Сергей Николенко СПбГУ— Санкт-Петербург 17 ноября 2025 г.



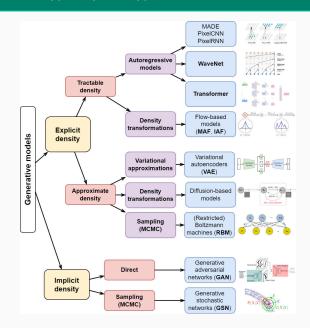


Random facts:

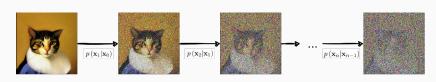
- 17 ноября Международный день студентов (в память об аресте >1200 пражских студентов нацистами 17 ноября 1939 года); кроме того, 17 ноября — Международный день недоношенных детей, а в России — День участковых уполномоченных полиции
- 17 ноября 1875 г. Елена Блаватская и Генри Стил Олкотт основали в Нью-Йорке
 Теософское общество, в которое входили Томас Эдисон, отец Джавахарлала Неру
 Мотилал, Уильям Батлер Йейтс и другие; общество пришло в упадок после смерти
 Блаватской в 1891 г., но созданный в 1907 г. российский филиал процветал до 1918-го
- 17 ноября 1959 г. южноафриканская компания De Beers начала производство искусственных алмазов
- 17 ноября 1853 г. произошло первое в истории сражение пароходофрегатов: русский «Владимир» победил турецкий «Перваз-Бахри»
- 17 ноября 2018 г. американский миссионер Джон Чау погиб от рук аборигенов Андаманских островов; Чау считал их последним оплотом Сатаны на Земле, но вскоре после высадки был застрелен из лука, посмертно получив за это премию Дарвина
- · 17 ноября 2019 г. в Китае выявили первого инфицированного COVID-19

Порождающие модели в глубоком обучении

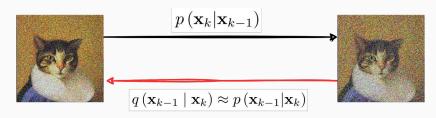
Глубокие порождающие модели



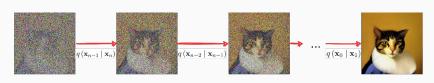
- (Sohl-Dickstein et al., 2015): Deep Unsupervised Learning using Nonequilibrium Thermodynamics
- Идея из статистической физики: давайте обучать марковскую цепь, которая постепенно сделает из случайного шума нужное распределение
- Для этого начнём с цепи, которая сделает из нужного распределения случайный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



• А потом попробуем обратить каждый шаг:



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

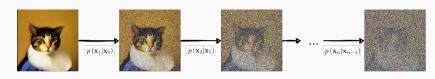


• Forward diffusion: начнём с $p\left(\mathbf{x}_{0}\right)$ и будем добавлять гауссовский шум:

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

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

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



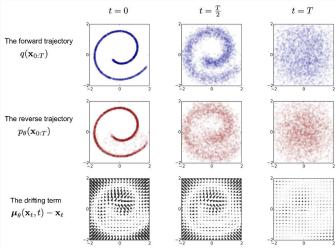
• И для каждого t мы можем просэмплировать \mathbf{x}_t в замкнутой форме через reparametrization trick: пусть $\epsilon \sim N\left(0,\mathbf{I}\right)$ и обозначим $\alpha_t = 1 - \beta_t$, $A_t = \prod_{i=1}^T \alpha_i$; тогда

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

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 β_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

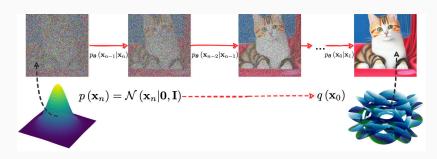
$$q\left(\mathbf{x}_{t-1}\mid\mathbf{x}_{t},\mathbf{x}_{0}
ight)=N\left(\mathbf{x}_{t-1}\middle|\tilde{\mu}\left(\mathbf{x}_{t},\mathbf{x}_{0}
ight),\tilde{eta}_{t}\mathbf{I}
ight),$$
 где

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

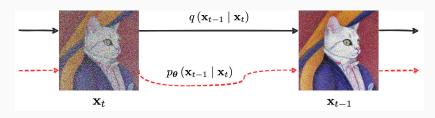
и поскольку мы знаем, что $\mathbf{x}_t = \sqrt{A_t}\mathbf{x}_0 + \sqrt{1-A_t}\epsilon$,

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

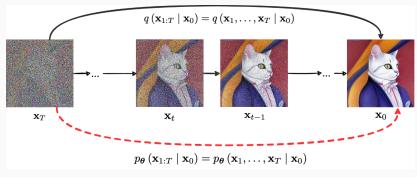
• В конечном счёте хочется, чтобы обратная диффузия реконструировала $q\left(\mathbf{x}_{0}\right)$ из стандартного входа в \mathbf{x}_{n} ; примерно так:



- 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 $q\left(\mathbf{z}\right)$:

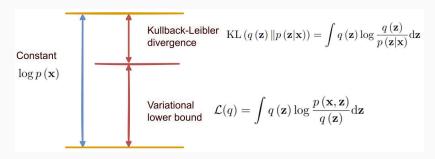
$$\mathbb{E}_{q(\mathbf{z})}\left[\log p\left(\mathbf{x}\right)\right] = \mathbb{E}_{q(\mathbf{z})}\left[\log p\left(\mathbf{x},\mathbf{z}\right)\right] - \mathbb{E}_{q(\mathbf{z})}\left[\log p\left(\mathbf{z}|\mathbf{x}\right)\right],$$

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} \mid \mathbf{x}_{t-1}, \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\right)} \frac{q\left(\mathbf{x}_{t} \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} \mid \mathbf{x}_{t-1}, \mathbf{x}_{0}\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)}{q\left(\mathbf{x}_{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} \mid \mathbf{x}_{t-1}, \mathbf{x}_{0}\right)}{p_{\theta}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}\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

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

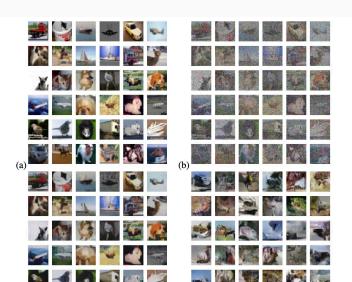
и пытаемся её параметры совместить с $q\left(\mathbf{x}_{t} \mid \mathbf{x}_{t-1}, \mathbf{x}_{0}\right)$. Для среднего, например,

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

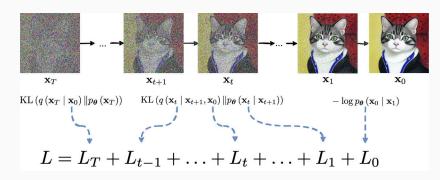
и поскольку мы знаем \mathbf{x}_t при обучении, мы можем параметризовать шум напрямую:

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

• (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:



- · 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

- 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:



 Для DDIM мы определим диффузионный процесс в терминах его обратной формы:

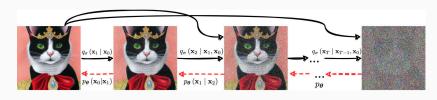
$$q_{\sigma}\left(\mathbf{x}_{1:T} \mid \mathbf{x}_{0}\right) = q_{\sigma}\left(\mathbf{x}_{T} \mid \mathbf{x}_{0}\right) \prod_{t=2}^{T} q_{\sigma}\left(\mathbf{x}_{t-1} \mid \mathbf{x}_{t}, \mathbf{x}_{0}\right).$$

• Теперь можно выразить распределения прямой диффузии по теореме Байеса:

$$q_{\sigma}\left(\mathbf{x}_{t}\mid\mathbf{x}_{t-1},\mathbf{x}_{0}\right) = \frac{q_{\sigma}\left(\mathbf{x}_{t-1}\mid\mathbf{x}_{t},\mathbf{x}_{0}\right)q_{\sigma}\left(\mathbf{x}_{t}\mid\mathbf{x}_{0}\right)}{q_{\sigma}\left(\mathbf{x}_{t-1}\mid\mathbf{x}_{0}\right)}.$$

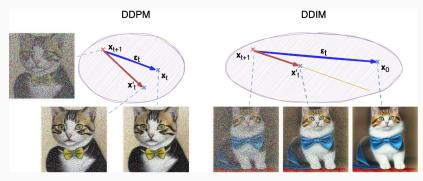
7

• И теперь можно показать, что у полученного процесса ровно те же маргиналы, так что обратную диффузию можно обучать с той же функцией ошибки, и она будет марковской цепью:



- Пока это не очень помогает: мы расширили класс распределений для прямой диффузии, но сэмплирование новой картинки всё равно должно пройти через все шаги обратной диффузии
- Но теперь, вместо того чтобы приближать случайный шум ϵ_t , который ведёт от \mathbf{x}_t к \mathbf{x}_{t+1} , можно приближать случайный шум ϵ_t , который приведёт сразу от \mathbf{x}_0 к \mathbf{x}_{t+1} !
- И когда мы идём в обратном направлении, мы приближаем направление не к \mathbf{x}_t , а сразу к \mathbf{x}_0 , и двигаться в эту сторону можно быстрее

• 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 один шаг с увеличенным $\epsilon!$

7

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



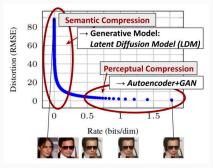
- Порождение в DDIMs тоже не обязано быть стохастическим; можно дисперсию обратной диффузии установить в ноль при порождении
- И теперь латентный код в пространстве \mathbf{x}_T соответствует ровно одной картинке, а DDIM хорошая модель для латентных представлений; вот, например, интерполяции в латентном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



STABLE DIFFUSION

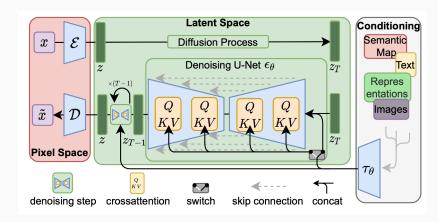
STABLE DIFFUSION

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



STABLE DIFFUSION

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

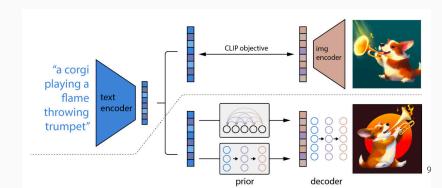


STABLE DIFFUSION

· 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:



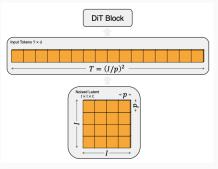
DALL-E 2



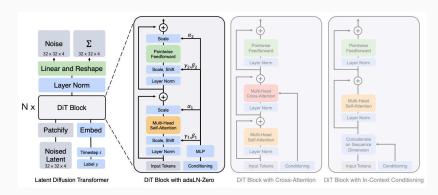
MIDJOURNEY



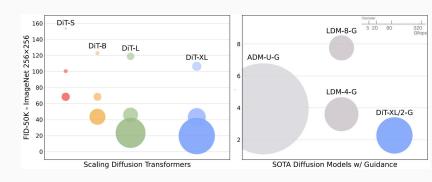
- Уже вполне современная архитектура DiT (Peebles, Xie, 2022)
- Это модель латентной диффузии, но вместо U-Net-подобной архитектуры для латентных кодов используют трансформер
- Вход трансформера тензор $I \times I \times C$, "patchified" в последовательность $p \times p$ патчей длины $T = (I/p)^2$; p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



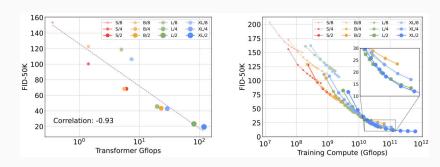
• Затем его обрабатывают блоками DiT; интересная часть здесь – adaptive layer norm с параметрами, взятыми из условия



· Получается, что DiT куда более вычислительно эффективны:

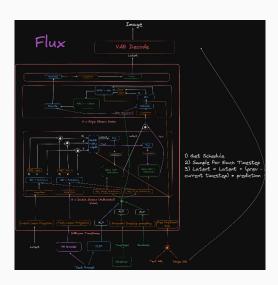


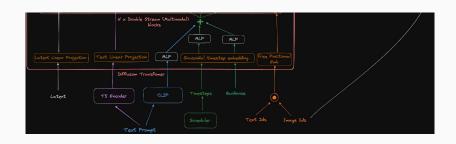
• Transformer GFlops сильно скоррелированы с качеством, а большие модели DiT используют вычисления более эффективно:

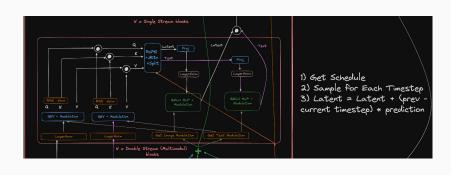


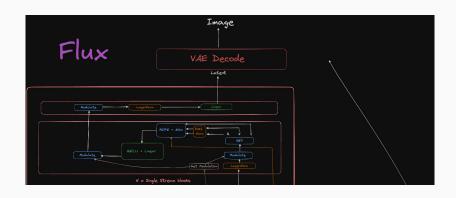
• Хорошие сэмплы:































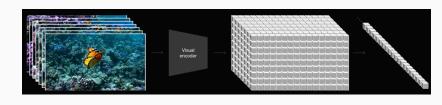
OPENAL SORA

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



OPENAI SORA

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





Спасибо!

Спасибо за внимание!



