\mathbf{x},t)$ от истинного $\log p_t(\mathbf{x}(t)|\mathbf{x}(0))$
- Благодаря тому, что процесс можно считать детерминированным, появляются полезные эффекты:
 - удобные латентные коды и интерполяции;
 - пропуски шагов, как в DDIM: можно брать разреженную сетку по t и ускорять процесс;
 - траектория к \mathbf{x}_0 «более прямая», чем в стохастической обратной SDE.







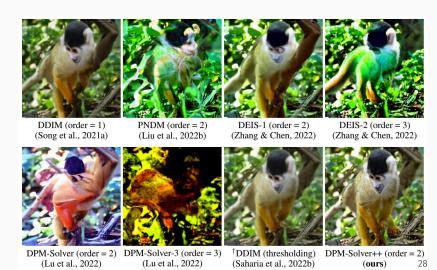
• А ещё, кстати, всё это работает и для управляемого порождения: если нам известно $p_t\left(\mathbf{y}|\mathbf{x}(t)\right)$ для некоторого условия \mathbf{y} (например, метки класса или частично заданного изображения), то мы можем так же сэмплировать из $p_0\left(\mathbf{x}(0)|y\right)$



Что дальше?

- Что произошло после DiT в диффузионных моделях?
- · DPM-Solver++ (Lu et al., 2022): как дальше ускорить DDIM?
- Как мы видели выше, DDIM это шаг первого порядка для интегрирования некоторого ОДУ
- Но ведь есть большая наука о том, как решать ОДУ лучше!
- Давайте выберем какую-нибудь многошаговую схему второго или третьего порядка, она будет понижать локальную ошибку, и в итоге можно будет шагать ещё шире, и шагов понадобится ещё меньше

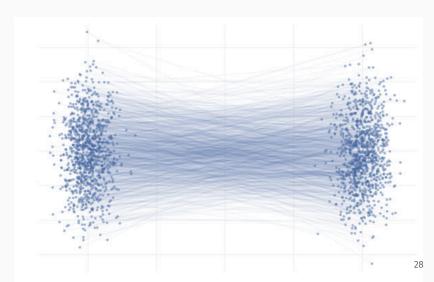
• Например, DPM-Solver++ действительно стабильнее и лучше работает, чем другие варианты (о которых не будем):



- UniPC (Zhao et al., 2023): метод типа «предиктор-корректор»
- Выводят универсальный корректирующий шаг (UniC), который можно добавить после любого базового предиктора, чтобы увеличить порядок точности без дополнительных вызовов модели, и выводят сопряжённый UniP-предиктор



• A самое интересное дальнейшее развитие событий – это, конечно, flow matching:



• На нём работают многие современные порождающие модели:



Figure 1: Protein generated by RFDiffusion (Watson et al., 2023).



Figure 2: Image from DALL-E 3 (Betker et al., 2023).

• Но это уже совсем другая история...

Спасибо!

Спасибо за внимание!



