

This document is a summary of the *Machine Perception* course at ETH Zürich. This summary was created during the spring semester of 2024. Due to updates to the syllabus content, some material may no longer be relevant for future versions of the course. I do not guarantee correctness or completeness, nor is this document endorsed by the lecturers. The order of the chapters is not necessarily the order in which they were presented in the course. For the full L<sup>A</sup>T<sub>E</sub>X source code, visit [github.com/Jovvik/eth-cheatsheets](https://github.com/Jovvik/eth-cheatsheets). All figures are created by the author, and, assuming the rules have not been changed, are not allowed to be used as a part of a cheat sheet during the exam.

## 1 CNN

$T$  is linear if  $T(\alpha \mathbf{u} + \beta \mathbf{v}) = \alpha T(\mathbf{u}) + \beta T(\mathbf{v})$ , invariant to  $f$  if  $T(f(\mathbf{u})) = T(\mathbf{u})$ , equivariant to  $f$  if  $T(f(\mathbf{u})) = f(T(\mathbf{u}))$ . Any linear shift-equivariant  $T$  can be written as a convolution. Convolution:  $I'(i, j) = \sum_{m=-k}^k \sum_{n=-k}^k K(-m, -n) I(m+i, n+j)$ . Correlation:  $I'(i, j) = \sum_{m=-k}^k \sum_{n=-k}^k K(m, n) I(m+i, n+j)$ . Conv. forward:  $z^{(l)} = w^{(l)} * z^{(l-1)} + b^{(l)} = \sum_m \sum_n w_{m,n} z_{i-m, j-n}^{(l-1)} + b^{(l)}$ . Backward inputs:  $\delta^{(l-1)} = \frac{\partial C}{\partial z_{i,j}^{(l-1)}} = \delta^{(l)} * \text{ROT}_{180}(w^{(l)})$ , backward kernel:  $\frac{\partial C}{\partial w_{m,n}^{(l)}} = \delta^{(l)} * \text{ROT}_{180}(z^{(l-1)})$ . Width or height after conv or pool:  $(\text{in} + 2 \cdot \text{pad} - \text{dil} \cdot (\text{kern} - 1) - 1) / \text{stride} + 1$ , rounded down. Channels = number of kernels.

1D conv as matmul:

$$\begin{bmatrix} k_1 & 0 & \dots & 0 \\ k_2 & k_1 & & \vdots \\ k_3 & k_2 & k_1 & 0 \\ 0 & k_3 & k_2 & 0 \\ \vdots & \vdots & \ddots & \vdots \end{bmatrix}$$

Backprop example (rotate K):

$$\begin{bmatrix} 0 & 0 & 1 & 1 & 0 \\ 0 & 1 & 0 & 1 & 1 \\ 1 & 1 & 1 & 1 & 0 \\ 1 & 0 & 1 & 1 & 0 \\ 1 & 1 & 0 & 1 & 1 \end{bmatrix} \mathbf{X} \rightarrow \begin{bmatrix} 1 & 0 & 1 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} \mathbf{K} \rightarrow \begin{bmatrix} 3 & 3 & 3 \\ 4 & 2 & 3 \\ 2 & 3 & 3 \end{bmatrix} \mathbf{Y} = \mathbf{X} * \mathbf{K} \rightarrow$$

$$\begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{[4]} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{[1]} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \mathbf{Y}' = \text{Pool}(\mathbf{Y}) \mid \frac{\partial E}{\partial \mathbf{Y}'} \rightarrow \frac{\partial E}{\partial \mathbf{Y}} \rightarrow$$

$$\begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} \rightarrow \frac{\partial E}{\partial \mathbf{K}} \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \frac{\partial E}{\partial \mathbf{V}}$$

$$\begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{[4]} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{[1]} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \mathbf{Y}' = \text{Pool}(\mathbf{Y}) \mid \frac{\partial E}{\partial \mathbf{Y}'} \rightarrow \frac{\partial E}{\partial \mathbf{Y}} \rightarrow$$

$$\begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} \rightarrow \frac{\partial E}{\partial \mathbf{K}} \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \frac{\partial E}{\partial \mathbf{V}}$$

$$\begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{[4]} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{[1]} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \mathbf{Y}' = \text{Pool}(\mathbf{Y}) \mid \frac{\partial E}{\partial \mathbf{Y}'} \rightarrow \frac{\partial E}{\partial \mathbf{Y}} \rightarrow$$

$$\begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} \rightarrow \frac{\partial E}{\partial \mathbf{K}} \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \frac{\partial E}{\partial \mathbf{V}}$$

$$\begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{[4]} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \xrightarrow{[1]} \begin{bmatrix} 0 & 0 & 0 \\ 1 & 0 & 0 \\ 0 & 0 & 0 \end{bmatrix} \rightarrow \mathbf{Y}' = \text{Pool}(\mathbf{Y}) \mid \frac{\partial E}{\partial \mathbf{Y}'} \rightarrow \frac{\partial E}{\partial \mathbf{Y}} \rightarrow$$

$$\begin{bmatrix} 1 & 0 & 1 \\ 1 & 1 & 1 \\ 0 & 1 & 0 \end{bmatrix} \rightarrow \frac{\partial E}{\partial \mathbf{K}} \rightarrow \begin{bmatrix} 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 \end{bmatrix} \rightarrow \frac{\partial E}{\partial \mathbf{V}}$$

output is copies of filter weighted by input, summed on overlaps.

## 2 RNN

**Vanilla RNN:**  $\hat{y}_t = \mathbf{W}_{hy} \mathbf{h}_t, \mathbf{h}_t = f(\mathbf{h}_{t-1}, \mathbf{x}_t, \mathbf{W})$ , usually  $\mathbf{h}_t = \tanh(\mathbf{W}_{hh} \mathbf{h}_{t-1} + \mathbf{W}_{xh} \mathbf{x}_t)$ .

**BPTT:**  $\frac{\partial L}{\partial \mathbf{W}} = \sum_t \frac{\partial L_t}{\partial \mathbf{W}}$ , treat unrolled model as multi-layer.  $\frac{\partial L_t}{\partial \mathbf{W}}$  has a term of  $\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k} = \prod_{i=k+1}^t \frac{\partial \mathbf{h}_i}{\partial \mathbf{h}_{i-1}} = \prod_{i=k+1}^t \mathbf{W}_{hh}^T \text{diag} f'(\mathbf{h}_{i-1})$ .

**Exploding/vanishing gradients:**  $\mathbf{h}_t = \mathbf{W}^T \mathbf{h}_1$ . If  $\mathbf{W}$  is diagonaliz.,  $\mathbf{W} = \mathbf{Q} \text{diag} \lambda \mathbf{Q}^T = \mathbf{Q} \mathbf{A} \mathbf{Q}^T, \mathbf{Q} \mathbf{Q}^T = \mathbf{I} \Rightarrow \mathbf{h}_t = (\mathbf{Q} \mathbf{A} \mathbf{Q}^T)^t \mathbf{h}_1 = (\mathbf{Q} (\text{diag} \lambda)^t \mathbf{Q}^T) \mathbf{h}_1 \Rightarrow \mathbf{h}_t$  becomes the dominant eigenvector of  $\mathbf{W}$ .  $\frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_k}$  has this issue.

Long-term contributions vanish, too sensitive to recent distrations. **Truncated BPTT:** take the sum only over the last  $\kappa$  steps. **Gradient clipping**  $\frac{\text{threshold}}{\|\nabla\|} \nabla$  fights exploding gradients.

**2.1 LSTM** We want constant error flow, not multiplied by  $W^t$ .

- Input gate: which values to write,
- forget gate: which values to reset,
- output gate: which values to read,
- gate: candidate values to write to state.

$$\begin{pmatrix} \mathbf{i} \\ \mathbf{f} \\ \mathbf{o} \\ \mathbf{g} \end{pmatrix} = \begin{pmatrix} \sigma \\ \sigma \\ \sigma \\ \tanh \end{pmatrix} \mathbf{W} \begin{pmatrix} \mathbf{x}_t \\ \mathbf{h}_{t-1} \end{pmatrix}$$

$$\mathbf{c}_t = \mathbf{f} \odot \mathbf{c}_{t-1} + \mathbf{i} \odot \mathbf{g}$$

$$\mathbf{h}_t = \mathbf{o} \odot \tanh(\mathbf{c}_t)$$

## 3 Generative modelling

Learn  $p_{\text{model}} \approx p_{\text{data}}$ , sample from  $p_{\text{model}}$ .

- Explicit density:
  - Approximate:
    - \* Variational: VAE, Diffusion
    - \* Markov Chain: Boltzmann machine
  - Tractable:
    - \* Autoregressive: FVSBN/NADE/MADE, Pixel(c/r)NN, WaveNet/TCN, Autor. Transf.,
    - \* Normalizing Flows
  - Implicit density:
    - Direct: Generative Adversarial Networks
    - MC: Generative Stochastic Networks
- Autoencoder:  $X \rightarrow Z \rightarrow X, g \circ f \approx \text{id}, f$  and  $g$  are NNs. Optimal linear autoencoder is PCA.

Undercomplete:  $|Z| < |X|$ , else overcomplete. Overcomp. is for denoising, inpainting. Latent space should be continuous and interpolable. Autoencoder spaces are neither, so they are only good for reconstruction.

## 4 Variational AutoEncoder (VAE)

Sample  $z$  from prior  $p_\theta(z)$ , to decode use conditional  $p_\theta(x | z)$  defined by a NN.

$D_{\text{KL}}(P \| Q) := \int_x p(x) \log \frac{p(x)}{q(x)} dx$ : KL divergence, measure similarity of prob. distr.

$D_{\text{KL}}(P \| Q) \neq D_{\text{KL}}(Q \| P), D_{\text{KL}}(P \| Q) \geq 0$

$z$  can also be categorical. Likelihood  $p_\theta(x) = \int_z p_\theta(x | z) p_\theta(z) dz$  is hard to maximize, let encoder NN be  $q_\phi(z | x)$ ,  $\log p_\theta(x^i) = \mathbb{E}_z [\log p_\theta(x^i | z)] - D_{\text{KL}}(q_\phi(z | x^i) \| p_\theta(z)) + D_{\text{KL}}(q_\phi(z | x^i) \| p_\theta(z | x^i))$ . Red is intractable, use  $\geq 0$  to ignore it; Orange is reconstruction loss, clusters similar samples; Purple makes posterior close to prior, adds cont. and interp.

Orange - Purple is ELBO, maximize it.

$x \xrightarrow{\text{enc}} \mu_{z|x}, \Sigma_{z|x} \xrightarrow{\text{sample}} z \xrightarrow{\text{dec}} \mu_{x|z}, \Sigma_{x|z} \xrightarrow{\text{sample}} \hat{x}$

Backprop through sample by reparametr.:  $z = \mu + \sigma \epsilon$ . For inference, use  $\mu$  directly. Disentanglement: features should correspond to distinct factors of variation. Can be done with semi-supervised learning by making  $z$  conditionally independent of given features  $y$ .

**4.1  $\beta$ -VAE**  $\max_{\theta, \phi} \mathbb{E}_x [\mathbb{E}_{z \sim q_\phi} \log p_\theta(x | z)]$  to disentangle s.t.  $D_{\text{KL}}(q_\phi(z | x) \| p_\theta(z)) < \delta$ , with KKT:  $\max \text{Orange} - \beta \text{Purple}$ .

## 5 Autoregressive generative models

Autoregression: use data from the same input variable at previous time steps

Discriminative:  $P(Y | X)$ , generative:  $P(X, Y)$ , maybe with  $Y$  missing. Sequence models are generative: from  $x_1 \dots x_{i+k}$  predict  $x_{i+k+1}$ .

Tabular approach:  $p(\mathbf{x}) = \prod_i p(x_i | \mathbf{x}_{<i})$ , needs  $2^{i-1}$  params. Independence assumption is too strong. Let  $p_{\theta_i}(x_i | \mathbf{x}_{<i}) = \text{Bern}(f_i(\mathbf{x}_{<i}))$ , where  $f_i$  is a NN. **Fully Visible Sigmoid Belief Networks:**  $f_i = \sigma(\alpha_0^{(i)} + \alpha^{(i)} \mathbf{x}_{<i}^T)$ , complexity  $n^2$ , but model is linear.

**Neural Autoregressive Density Estimator:** add hidden layer.  $\mathbf{h}_i = \sigma(\mathbf{b} + \mathbf{W}_{:, <i} \mathbf{x}_{<i})$ ,  $\hat{x}_i = \sigma(c_i + \mathbf{V}_{i, \cdot} \mathbf{h}_i)$ . Order of  $\mathbf{x}$  can be arbitrary but fixed. Train by max log-likelihood in  $O(TD)$ ,

can use 2nd order optimizers, can use **teacher forcing**: feed GT as previous output.

Extensions: Convolutional; Real-valued: conditionals by mixture of gaussians; Order-less and deep: one DNN predicts  $p(x_k | x_{i_1} \dots x_{i_j})$ .

**Masked Autoencoder Distribution Estimator:** mask out weights s.t. no information flows from  $x_d \dots$  to  $\hat{x}_d$ . Large hidden layers needed. Trains as fast as autoencoders, but sampling needs  $D$  forward passes.

**PixelRNN:** generate pixels from corner, dependency on previous pixels is by RNN (LSTM). **PixelCNN:** also from corner, but condition by CNN over context region (perceptive field)  $\Rightarrow$  parallelize. For conditionals use masked convolutions. Channels: model R from context, G from R + cont., B from G + R + cont. Training is parallel, but inference is sequential  $\Rightarrow$  slow. Use conv. stacks to mask correctly.

NLL is a natural metric for autoreg. models, hard to evaluate others.

**WaveNet:** audio is high-dimensional. Use dilated convolutions to increase perceptive field with multiple layers.

AR does not work for high res images/video, convert the images into a series of tokens with an AE: Vector-quantized VAE. The codebook

is a set of vectors.  $x \xrightarrow{\text{enc}} z \xrightarrow{\text{codebook}} z_q \xrightarrow{\text{dec}} \hat{x}$ .

We can run an AR model in the latent space.

**5.1 Attention**  $\mathbf{x}_t$  is a convex combination of the past steps, with access to all past steps. For  $X \in \mathbb{R}^{T \times D}$ :  $K = XW_K, V = XW_V, Q = XW_Q$ . Check pairwise similarity between query and keys via dot product: let attention weights be  $\alpha = \text{Softmax}(QK^T / \sqrt{D})$ ,  $\alpha \in \mathbb{R}^{1 \times T}$ . Adding mask  $M$  to avoid looking into the future:

$$X = \text{Softmax} \left( \frac{(XW_Q)(XW_K)^T}{\sqrt{D}} + M \right) (XW_V)$$

Multi-head attn. splits  $W$  into  $h$  heads, then concatenates them. Positional encoding injects information about the position of the token. Attn. is  $O(T^2 D)$ .

## 6 Normalizing Flows

VAs dont have a tractable likelihood, AR models have no latent space. Want both. Change of variable for  $x =$

$f(z): p_x(x) = p_z(f^{-1}(x)) \left| \det \frac{\partial f^{-1}(x)}{\partial x} \right| = p_z(f^{-1}(x)) \left| \det \frac{\partial f(z)}{\partial z} \right|^{-1}$ . Map  $Z \rightarrow X$  with a deterministic invertible  $f_\theta$ . This can be a NN, but computing the determinant is  $O(n^3)$ . If the Jacobian is triangular, the determinant is  $O(n)$ . To do this, add a coupling layer:

$\begin{pmatrix} y^A \\ y^B \end{pmatrix} = \begin{pmatrix} h(x^A, \beta(x^B)) \\ x^B \end{pmatrix}$ , where  $\beta$  is any model, and  $h$  is elementwise.

$$\begin{pmatrix} x^A \\ x^B \end{pmatrix} = \begin{pmatrix} h^{-1}(y^A, \beta(y^B)) \\ y^B \end{pmatrix}, J = \begin{pmatrix} h' & h' \beta' \\ 0 & 1 \end{pmatrix}$$

Stack these for expressivity,  $f = f_k \circ \dots \circ f_1$ .  
 $p_x(x) = p_z(f^{-1}(x)) \prod_k \left| \det \frac{\partial f_k^{-1}(x)}{\partial x} \right|$ .

Sample  $z \sim p_z$  and get  $x = f(z)$ .

• Squeeze: reshape, increase chan.

• ActNorm: batchnorm with init. s.t. output  $\sim \mathcal{N}(0, I)$  for first minibatch.  $y_{i,j} = s \odot x_{i,j} + b$ ,  $x_{i,j} = (y_{i,j} - b)/s$ ,  $\log \det = H \cdot W \cdot \sum_i \log |s_i|$ : linear.

•  $1 \times 1$  conv: permutation along channel dim. Init  $W$  as rand. orthogonal  $\in \mathbb{R}^{C \times C}$  with  $\det W = 1$ .  $\log \det = H \cdot W \cdot \log |\det W|$ :  $O(C^3)$ . Faster:  $W := PL(U + \text{diag}(s))$ , where  $P$  is a random fixed permut. matrix,  $L$  is lower triang. with 1s on diag.,  $U$  is upper triang. with 0s on diag.,  $s$  is a vector. Then  $\log \det = \sum_i \log |s_i|$ :  $O(C)$  Conditional coupling: add parameter  $w$  to  $\beta$ .

**SRFlow**: use flows to generate many high-res images from a low-res one. Adds affine injector between conv. and coupling layers.  $h^{n+1} = \exp(\beta_{\theta,s}^n(u)) \cdot h^n + \beta_{\theta,b}(u)$ ,  $h^n = \exp(-\beta_{\theta,s}^n(u)) \cdot (h^{n+1} - \beta_{\theta,b}^n(u))$ ,  $\log \det = \sum_{i,j,k} \beta_{\theta,s}^n(u_{i,j,k})$ .

**StyleFlow**: Take StyleGAN and replace the network  $z \rightarrow w$  (aux. latent space) with a normalizing flow conditioned on attributes.

**C-Flow**: condition on other normalizing flows: multimodal flows. Encode original image  $x_B^1: z_B^1 = f_\phi^{-1}(x_B^1 | x_A^1)$ ; encode extra info (image, segm. map, etc.)  $x_A^2: z_A^2 = g_\theta^{-1}(x_A^2)$ ; generate new image  $x_B^2: x_B^2 = f_\phi(z_B^1 | z_A^2)$ . Flows are expensive for training and low res. The latent distr. of a flow needn't be  $\mathcal{N}$ .

## 7 Generative Adversarial Networks (GANs)

Log-likelihood is not a good metric. We can have high likelihood with poor quality by mixing in noise and not losing much likelihood; or low likelihood with good quality by remembering input data and having sharp peaks there.

**Generator**  $G: \mathbb{R}^Q \rightarrow \mathbb{R}^D$  maps noise  $z$  to data, **discriminator**  $D: \mathbb{R}^D \rightarrow [0, 1]$  tries to decide if data is real or fake, receiving both gen. outputs and training data. Train  $D$  for  $k$  steps for each step of  $G$ .

Training GANs is a min-max process, which are hard to optimize.  $V(G, D) = \mathbb{E}_{x \sim p_d} \log(D(x)) + \mathbb{E}_{\hat{x} \sim p_m} \log(1 - D(\hat{x}))$

For  $G$  the opt.  $D^* = p_d(x)/(p_d(x) + p_m(x))$ . Jensen-Shannon divergence (symmetric):  $D_{JS}(p \| q) = \frac{1}{2} D_{KL}(p \| \frac{p+q}{2}) + \frac{1}{2} D_{KL}(q \| \frac{p+q}{2})$ . Global minimum of  $D_{JS}(p_d \| p_m)$  is the glob. min. of  $V(G, D)$ ,  $V(G, D^*) = -\log(4)$  and at optimum of  $V(D^*, G)$  we have  $p_d = p_m$ .

If  $G$  and  $D$  have enough capacity, at each update step  $D$  reaches  $D^*$  and  $p_m$  improves  $V(p_m, D^*) \propto \sup_D \int_x p_m(x) \log(-D(x)) dx$ , then  $p_m \rightarrow p_d$  by convexity of  $V(p_m, D^*)$  wrt.  $p_m$ . These assumptions are too strong.

If  $D$  is too strong,  $G$  has near zero gradients and doesn't learn ( $\log'(1 - D(G(z))) \approx 0$ ). Use gradient ascent on  $\log(D(G(z)))$  instead.

Model collapse:  $G$  only produces one sample or one class of samples. Solution: **unrolling** — use  $k$  previous  $D$  for each  $G$  update.

GANs are hard to compare, as likelihood is intractable. FID is a metric that calculates the distance between feature vectors calculated for real and generated images.

DCGAN: pool  $\rightarrow$  strided convolution, batchnorm, no FC, ReLU for  $G$ , LeakyReLU for  $D$ .

Wasserstein GAN: different loss, gradients don't vanish. Adding gradient penalty for  $D$  stabilizes training. Hierarchical GAN: generate low-res image, then high-res during training. StyleGAN: learn intermediate latent space  $\mathcal{W}$  with FCs, batchnorm with scale and mean from  $\mathcal{W}$ , add noise at each layer.

**GAN inversion**: find  $z$  s.t.  $G(z) \approx x \Rightarrow$  manipulate images in latent space, inpainting. If  $G$  predicts image and segmentation mask, we can use inversion to predict mask for any image, even outside the training distribution.

**7.1 3D GANs** 3D GAN: voxels instead of pixels. PlatonicGAN: 2D input, 3D output differentially rendered back to 2D for  $D$ .

HoloGAN: 3D GAN + 2D superresolution GAN

GRAF: radiance fields more effic. than voxels

GIRAFFE: GRAF + 2D conv. upscale

EG3D: use 3 2D images from StyleGAN for features, project each 3D point to tri-planes.

**7.2 Image Translation** E.g. sketch  $X \rightarrow$  image  $Y$ . Pix2Pix:  $G: X \rightarrow Y$ ,  $D: X, Y \rightarrow [0, 1]$ . GAN loss +  $L_1$  loss between sketch and image. Needs pairs for training.

CycleGAN: unpaired. Two GANs  $F: X \rightarrow Y$ ,  $G: Y \rightarrow X$ , cycle-consistency loss  $F \circ G \approx \text{id}$ ;  $G \circ F \approx \text{id}$  plus GAN losses for  $F$  and  $G$ .

BicycleGAN: add noise input.

Vid2vid: video translation.

## 8 Diffusion models

High quality generations, better diversity, more stable/scalable.

Diffusion (forward) step  $q$ : adds noise to  $x_t$  (not learned). Denoising (reverse) step  $p_\theta$ : removes noise from  $x_t$  (learned).

$$q(x_t | x_{t-1}) = \mathcal{N}(\sqrt{1 - \beta} x_{t-1}, \beta I)$$

$p_\theta(x_{t-1} | x_t) = \mathcal{N}(\mu_\theta(x_t, t), \sigma_\theta^2 I)$   
 $\beta_t$  is the variance schedule (monotone  $\uparrow$ ). Let  $\alpha_t := 1 - \beta_t$ ,  $\bar{\alpha}_t := \prod \alpha_i$ , then  $q(x_t | x_0) = \mathcal{N}(\sqrt{\bar{\alpha}_t} x_0, (1 - \bar{\alpha}_t) I) \Rightarrow x_t = \sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon$ . Denoising is not tractable naively:  $q(x_{t-1} | x_t) = q(x_t | x_{t-1}) q(x_{t-1}) / q(x_t)$ ,  $q(x_t) = \int q(x_t | x_0) q(x_0) dx_0$ .

Conditioning on  $x_0$  we get a Gaussian. Learn model  $p_\theta(x_{t-1} | x_t) \approx q(x_{t-1} | x_t, x_0)$  by predicting the mean.

$\log p(x_0) \geq \mathbb{E}_{q(x_{1:T} | x_0)} \log \left( \frac{p(x_{0:T})}{q(x_{1:T} | x_0)} \right) = \mathbb{E}_{q(x_1 | x_0)} \log p_\theta(x_0 | x_1) - D_{KL}(q(x_T | x_0) \| p(x_T)) - \sum_{t=2}^T \mathbb{E}_{q(x_t | x_0)} D_{KL}(q(x_{t-1} | x_t, x_0) \| p_\theta(x_{t-1} | x_t))$ , where orange and purple are the same as in VAEs, and blue are the extra loss functions. In a sense VAEs are 1-step diffusion models.

$t$ -th denoising is just  $\arg \min_\theta \frac{1}{2\sigma_\theta^2(t)} \|\mu_\theta - \mu_q\|_2^2$ , so we want  $\mu_\theta(x_t, t) \approx \mu_q(x_t, x_0)$ .  $\mu_q(x_t, x_0)$  can be written as  $\frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \epsilon_0$ , and  $\mu_\theta(x_t, t) = \frac{1}{\sqrt{\alpha_t}} x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t} \sqrt{\alpha_t}} \hat{\epsilon}_\theta(x_t, t)$ , so the NN learns to predict the added noise.

Training:  $\text{img } x_0, t \sim \text{Unif}(1 \dots T), \epsilon \sim \mathcal{N}(0, I)$ ,  $\text{GD on } \nabla_\theta \|\epsilon - \epsilon_\theta(\sqrt{\bar{\alpha}_t} x_0 + \sqrt{1 - \bar{\alpha}_t} \epsilon, t)\|^2$ .

Sampling:  $x_T \sim \mathcal{N}(0, I)$ , for  $t = T$  downto 1:  $z \sim \mathcal{N}(0, I)$  if  $t > 1$  else  $z = 0$ ;

$$x_{t-1} = \frac{1}{\sqrt{\alpha_t}} (x_t - \frac{1 - \alpha_t}{\sqrt{1 - \bar{\alpha}_t}} \epsilon_\theta(x_t, t)) + \sigma_t z.$$

$\sigma_t^2 = \beta_t$  in practice.  $t$  can be continuous.

**8.1 Conditional generation** Add input  $y$  to the model.

**ControlNet**: don't retrain model, add layers that add something to block outputs.

(Classifier-free) **guidance**: mix predictions of a conditional and unconditional model, because conditional models are not diverse.  $\eta_{\theta_1}(x, c; t) = (1 + \rho) \eta_{\theta_1}(x, c; t) - \rho \eta_{\theta_2}(x; t)$ .

**8.2 Latent diffusion models** High-res images are expensive to model. Predict in latent space, decode with a decoder.

## 9 Reinforcement learning

Environment is a Markov Decision Process: states  $S$ , actions  $A$ , reward  $r: S \times A \rightarrow \mathbb{R}$ , transition  $p: S \times A \rightarrow S$ , initial  $s_0 \in S$ , discount factor  $\gamma$ .  $r$  and  $p$  are deterministic, can be a distribution. Learn policy  $\pi: S \rightarrow A$ . Value  $V_\pi: S \rightarrow \mathbb{R}$ , the reward from  $s$  under  $\pi$ .

**Bellman eq.:**  $G_t := \sum_{k=0}^\infty \gamma^k R_{t+k+1}$ ,  $v_\pi(s) := \mathbb{E}_\pi[G_t | S_t = s] = \mathbb{E}_\pi[R_{t+1} + \gamma G_{t+1} | S_t = s] = \sum_a \pi(a | s) \sum_{s'} \sum_r p(s', r | s, a) [r + \gamma \mathbb{E}_\pi[G_{t+1} | S_{t+1} = s']] = \sum_a \pi(a | s) \sum_{s', r} p(s', r | s, a) [r + \gamma v_\pi(s')]$ . Can be solved via dynamic programming (needs knowledge of  $p$ ), Monte-Carlo or Temporal Difference learning.

**9.1 Dynamic programming** Value iteration: compute optimal  $v_*$ , then  $\pi_*$ .

Policy iteration: compute  $v_\pi$  and  $\pi$  together. For any  $V_\pi$  the greedy policy (optimal) is  $\pi'(s) = \arg \max_{a \in A} (r(s, a) + \gamma V_\pi(p(s, a)))$ .

**Bellman optimality:**  $v_*(s) = \max_a q_*(s, a) = \max_a \sum_{s', r} p(s', r | s, a) [r + \gamma v_*(s')] \Rightarrow$  update step:  $V_{\text{new}}^*(s) = \max_{a \in A} (r(s, a) + \gamma V_{\text{old}}^*(s'))$ , when  $V_{\text{old}}^* = V_{\text{new}}^*$ , we have optimal policy.

Converges in finite steps, more efficient than policy iteration. But needs knowledge of  $p$ , iterates over all states and  $O(|S|)$  memory.

**9.2 Monte Carlo sampling** Sample trajectories, estimate  $v_\pi$  by averaging returns. Doesn't need full  $p$ , is unbiased, but high variance, exploration/exploitation dilemma, may not reach term. state.



**9.3 Temporal Difference learning** For each  $s \rightarrow s'$  by action  $a$  update:  $\Delta V(s) = r(s, a) + \gamma V(s') - V(s)$ .  **$\epsilon$ -greedy policy:** with prob.  $\epsilon$  choose random action, else greedy.

**9.4 Q-learning**  $Q$ -value f.:  $q_\pi(s, a) := \mathbb{E}_\pi[G_t | S_t = s, A_t = a]$ .

**SARSA** (on-policy): For each  $S \rightarrow S'$  by action  $A$  update:  $\Delta Q(S, A) = r(S, A) + \gamma Q(S', A) - Q(S, A)$ ,  $Q(S, A) += \alpha \Delta Q(S, A)$ ,  $\alpha$  is LR.

**Q-learning** (off-policy/offline):  $\Delta Q(S, A) = R_{t+1} + \gamma \max_a Q(S', a) - Q(S, A)$   
All these approaches do not approximate values of states that have not been visited.

**9.5 Deep Q-learning** Use NN to predict  $Q$ -values. Loss is  $(R + \gamma \max_{a'} Q_\theta(S', a') - Q_\theta(S, A))^2$ , backprop only through  $Q_\theta(S, A)$ . Store history in replay buffer, sample from it for training  $\Rightarrow$  no correlation in samples.

**9.6 Deep Q-networks** Encode state to low dimensionality with NN.

**9.7 Policy gradients**  $Q$ -learning does not handle continuous action spaces. Learn a policy directly instead,  $\pi(a_t | s_t) = \mathcal{N}(\mu_t, \sigma_t^2 | s_t)$ . Sample trajectories:  $p(\tau) = p(s_1, a_1, \dots, s_T, a_T) = p(s_1) \prod \pi(a_t | s_t) p(s_{t+1} | a_t, s_t)$ . This is on-policy.

Eval:  $J(\theta) := \mathbb{E}_{\tau \sim p_\theta(\tau)} [\sum_t \gamma^t r(s_t, a_t)]$ . To optimize, need to compute  $\mathbb{E}$  (see proofs).

**REINFORCE:** MC sampling of  $\tau$ . To reduce variance, subtract baseline  $b(s_t)$  from reward.

**9.8 Actor-Critic**  $\nabla_\theta J(\theta) = \frac{1}{N} \sum_i \sum_t \nabla \log \pi_\theta(a_t^i | s_t^i) (r(s_t^i, a_t^i) + \gamma V(s_{t+1}^i) - V(s_t^i))$ .  $\pi$  = actor,  $V$  = critic. Est. value with NN, not traj. rollouts.

**9.9 Motion synthesis** **Data-driven:** bad perf. out of distribution, needs expensive mocap.

**DeepMimic:** RL to imitate reference motions while satisfying task objectives.

**SFV:** use pose estimation: videos  $\rightarrow$  train data.

## 10 Neural Implicit Representations

Voxels/volum. primitives are inefficient ( $n^3$  compl.). Meshes have limited granularity and have self-intersections. **Implicit representation:**  $S = \{x | f(x) = 0\}$ . Can be invertibly transformed without accuracy loss. Usually represented as signed distance function values on a grid, but this is again  $n^3$ . By UAT, approx.  $f$  with NN. **Occupancy networks:** predict probability that point is inside the shape.

**DeepSDF:** predict SDF. Both conditioned on input (2D image, class, etc.). Continuous, any topology/resolution, memory-efficient. NFs can model other properties (color, force, etc.).

**10.1 Learning 3D Implicit Shapes Inference:** to get a mesh, sample points, predict occupancy/SDF, use marching cubes.

**10.1.1 From watertight meshes** Sample points in space, compute GT occupancy/SDF, CE loss.

**10.1.2 From point clouds** Only have samples on the surface. Weak supervision: loss =  $|f_\theta(x_i)|^2 + \lambda \mathbb{E}_x (\|\nabla_x f_\theta(x)\| - 1)^2$ , edge points should have  $\|\nabla f\| \approx 1$  by def. of SDF,  $f \approx 0$ .

**10.1.3 From images** Need differentiable rendering 3D  $\rightarrow$  2D. **Differentiable Volumetric Rendering:** for a point conditioned on encoded image, predict occupancy  $f(x)$  and RGB color  $c(x)$ . **Forward:** for a pixel, ray-march and root find  $\hat{p} : f(\hat{p}) = 0$  with secant. Set pixel color to  $c(\hat{p})$ . **Backward:** see proofs.

## 10.2 Neural Radiance Fields (NeRF)

$(x, y, z, \theta, \phi) \xrightarrow{\text{NN}} (r, g, b, \sigma)$ . Density is predicted before adding view direction  $\theta, \phi$ , then one layer for color. **Forward:** shoot ray, sample points along it and blend:  $\alpha := 1 - \exp(-\sigma_i \delta_i)$ ,  $\delta_i := t_{i+1} - t_i$ ,  $T_i := \prod_{j=1}^{i-1} (1 - \alpha_j)$ , color is  $c = \sum_i T_i \alpha_i c_i$ . Optimized on many views of the scene. Can handle transparency/thin structure, but worse geometry. Needs many (50+) views for training, slow rendering for high res, only models static scenes.

**10.2.1 Positional Encoding for High Frequency Details** Replace  $x, y, z$  with pos. enc. or rand. Fourier feats. Adds high frequency feats.

**10.2.2 NeRF from sparse views** Regularize geometry and color.

**10.2.3 Fast NeRF render. and train.** Replace deep MLPs with learn. feature hash table + small MLP. For  $x$  interp. features between corners.

**10.3 3D Gaussian Splatting Alternative parametr.** Find a cover of object with primitives, predict inside. Or sphere clouds. Both ineff. for thin structures. Ellipsoids are better.

Initialize point cloud randomly or with an approx. reconstruction. Each point has a 3D Gaussian. Use camera params. to project ("splat") Gaussians to 2D and differentially render them. Adaptive density control

moves/clones/merges points.

Rasterization: for each pixel sort Gaussians by depth, opacity  $\alpha = o \cdot \exp(-0.5(x - \mu')^\top \Sigma'^{-1}(x - \mu'))$ , rest same as NeRF.

Each Gaussian primitive has a center  $\mu \in \mathbb{R}^3$ , a covariance  $\Sigma \in \mathbb{R}^{3 \times 3}$ , and a color  $c \in \mathbb{R}^3$  and an opacity  $o \in \mathbb{R}$ .

To model view-dependent color, the color can be replaced with spherical harmonics, i.e.  $c \in \mathbb{R}^9$

To keep covariance semi-definite:  $\Sigma = RSS^\top R^\top$ , where  $S \in \mathbb{R}^{3 \times 3}$  is a diagonal scaling matrix and  $R \in \mathbb{R}^{3 \times 3}$  is a rotation matrix.

## 11 Parametric body models

**11.1 Pictorial structure** Unary terms and pairwise terms between them with springs.

**11.2 Deep features** Direct regression: predict joint coordinates with refinement.

Heatmaps: predict probability for each pixel, maybe Gaussian. Can do stages sequentially.

**11.3 3D** Naive 2D  $\rightarrow$  3D lift works. But can't define constraints  $\Rightarrow$  2m arms sometimes.

**Skinned Multi-Person Linear model** (SMPL) is the standard non-commercial model. 3D mesh, base mesh is  $\sim 7k$  vertices, designed by an artist. To produce the model, point clouds (scans) need to be aligned with the mesh. **Shape deformation subspace:** for a set of human meshes T-posing, vectorize their vertices  $T$  and subtract the mean mesh. With PCA represent any person as weighted sum of 10-300 basis people,  $T = S\beta + \mu$ .

For pose, use **Linear Blend Skinning**.  $\mathbf{t}'_i = \sum_k w_{ki} \mathbf{G}_k(\theta, \mathbf{J}) \mathbf{t}_i$ , where  $\mathbf{t}$  is the T-pose positions of vertices,  $\mathbf{t}'$  is transformed,  $w$  are weights,  $\mathbf{G}_k$  is rigid bone transf.,  $\theta$  is pose,  $\mathbf{J}$  are joint positions. Linear assumption produces artifacts. **SMPL:**  $\mathbf{t}'_i = \sum_k w_{ki} \mathbf{G}_k(\theta, \mathbf{J}(\beta)) (\mathbf{t}_i + \mathbf{s}_i(\beta) + \mathbf{p}_i(\theta))$ . Adds shape correctives  $\mathbf{s}(\beta) = S\beta$ , pose cor.  $\mathbf{p}(\theta) = P\theta$ ,  $\mathbf{J}$  dep. on shape  $\beta$ . Predicting human pose is just predicting  $\beta, \theta$  and camera parameters.

**11.3.1 Optimization-based fitting** Predict 2D joint locations, fit SMPL to them by argmin with prior regularization. Argmin is hard to find, learn  $F: \Delta\theta = F(\frac{\partial L_{\text{reproj}}}{\partial \theta}, \theta^t, x)$ . Issues: self-occlusion, no depth info, non-rigid deformation (clothes).

**11.3.2 Template-based capture** Scan for first frame, then track with SMPL.

**11.3.3 Animatable Neural Implicit Surfaces** Model base shape and  $w$  with 2 NISs.

## 12 ML

Perceptron converges in finite time iff data is linearly separable. **MAP**  $\theta^* \in \arg \max_\theta p(\theta | X, y)$ . **MLE**  $\theta \in \arg \max_\theta p(y | X, \theta)$  consistent, efficient. **Binary cross-entropy**  $L(\theta) = -y_i \log(\hat{y}_i) - (1 - y_i) \log(1 - \hat{y}_i)$ . Cross-entropy  $H(p_d, p_m) = H(p_d) + D_{\text{KL}}(p_d \| p_m)$ . For any continuous  $f \exists \text{NN } g(x), |g(x) - f(x)| < \epsilon$ . 1 hidden layer is enough, activation function needs to be nonlinear.

MLP backward inputs:  $\delta^{(l)} = \delta^{(l+1)} \cdot \frac{\partial \mathbf{z}^{(l+1)}}{\partial \mathbf{z}^{(l)}}$ ,

backward weights:  $\left[ \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{w}_{ij}^{(l)}} \right]_k = f'(\mathbf{a})_k \cdot \mathbf{z}_j^{(l)}$ .

$[k = i], \frac{\partial^2 L}{\partial \mathbf{z}^{(l)} \partial \mathbf{z}^{(l)}} = \delta^{(l)} \frac{\partial \mathbf{z}^{(l)}}{\partial \mathbf{w}_{ij}^{(l)}}$ , backward bias:

$\frac{\partial L}{\partial \mathbf{b}_i^{(l)}} = \text{same, but no } \mathbf{z}$ .

## 12.1 Activation functions

name	$f(x)$	$f'(x)$	$f(X)$
sigmoid	$\frac{1}{1+e^{-x}}$	$\sigma(x)(1 - \sigma(x))$	$(0, 1)$
tanh	$\frac{e^x - e^{-x}}{e^x + e^{-x}}$	$1 - \tanh(x)^2$	$(-1, 1)$
ReLU	$\max(0, x)$	$[x \geq 0]$	$[0, \infty)$

Finite range: stable training, mapping to prob. space. Sigmoid, tanh saturate (value with large mod have small gradient)  $\Rightarrow$  vanishing gradient, Tanh is linear around 0 (easy learn), ReLU can blow up activation; piecewise linear  $\Rightarrow$  faster convergence.

**12.2 GD algos** **SGD:** use 1 sample. For sum structured loss is unbiased. High variance, efficient, jumps a lot  $\Rightarrow$  may get out of local min., may overshoot. **Mini-batch:** use  $m < n$  samples. More stable, parallelized. **Polyak's momentum:** velocity  $\mathbf{v} := \alpha \mathbf{v} - \epsilon \nabla_\theta L(\theta)$ ,  $\theta := \theta + \mathbf{v}$ . Move faster when high curv., consistent or noisy grad. **Nesterov's momentum:**  $\mathbf{v} := \alpha \mathbf{v} - \epsilon \nabla_\theta L(\theta + \alpha \mathbf{v})$ . Gets grad. at future point. **AdaGrad:**  $\mathbf{r} := \mathbf{r} + \nabla \odot \nabla$ ,  $\Delta\theta = -\epsilon / (\delta + \sqrt{\mathbf{r}}) \odot \nabla$ . Grads decrease fast for variables with high historical gradients, slow for low. But can decrease LR too early/fast. **RMSProp:**  $\mathbf{r} := \rho \mathbf{r} + (1 - \rho) \nabla \odot \nabla$ , use weighted moving average  $\Rightarrow$  drop history from distant past, works better for noncon-

