

ADL4CV: Lecture Notes

Philipp Wulff / Summer semester 2022

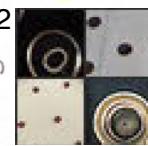
ADL4CV - 1. Visualization

(1)

- Visualize the input space: the receptive field

→ Extract the input image patches that result in the highest activation for a feature map pixel/unit

⇒ In later layers the input image patches are more specialized (objects instead of lines)



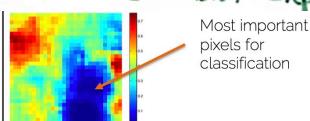
- Visualize importance (of features): the input space

Layer 5



→ Block different parts of the image and observe classification score change

Occlusion experiment ↗



Most important pixels for classification
⇒ Create a map (each pixel is the classification prob. if an occlusion square is in this region)

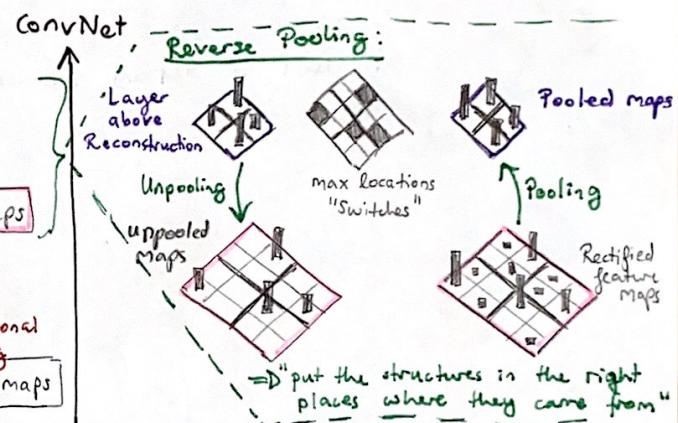
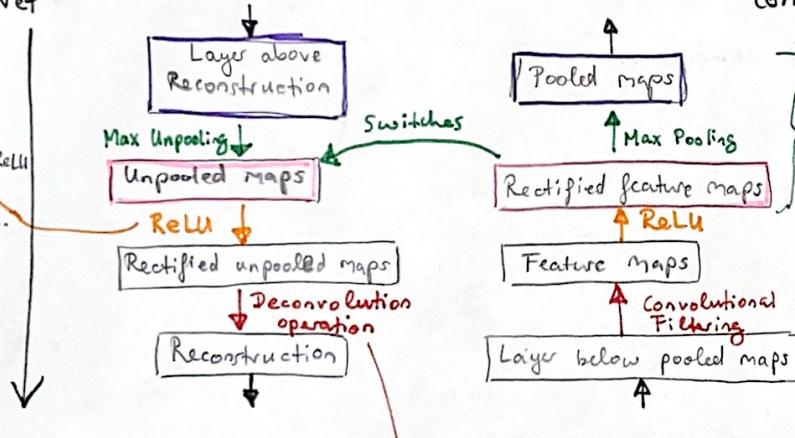
- Visualize features: Map activations back to the image space / Use DeconvNet to visualize features at a certain layer

① Ex: If filter 15 of the 3rd layer is highly activated

⇒ To visualize it: zero out other filters of layer 3 & pass through DeconvNet

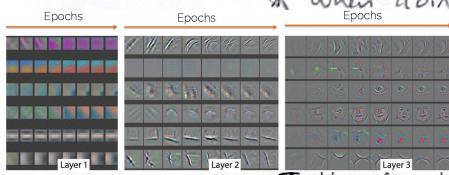
DeconvNet

We can use ReLU
as we also
visualize pos.
features



Convolve with the transpose of the learned filter!
⇒ "transposed" = vertically & horizontally flipped

⇒ print the reconstruction

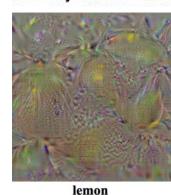


* When doing this for every filter of every layer ⇒ layer 1: low-level geom. features
⇒ deeper layers: reconstruct patterns/ faces/cats/...

After more epochs: Reconstructions become more detailed

* Visualization shows that large (e.g. 11x11) kernels produce inactive filters
(→ empty/blurry reconstructions)

② With gradient ascent: generate a synthetic image that max. activates a filter.



E.g. find image that max. score of some class

1. Take trained CNN (data was zero mean)

2. Pass empty zero image through CNN

3. Max. Score of class c with G.A. + Backprop

4. Update the image

* Improve through regularization (Gaussian blur, clip small values in image or gradients)

⇒ Now, max. activations of feature maps instead.

$$\operatorname{argmax}_I S_c(I) + \lambda \|I\|_2^2$$

L2-norm
to avoid
only a few
large pixels

③ Amplify the feature activations at some layer in the network (DeepDream)

(-> on a trained classifier)

1. Forward pass of an image up to layer L

2. Set gradient of layer L = activations

→ large gradients for whatever was detected

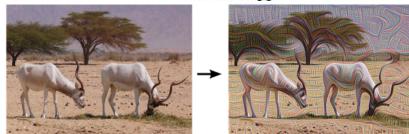
3. Backprop (-> calculate the gradient of the activations of layer L w.r.t. the input image pixels)

⇒ Show more "dogs"

4. Update image

⇒ Small changes to image according to a specific layer

shallow layer:



Deep layer:



Low layers: Amplify basic geometries

Deeper layers: Amplify objects

- Visualization with t-SNE:

→ Visualize FC-layers (e.g. last FC-layer in AlexNet has dim. 4096)

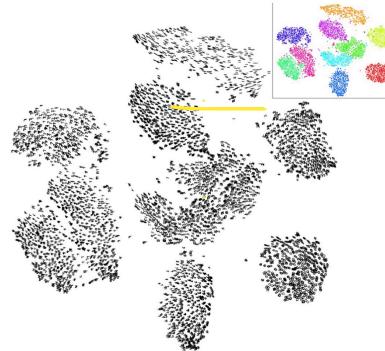
↳ show clusters of images in feature space with t-SNE

* Maps high-dim embedding to 2D Map which preserves pairwise distances

the points

⇒ Used to debug the network

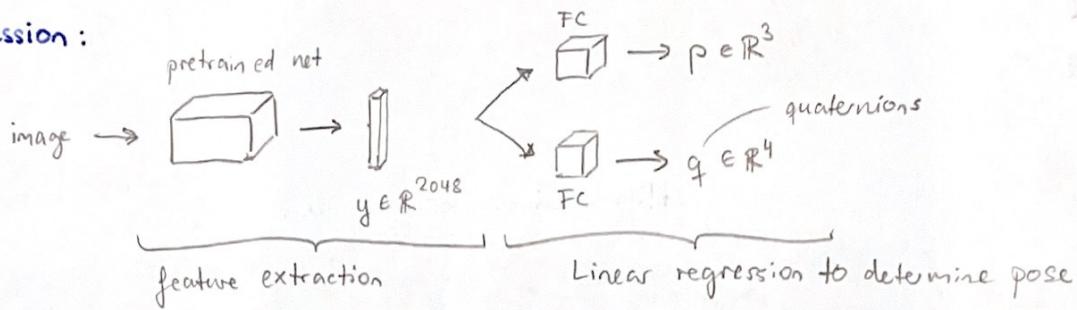
t-SNE Visualization: MNIST



ADL4CV - 2. Siamese Neural Networks & Similarity Learning

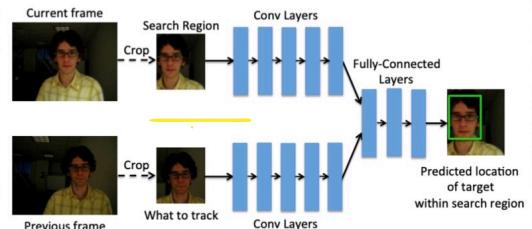
(2)

- Pose regression:



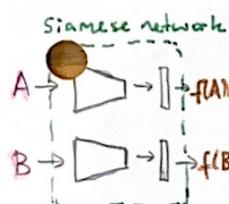
- Bounding box regression:

* Process current + previous frame with two feature extractors in parallel and achieve combined vector embedding → pass to FC-layers



- Similarity Learning: → Instead of discrete classification learn similarity function

→ Requires retraining whenever we add new samples



* Pass two images through Siamese networks (shared weights)

⇒ compare the encodings/embeddings

1) If A and B depict the same person :

$$\text{minimize } L(A, B) = \|f(A) - f(B)\|^2$$

2) If A and B depict different persons :

$$\text{minimize } L(A, B) = \max(0, m^2 - \|f(A) - f(B)\|^2)$$

$$y^* = \begin{cases} 1, & \text{if } A \text{ and } B \text{ depict the same person} \\ 0, & \text{otherwise} \end{cases}$$

combines 1) and 2)

$$\Rightarrow \text{Contrastive loss: } L(A, B) = y^* \|f(A) - f(B)\|^2 + (1 - y^*) \max(0, m^2 - \|f(A) - f(B)\|^2)$$

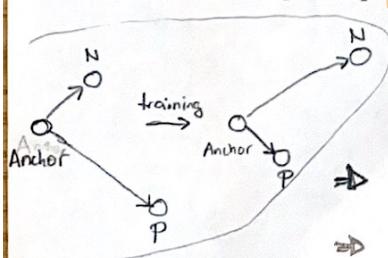
If A and B are already m^2 apart, do not pay anymore attention.

Or use the triplet loss:

$$\|f(A) - f(P)\|^2 \leq \|f(A) - f(N)\|^2$$

positive pair
(depicts same person as A)

negative pair



$$\|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + m \leq 0$$

$$\Rightarrow L(A, P, N) = \max(0, \|f(A) - f(P)\|^2 - \|f(A) - f(N)\|^2 + m)$$

Problems: - slow training when A and P are already similar

↳ Train only for a few epochs with the hard cases where $d(A, P) \approx d(A, N)$

↳ "Hard negative mining"
find from Nearest Neighbor search

* Random sampling does not work well → $O(n^3)$ possible triplets

⇒ triplet loss is more efficient & leads to better results for similarity learning

Since we learned a ranking → at test time we only need to do nearest neighbour search

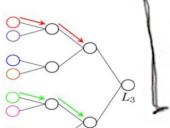
- Improving similarity learning: build a hierarchical tree for hard cases

Loss: * Use triplet loss

Sampling: * Choose best triplets in training (diversity of classes + hard cases)
 ↳ In total $O(n^3)$ combinations

→ Sampling method: Hierarchical triplet loss

1. Build the tree
 - Build tree where the leaves are the different image classes
 - Leaves are ordered s.t. its neighbor leaves are closest



$$d(p, q) = \frac{1}{n_p n_q} \sum_{i \in p} \sum_{j \in q} \|r_i - r_j\|^2$$

cardinality of classes p and q

feature embeddings of classes p and q

2. Sample diverse classes: Randomly select l^1 leaf nodes
3. For each of the l^1 chosen classes select the $m-1$ nearest neighbor
 → learn hard cases
4. From each class randomly pick t images
 ⇒ $t \cdot m \cdot l^1$ images per mini-batch
5. Construct all possible triplets from the mini-batch
6. Compute the triplet loss (the margin can be adaptable to the distances between classes on the hierarchical tree)

$$\mathcal{L}_M = \frac{1}{2Z_M} \sum_{T \in T^M} [\|x_a^z - x_p^z\| - \|x_a^z - x_n^z\| + \alpha_z]_+$$

all the triplets

"Divide & Conquer"

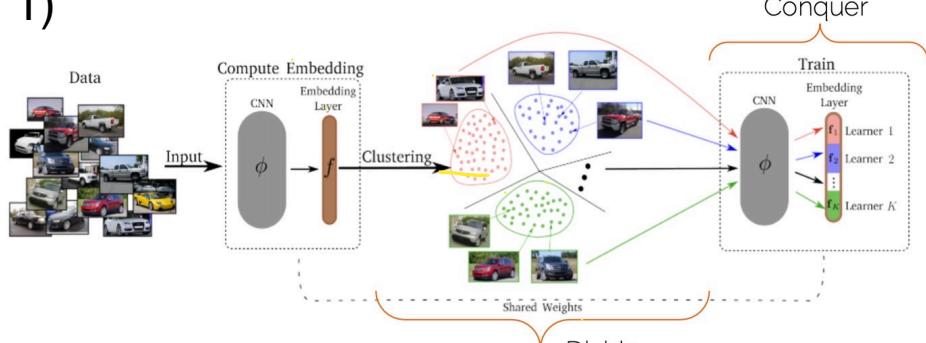
Ensembles: * Divide { 1) Cluster the embedding space into K clusters (e.g. with k-means)

2) Build k independent learners

3) Train each model on $\frac{\text{all}}{K}$ mini-batches from its cluster

Conquer { 4) After convergence, finetune using all models at the same time

5) Since we updated the network, the embedding may have changed
 → repeat from 1)

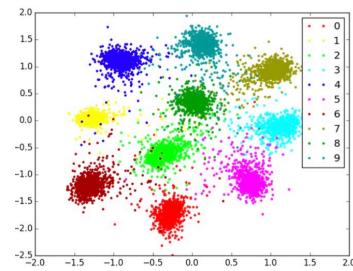


- Tips & Tricks:

- Stronger backbone (feature extractor / dense net) improves results
- Out-of-the-box: use ProxyNCA and SoftTripleLoss
- Even naive ensembles boost performance (using the same architecture)
- Sampling strategy and choice of loss function are important

- Applications in vision:

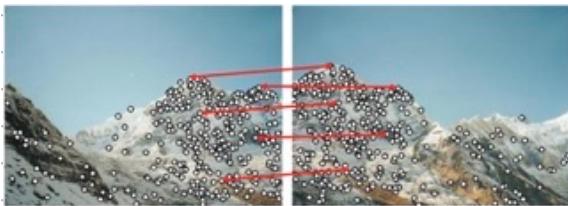
Classification: Siamese network on MNIST



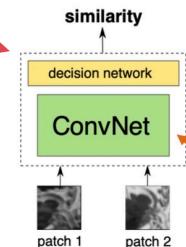
Establishing image correspondences

Classical method pipeline:

- Extract manually designed feature descriptors
 - Harris, SIFT, SURF: based on image gradients
 - suffer under illumination or viewpoint changes
 - Slow to extract dense features
- Match descriptors from the two images
 - Many descriptors are similar, one needs to filter out possible double matches and keep only reliable ones.

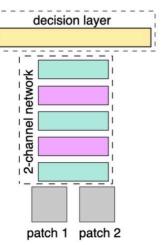


End-to-end learning for patch similarity



- Fast to allow dense extraction
- Invariant to a wide array of transformations (illumination, viewpoint)

- Classic Siamese architecture
 - Shared layers
 - Simulated feature extraction
 - One decision layer
 - Simulates the matching

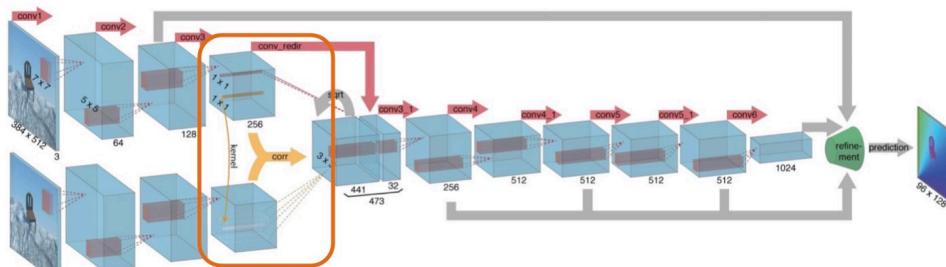


Optical Flow

- Input: 2 consecutive images (e.g. from a video)
- Output: per-pixel displacement from image A to B
- Result: “perceived” 2D motion (not the real motion)

-> e.g. FlowNet

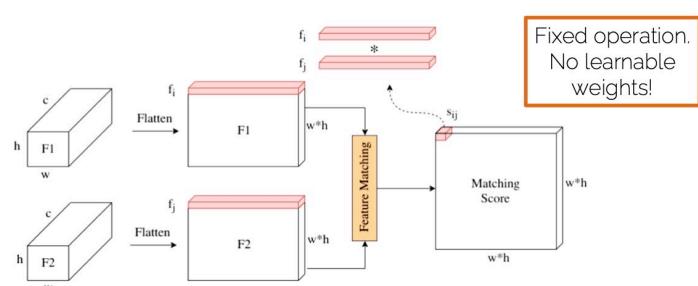
- Architecture 1: Stack 2 images
- Architecture 2: Siamese network



How to combine the information from both images?

Correlation Layer:

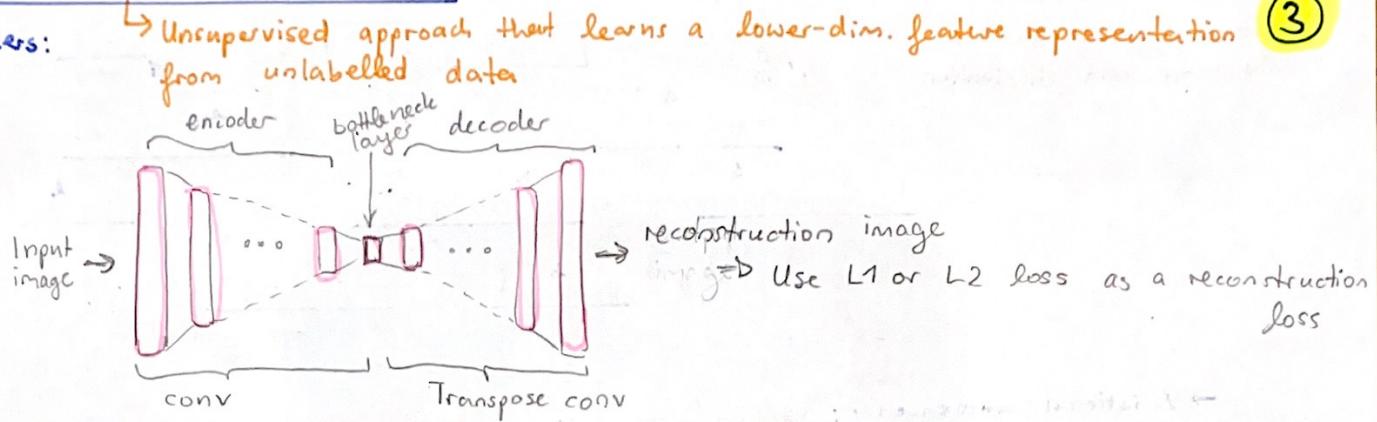
- Multiplies a feature vector with another feature vector
- The matching score represents how correlated these two feature vectors are



Fixed operation.
No learnable weights!

ADL4CV 3 - Autoencoders & VAE

- Autoencoders:



Use-cases:

- * Use the latent representation for clustering

* For automatic pretraining
↓
encoder as feature extractor

1. Pretrain autoencoder on large unlabelled dataset
2. Fine-tune the encoder with additional FC-layers on small labelled dataset

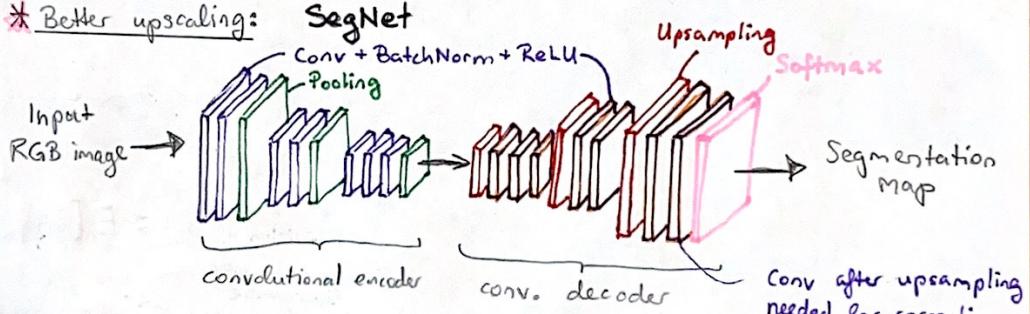
* For pixel-wise predictions:

→ Semantic segmentation:

* Previously used: fully-convolutional networks (FCN)

1. Encode image into latent representation
2. Upscale it to the original resolution with a feature map for each class

* Better upscaling: SegNet

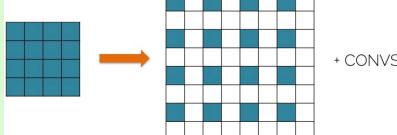


Upsampling Methods:

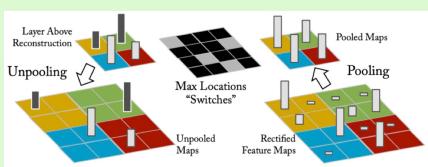
1) Interpolation: Nearest Neighbor/Lin./Bicubic
-> few artefacts



2) Fixed Unpooling (+ Conv to refine) -> efficient



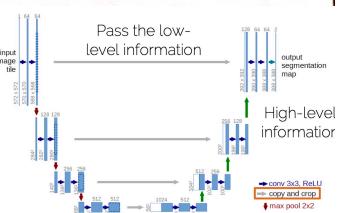
3) à la DeconvNet (save max. locations + Conv)
(unlike in DeconvNet, Conv-weights are learned)
-> keeps details of the structure



* The encoder compresses high-level info in the channels of the bottleneck layer

* Pass low-level info through skip-connections
⇒ crop feature maps from skip-connections such that dimensions align ⇒ concatenate

⇒ Results in sharp boundaries for segmentation



⇒ Unsupervised monocular depth estimation

1. Predict inverse depth map from left image

2. Use inverse warping to reconstruct the left image from right image + inv. depth map (general CV technique)

3. Min. reconstruction error

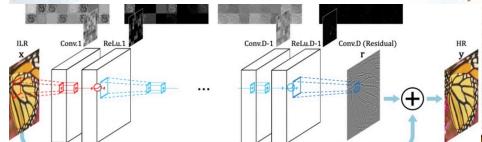
requires two cameras during training
only one at test time

⇒ Image Super Resolution

* Add skip connection from input to last layer

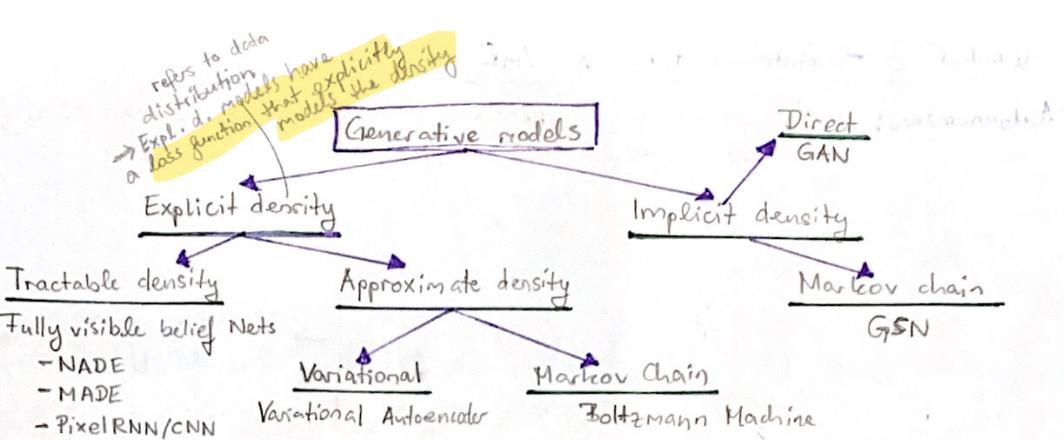
⇒ the network only learns the residual needed for high-resolution.

⇒ No need to learn the image structure

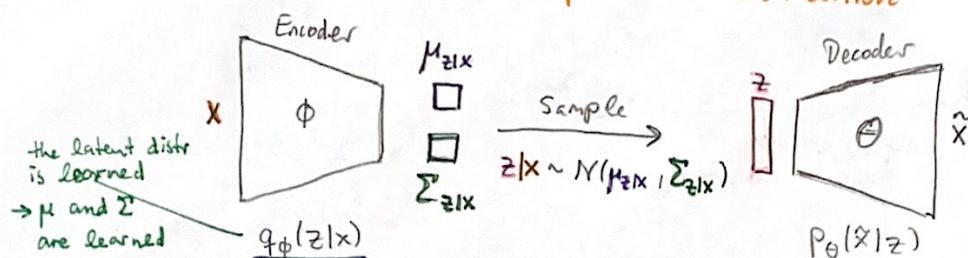


- Generative models:

↳ create new samples from the same data distribution



→ Variational autoencoder: → Latent space is a distribution



At test time: Sample from latent space + pass through decoder

At training time: Optimize parameters ϕ and θ

$$\begin{aligned}
 \text{Loss function: } \log p_\theta(x_i) &= \mathbb{E}_{z \sim q_\phi(z|x_i)} [\log p_\theta(x_i)] \\
 \text{from Bayes' rule } p_\theta(z|x) &= \frac{p_\theta(x|z)p_\theta(z)}{p_\theta(x)} \quad \xrightarrow{\text{maximize the exp.}} \text{log prob. of seeing the data point given the model parameters} \\
 &= \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[\frac{p_\theta(x_i|z)p_\theta(z)}{p_\theta(z|x_i)} \right] \\
 &= \mathbb{E}_{z \sim q_\phi(z|x_i)} \left[\frac{p_\theta(x_i|z)p_\theta(z)}{p_\theta(z|x_i)} \frac{q_\phi(z|x_i)}{q_\phi(z|x_i)} \right]
 \end{aligned}$$

$$\begin{aligned}
 &= \mathbb{E}_z [\log p_\theta(x_i|z)] - \mathbb{E}_z \left[\frac{q_\phi(z|x_i)}{p_\theta(z)} \right] + \mathbb{E}_z \left[\frac{q_\phi(z|x_i)}{p_\theta(z|x_i)} \right] \\
 &\text{definition of Kullback-Leibler divergence} \\
 &\rightarrow \text{measures difference between two distributions} \\
 &= \underbrace{\mathbb{E}_z [\log p_\theta(x_i|z)]}_{\text{likelihood}} - \underbrace{\text{KL}(q_\phi(z|x_i) \parallel p_\theta(z))}_{\text{posterior prior}} + \underbrace{\text{KL}(q_\phi(z|x_i) \parallel p_\theta(z|x_i))}_{\text{Measures difference between latent distr. And Gaussian prior}} \\
 &\rightarrow \text{Log-likelihood of reconstructing the data point given } z \\
 &= \mathcal{L}(x_i|\phi, \theta) \quad \text{is a lower bound on the total loss}
 \end{aligned}$$

$$\Rightarrow \text{Optimize: } \phi^*, \theta^* = \underset{\phi, \theta}{\operatorname{argmax}} \sum_{i=1}^N \mathcal{L}(x_i|\phi, \theta)$$

Things to note: * Latent space dimensions are interpretable
→ e.g. head pose, smile

* VAE Use-case: For generative models (→ sample from model)

Degree of smile



Head pose

- Image Synthesis without GANs:



* Transform semantic segmentation image into real image

→ Can't use L2-loss between predicted image and ground truth image
(it will penalize e.g. coloring white cars black, even though we only care about realistic colorization) → could penalize realistic results

⇒ Use perceptual loss: measures the "content" of the image

(a.k.a. content loss / feature reconstruction loss)

1. Take pretrained VGG network trained for image classification (\rightarrow ImageNet)
 2. Pass synthesized & ground truth image through the network
 3. Compare the feature maps/layers throughout the network

$$\ell_{\text{feat}}^{\phi_j}(\hat{y}, y) = \underbrace{\frac{1}{C_j H_j W_j}}_{\text{feature Map size}} \|\phi_j(\hat{y}) - \phi_j(y)\|_2^2$$

↑
feature maps of the g.t. image
at layer j

⇒ Intuition: the image classification network will produce "similar" features in deeper levels for images with the same object (e.g. car) independent of fine details



- Style transfer: Content image + Style image \rightarrow style transf. image

⇒ Use style loss! → measures the style but not the content

1. $-''$ } as in content loss
2. $\frac{1}{2}''$ }

3. Compare the Gram matrices of feature maps throughout the network

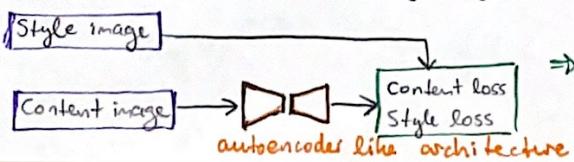
$$l_{\text{style}}^{\phi, i}(\hat{y}, y) = \| G_j^\phi(\hat{y}) - G_j^\phi(y) \|_F^2 \quad \text{Frobenius Norm}$$

with $G_j^\phi(x)_{c,c} = \frac{1}{C_j H_j W_j} \sum_{h=1}^{H_j} \sum_{w=1}^{W_j} \phi_j(x)_{h,w,c} \cdot \phi_j(x)_{h,w,c}$

\Rightarrow Intuition: Captures information about which features tend to activate together

* Option to combine loss functions: Content loss + Style loss
 w.r.t. content image w.r.t. style image \Rightarrow iteratively update white noise image to take on

* Fast Neural Style Transfer:
↳ In training: Use same style
image & different
content images



\Rightarrow At test time:
Only one forward pass needed.

ADL4CV 4 - Graph Neural Networks (GNN)

(5)

order of pixels matters

* Different data domain:

- e.g. ~~convolutions~~ require a certain structure (\rightarrow domain is regular)

\Rightarrow GNNs work on new domain: - Irregular

e.g. - point clouds

- citation networks (node = paper, conn = citation)

- social networks

- Permutation invariant (order of points does not matter)

- Transformation invariant (e.g. the prediction is the same for a rotated point cloud)

* Key challenges:

Variable input size (# of nodes and edges)

Need invariance to node permutations (order in which nodes are processed)

* General idea:

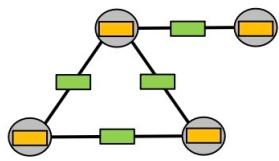
Encode contextual graph info in node embeddings by iteratively combining neighboring nodes' features

"Each update step is understood as a layer in common NNs"

Graph: $G = (V, E)$

Initial node embeddings: $h_i^{(0)}$ with $i \in V$

Initial edge embeddings: $h_{i,j}^{(0)}$ with $(i, j) \in E$



Initial Graph

At iteration step l : (Method only w/ node embeddings)

collect info from all neighbors

$$\textcircled{1} \quad \text{Message passing: } m_v^{(l+1)} = \sum_{u \in N_v} M^{(l)}(h_u^{(l)}, h_v^{(l)}, h_{(u,v)}^{(l)})$$

Learnable function w/ shared weights across graph
aggregation over all neighbors (order invariant)

$$\textcircled{2} \quad \text{Embedding update: } h_v^{(l+1)} = U^{(l)}(h_v^{(l)}, m_v^{(l+1)})$$

Learnable function with shared weights across graph

There are different types of GNNs depending on the formulation of $M(\dots)$ and $U(\dots)$.

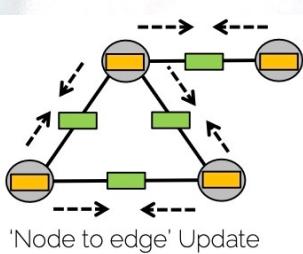
example: Graph Conv. NN

$$m_v^{(l+1)} = \sum_{u \in N_v, u \neq v} \frac{h_u^{(l)}}{\sqrt{|N_v| \cdot |N_u|}}$$

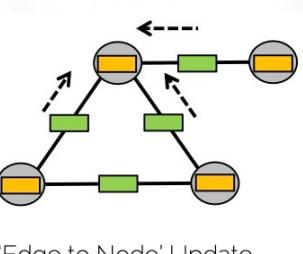
Normalization

$$h_v^{(l+1)} = \sigma(W^{(l+1)} m_v^{(l+1)})$$

non-lin. learnable matrix with size (# channels out, # channels in)



'Node to edge' Update

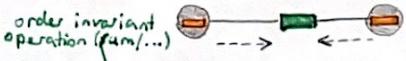


'Edge to Node' Update

Method with Node + Edge updates:

① Node-to-edge updates:

$$h_{(i,j)}^{(l+1)} = Ne(h_i^{(l)}, h_{(i,j)}^{(l)}, h_j^{(l)})$$

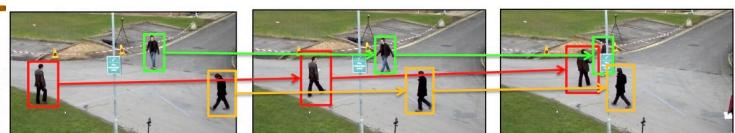


$$m_i^{(l+1)} = \bigoplus \{h_{(i,j)}^{(l+1)}\}_{j \in N_i}$$

edge embeddings
neighbors of node i

$$h_i^{(l+1)} = Ne(m_i^{(l+1)}, h_i^{(l)})$$

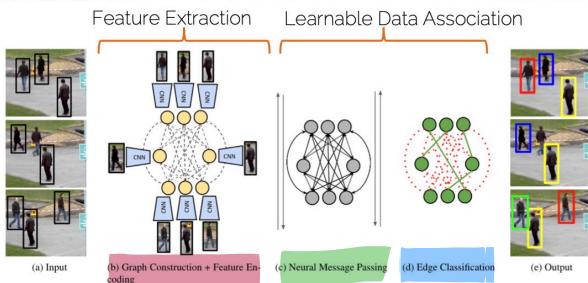
- Message passing networks for CV:



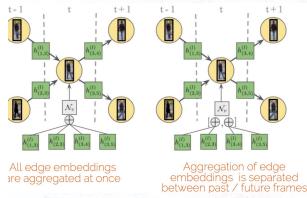
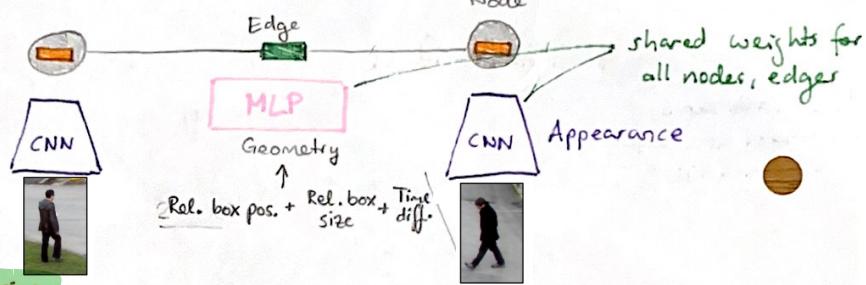
* Multi-object tracking:

1. Object detection
Every object is node and each edge can be a trajectory
2. Data Association
Connect nodes → Use the GNN to find the best connections

- In Step 2:
- (b) Encode appearance and scene geometry cues into node & edge embeddings (e.g. distance)
 - (c) Propagate cues
 - (d) Classify edge embeddings → predict solutions of the tracking problem



→ Feature encoding



→ Time-aware Message Passing
→ Separate the aggregation of edge embeddings between past / future frames

→ Classifying edges

"Classify into active and inactive edges"
→ use MLP with edge embeddings as input → predict whether an edge is active/inactive

Loss function:

$$\mathcal{L} = \frac{-1}{|E|} \sum_{l=l_0}^L \sum_{(i,j) \in E} \underbrace{\omega \cdot y_{(i,j)} \log(\hat{y}_{(i,j)}^{(l)}) + (1 - y_{(i,j)}) \log(1 - \hat{y}_{(i,j)}^{(l)})}_{\text{edge prediction w/ sigmoid at iteration } l}$$

sum over last steps

Weight to balance active/inactive edges

BCE - Loss

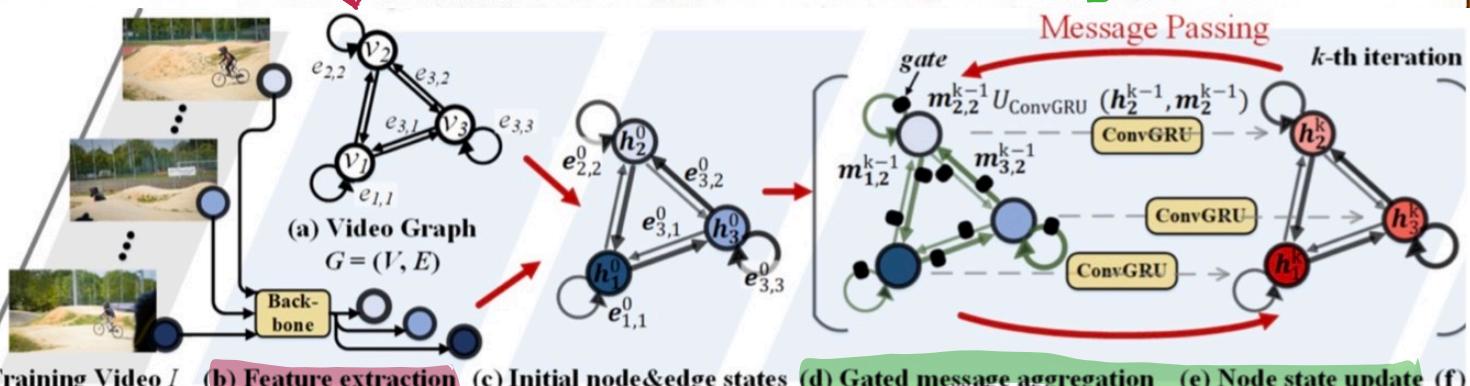
* Video object segmentation:

→ Goal: Generate accurate and temporally consistent pixel masks for objects in a video sequence.

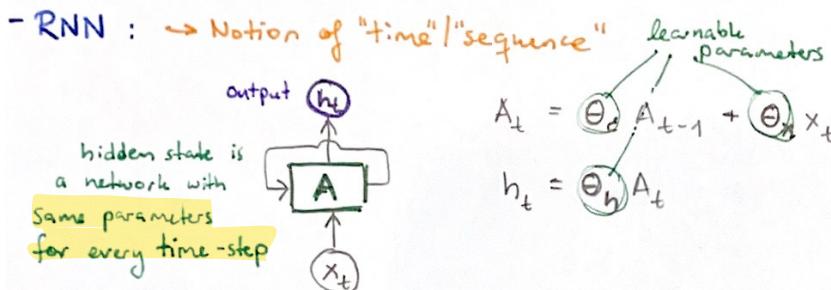
→ If only one object to segment ⇒ Each node is a frame

* Initial embeddings: extract features with DeepLabV3

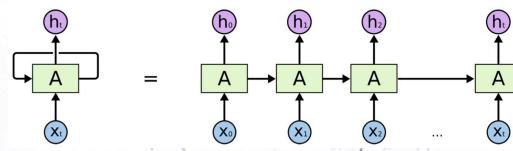
* Message passing with spatial information: conv. recurrent network to preserve embeddings (embeddings are feature maps) → keep pixel outputs



ATTENTION:

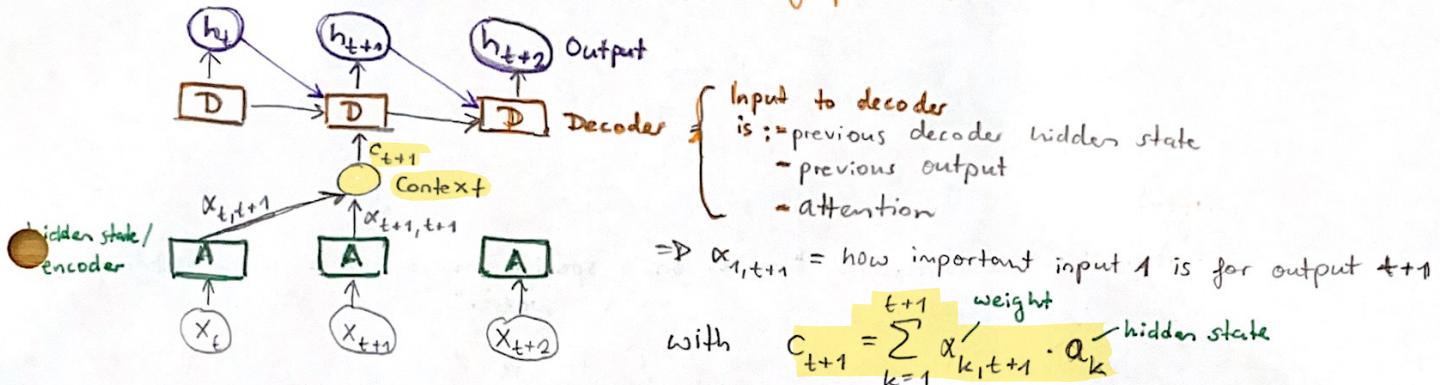


Unrolling a RNN:



\Rightarrow Problem with long-term dependencies (which input is important for output h_t)

- Attention: \rightarrow Aims to solve the long-term dependency problem



* Soft Attention = all attention masks / weights sum up to 1.

* To compute the attention mask: (Requires training a small NN)

Previous decoder state d_t
encoder hidden state a_1

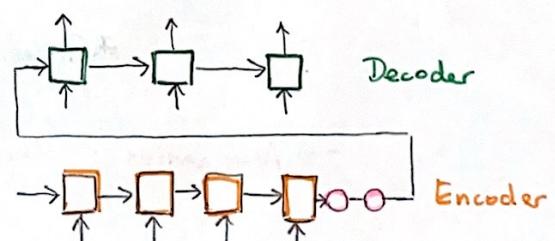
$$\text{NN} \rightarrow f_{1,t+1}$$

Normalize $\Rightarrow \alpha_{1,t+1} = \frac{\exp(f_{1,t+1})}{\sum_{k=1}^{t+1} \exp(f_{k,t+1})}$

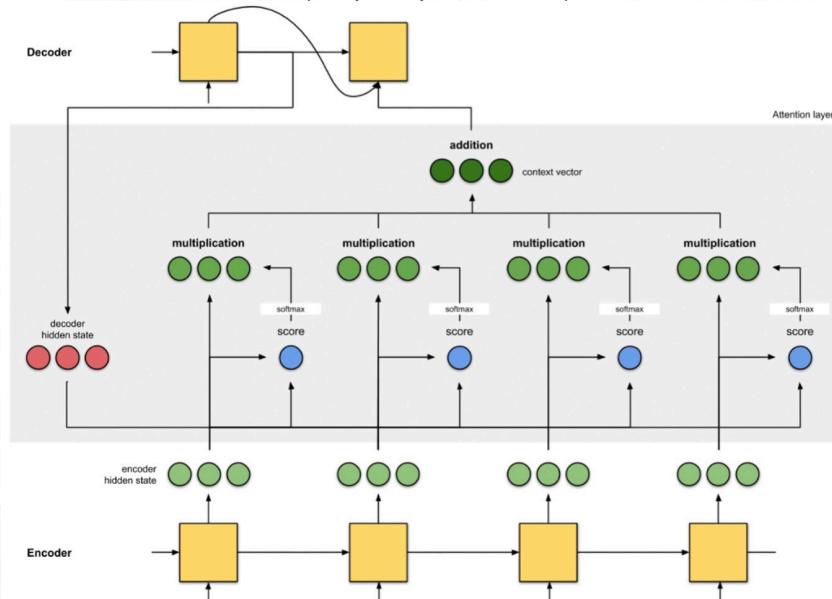
- Seq2Seq: \rightarrow Network for sentence translation

1. Read the whole sentence in language 1 and embed it

2. Translate it to language 2 using embedding + each word



\Rightarrow Also use Attention (Seq2Seq + Attention)

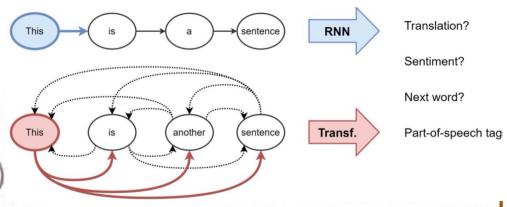


- Transforness: → No RNN, just attention.

* Current SotA in NLP

* Are basically a Graph Attention Network (GAT)

→ GNN where the aggregation step is done via attention (weighted sum)



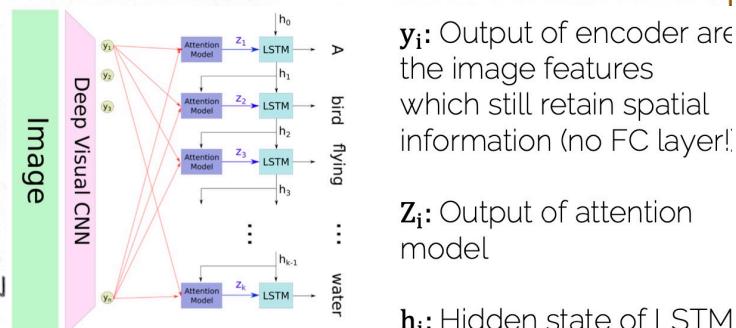
- Attention for CV:

* Image captioning

1. Input: image
2. Encode the image (feature extractor/CNN, e.g. AlexNet)
3. Decoder: Attention-based RNN
4. Output: Sentence describing the image

(Soft-)

→ Attention-mechanism: "focuses on a specific image patch when outputting a word."



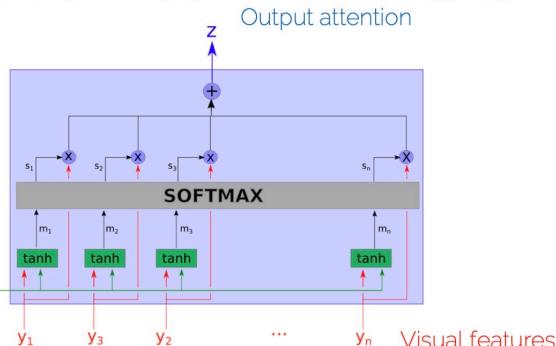
y_i : Output of encoder are the image features which still retain spatial information (no FC layer!)

z_i : Output of attention model

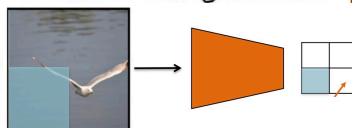
h_i : Hidden state of LSTM

Why do we need Attention?

Inputs to the Attention model: Feature descriptor for an image patch (→ feature map)

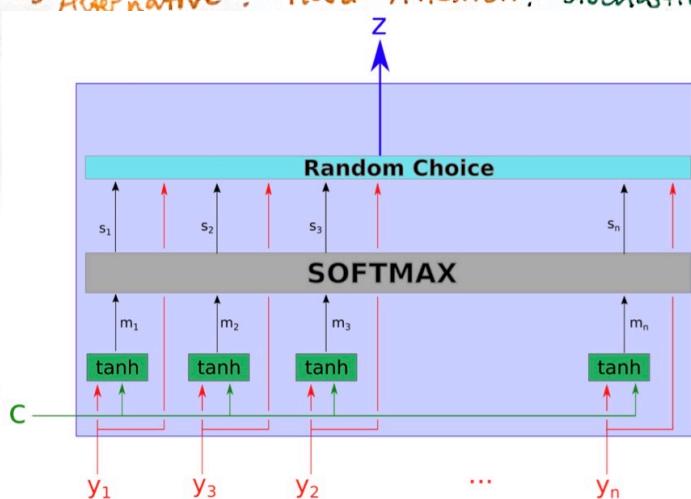


- tanh to create bounded outputs
- Softmax for attention values between 0 and 1
- Multiply by image features to rank them by importance



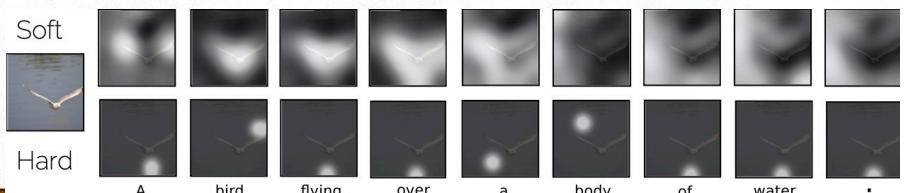
* Process is deterministic and can be backproped

→ Alternative: Hard-Attention: stochastic process; gradient via Monte Carlo sampling + estimation



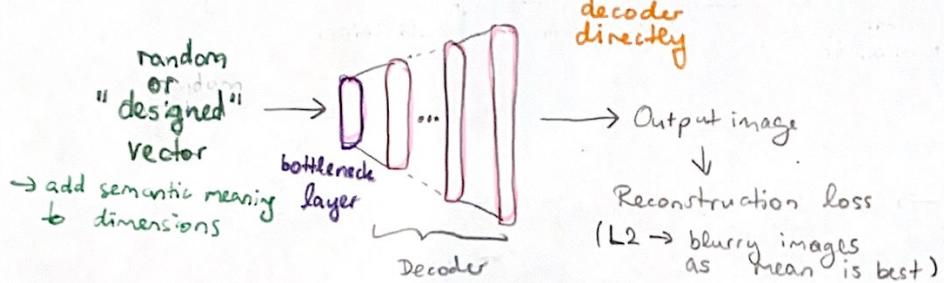
Choose features by sampling with probabilities s_i

• Soft vs hard attention



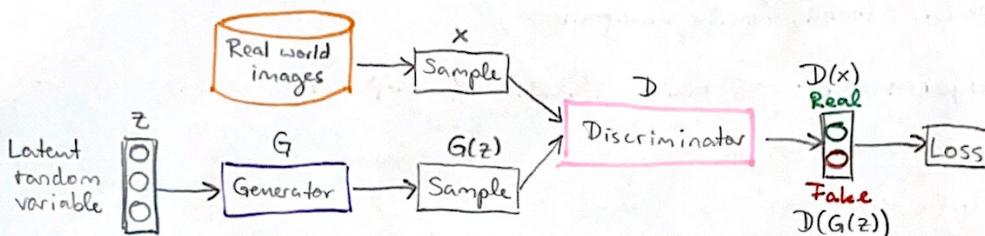
ADL4CV 5 - Generative Neural Networks

- Decoder as Generative model: → train decoder directly



- Generative Adversarial Networks (GANs): → Learn the loss function D

* is an implicit density model = the loss function does not explicitly match the data distribution



On real data:

$D(x)$ tries to be near 1

On fake data:

D tries to make $D(G(z))$ near 0, G tries to make $D(G(z))$ near 1

→ Alternating gradient updates:

"First debug step:
Only do 1). If $\mathcal{J}^{(D)}$
can't go to 0 something
is wrong."

1) Fix G and perform gradient descend on the discriminator loss:

from BCE loss

$$\mathcal{J}^{(D)} \downarrow = -\frac{1}{2} \mathbb{E}_{x \sim P_{\text{Data}}} \log D(x) - \frac{1}{2} \mathbb{E}_{z \sim P_z} \log (1 - D(G(z)))$$

2) Fix D and perform the gradient descend on the generator loss:

G min. the prob. that D is correct

$$\mathcal{J}^{(G)} = -\mathcal{J}^{(D)}$$

"Minimax game loss"

G max. the log prob. of D being mistaken

$$\mathcal{J}^{(G)} = -\frac{1}{2} \mathbb{E}_z \log D(G(z))$$

used in practice

→ G can still learn even if D rejects all the gen. samples
"Heuristic method loss"

→ Loss functions:

" D - and G -loss must be balanced
→ D provides supervision for G "

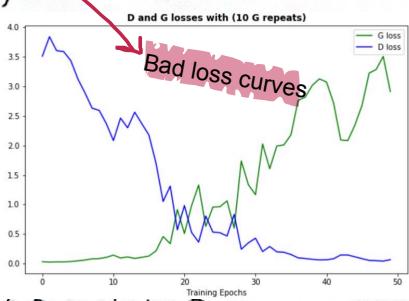
* If the D loss gets to 0, there is no recovery
* Good loss curves = plateau at constant value ($\neq 0$)

If D rejects all samples from G , the G -loss is always high and D -loss is always low → G cannot learn

* The larger the image, the more difficult the generator architecture becomes
* Cannot determine if training is done from loss
 ↳ Look at generated samples

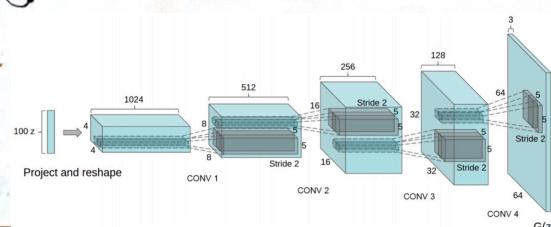
* Training schedules:

→ Adaptive schedules, e.g. $\begin{cases} \text{while } \text{loss-}D > t-D: \text{train } D \\ \text{while } \text{loss-}G > t-G: \text{train } G \end{cases}$



deep convolutional GAN

- DCGAN: → 2015 Paper; Is a good baseline architecture



- Mode Collapse: → is the Problem when the generated images fool D, but have little variety.

$$\max_D \min_G V(G, D) \neq \min_G \max_D V(G, D) = \min_G \max_D \mathbb{E}_{x \sim P_{\text{Data}}} \log D(x) + \mathbb{E}_z \log (1 - D(G(z)))$$

↓ ↓
G wants D wants
to min. the to max. it
value func.

value function definition

- loop ordering**
- * if D is in the inner loop → converges to correct distr.
 - * if G is in the inner loop → D is fixed while updating G → "learn to fool the same D"
→ easy to converge to one sample

Additional observations: (modes = faces, cars, etc. in the data)

* More modes → smaller recovery rate

* Larger latent space → more mode collapse

⇒ For these reasons, GANs are usually trained on specific domains (smaller distributions)

- Evaluation of GAN performance: → main difficulty of GANs is knowing how good they are

... since Convs are "local" operators

* GANs have problems with generating global structures (e.g. counting)

* Human evaluation: → check training curves

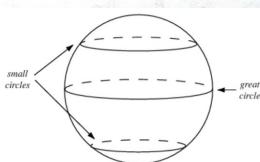
→ predict same latent codes every n steps ↗ variety
become more realistic

* Inception score (IS): 1) Train a classifier on the real images & classify generated images
2) Saliency: Check if images are classified with high confidence
3) Diversity: Check if we obtain images from all classes

* Look at discriminator: → Use D-layers for feature extractor in a classifier
→ fine-tune the last layer
→ if high class accuracy → we have good D and G

- GAN Hacks:

- 1) Normalize inputs between -1 and 1 + tanh as last generator layer
→ so that all inputs to the discriminator are in the same range
- 2) Don't sample z from uniform distr. / Sample z from Gaussian or from surface of hypersphere
→ Do interpolations via a great circle rather than linearly
- 3) Try BatchNorm ; Construct mini-batcher with all real or all fake images
↳ stabilizes training



- 4) Adam for generator, SGD for discriminator
- 5) One-sided label smoothing: $D^{(D)} = -\frac{1}{2} \mathbb{E}_{x \sim P_{\text{Data}}} \log D(x) - \frac{1}{2} \mathbb{E}_z \log (1 - D(G(z)))$
↳ If image is real ($y=1$) set $y=\lambda$ to some value $\lambda < 1$
→ Reduces confidence of D
- 6) Historical Generator batches: Construct batch from current + past generated images
→ pass to D ; stabilizes D
- 7) Avoid sparse gradients: → Use LeakyReLU instead of ReLU
→ Downsampling: Avg. pool or conv + stride
→ Upsampling: deconv + stride or PixelShuffle

Problem: D is noisy due to SGD

- 8) Exp. averaging of weights: keep second 'vector' of weights that are averaged -> almost no cost, average of weights from last n iters

- Other objective functions: → Don't change it too much; mostly the heuristic loss is fine (8)
The loss function alone will not make it suddenly work!

→ EBGAN: Energy-based Generative Adversarial Networks

* discriminator is an autoencoder trained on real images

$$D(x, z) = D(x) + [m - D(G(z))]^+ \quad \rightarrow_{>0} \text{Want reconstruction cost } D(x) = \| \text{Dec}(\text{Enc}(x)) - x \| \text{ to be low}$$

$$\mathcal{L}_G(z) = D(G(z)) \quad \rightarrow_0 \Rightarrow \text{if there is a good reconstruction, it is a real image}$$

here $[u]^+ = \max(0, u)$ - We want to penalize the discriminator if the reconstruction error for generated images drops below a value m.

→ BEGAN: Boundary equilibrium GAN

* Similar to EBGAN, but instead of a reconstruction loss measure the difference in the data distr. of real and generated images
↳ Wasserstein distance / Earth-mover distance

→ WGAN: Wasserstein GAN

* BEGAN with Wasserstein distance (comp. dist. between two distributions)

→ Problem: Wasserstein-dist. is expensive to compute

↳ Get upper bound from dual: $W(P_r, P_\theta) = \sup_{\substack{\|f\|_{L^1} \leq 1 \\ f \text{ is 1-Lipschitz function}}} \mathbb{E}_{x \sim P_r}[f(x)] - \mathbb{E}_{x \sim P_\theta}[f(x)]$

1-Lipschitz constraint: $|f(x_1) - f(x_2)| \leq |x_1 - x_2|$

⇒ Define Update with weight clipping to train a Lipschitz-function:

$$w \leftarrow w + \alpha \cdot \text{RMSProp}(w, g_w)$$

$$w \leftarrow \text{clip}(w, -c, c)$$

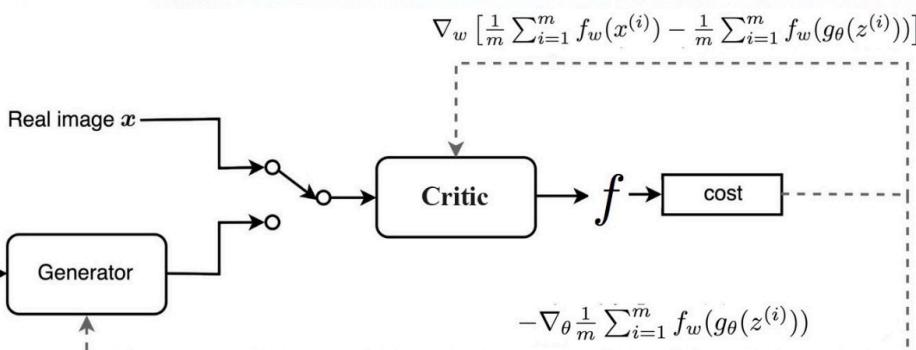
*f is difficult to define
→ approximate f by learning a NN*

→ f is the critic func

⇒ Main advantage of WGAN: Discriminator/Critic loss converges to zero instead of plateauing

- + mitigates mode collapse
- + generator still learns when critic performs well
- + actual convergence

- Enforcing Lipschitz constraint is difficult
- Weight clipping is "terrible"
 - > too high: takes long time to reach limit; slow training
 - > too small: vanishing gradients when layers are big



ADL4CV 6 - GAN Architectures & Conditional GANs

(9)

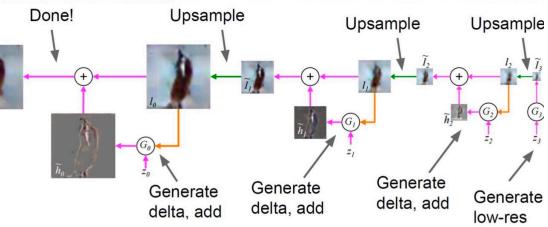
Multiscale GANs → deal w/ the problem where generating the global image structure is difficult
 ↳ from 2015

Idea: Multiple generators at different image scales w/ upscaling in between

↳ each higher-scale generator adds details to the current image

(Generators at low-res learns the global structure)

→ We need discriminators at every level/resolution



Progressive Growing GANs → from 2017; extends Multiscale GANs

1. Start by generating low-res images + train

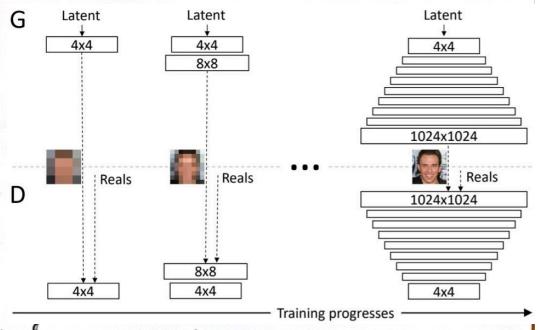
2. Append upsampling block to G and downsampling to D

↳ higher resolution; initialize weights in lower layers w/ previous weights

3. fixed schedule when to append the next layer

Additional Tricks:

- 1-layer feature maps (get RGB-image from series of 1x1-conv) ↳ makes appending upsampling blocks easier → "to-RGB" operator
- When adding new blocks: lin. crossfade between "old" & "new" Image reconstr. loss ↳ gradually move weight to larger layer ⇒ stabilizes the loss
- Analogous "from-RGB" operator for discriminator



Conditional GANs → for Modeling (adding semantic meaning to the output, "what to generate")
 (cGANs)

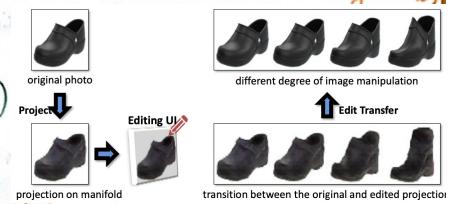
↳ iGAN: 1) project image onto learned latent manifold

2) Edit projection for example (→ edit latent code)

3) "Generate" image w/ modified latent code

→ How to project a real image onto the manifold?

$$\text{with optimization} \quad z^* = \underset{\text{latent vector}}{\arg\min} \underset{\text{reconstr. loss}}{L}(G(z), x^R) \quad \underset{\text{trained generator}}{G}$$



Naive approach:
 Use condition img as input to G

Problem: D only judges on real/fake but not if output img ≈ cond. img

Edit transfer in the OG paper: reproject edits in the generated image onto the real image → fewer imperfections

→ instead of solving the optimization problem, train an inverting network $z^* = P(x^R)$

$$\Theta_P^* = \underset{\Theta_P}{\arg\min} \sum_{X_n^R} L(G(P(X_n^R, \Theta_P)), X_n^R)$$

like an auto-encoder w/ a fixed decoder G

Hybrid Method: first use inverting network, then fine-tune by solving the optimization problem

→ In practice they add a manifold smoothness loss. stay close to some real image

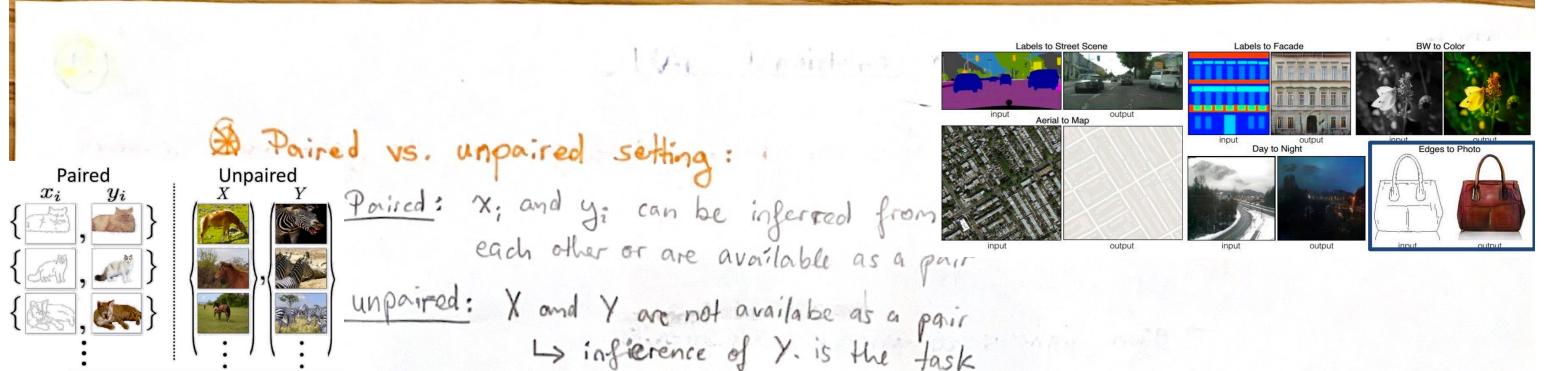
$$\text{Objective: } z^* = \underset{z}{\arg\min} \left\{ \underset{g}{\sum} L_g(G(z), v_g) + \lambda_s \|z - z_0\|_2^2 \right\}$$

guidance

v_g



z_0



Make the discriminator aware of the user guidance (s.t. it has an influence on the generated image)

↳ Need **paired data** (often called self-supervised)
(e.g. edge detection)

↳ true **conditional GAN**

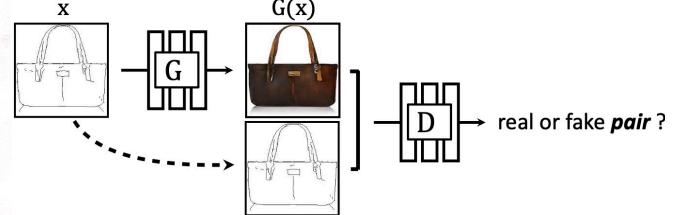
↳ match joint distribution

$$P(G(x), y) \sim P(x, y)$$

⇒ E.g. **Pix2Pix** (2017)

↳ Tricks in Pix2Pix:
- adds ^{L1-loss (GT & output image)} loss to GAN-loss, since the problem is underconstrained

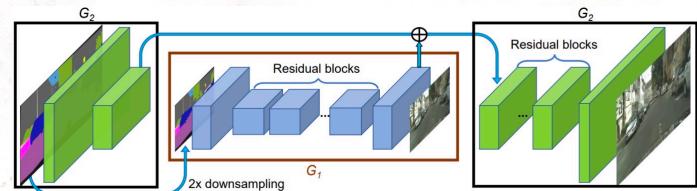
$$\min_G \max_D \mathbb{E}_{x,y} [\log D(x, G(x)) + \log(1 - D(x, y))]$$



- use dropout to add noise → probabilistic model

* Pix2PixHD (2018) → cGANs + Multi-scale

"Generator within a generator"
→ it is more efficient to use the weights in this hierarchical architecture (You can only fit so many weights in a GPU)



⇒ Conditioning at multiple scales (in the image at 2 resolutions)

* How to formulate a conditional GAN in the unpaired setting?

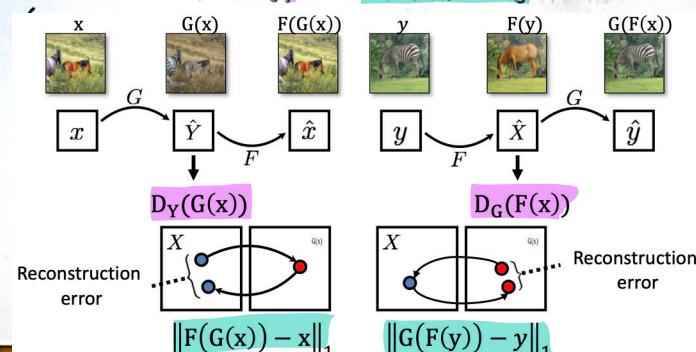
* We want a bijective mapping (each input points at exactly one output)

* Method: Enforce **Cycle-Consistency** (mapping from X to Y results in unique Y and the other direction gives unique X)

* Cycle Consistency Loss:

↳ Use another generator F to map Y to X

⇒ **Discriminator loss** + **Recontr. loss** (In each direction → 4 losses in total)

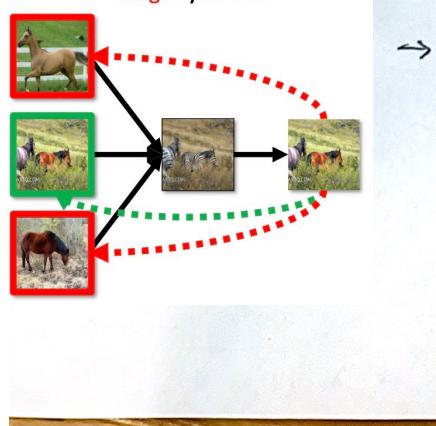


Problem: We have not input-output pairs anymore

→ how to make G care about the input?

This will probably lead to mode collapse

Sample cycle loss



ADL4CV 7 - Videos & Autoregression

(10)

StyleGAN (2019)

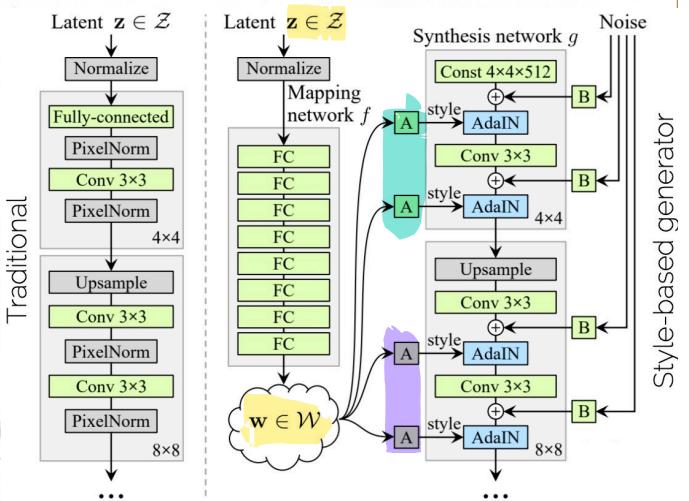
use style image as input to mapping net

↳ Difference to traditional GAN:

* Mapping network maps latent code z to an interpretable style vector w

* Feed in noise & style vectors every few layers (they use AdaIN) → added to feature map

* Style vectors added at low-res are coarse styles (thereafter middle+fine styles)
→ multi-resolution control (e.g. change face details independent of main features)

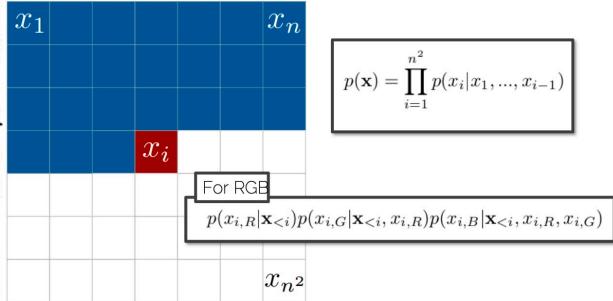


Autoregressive Models

- * GANs learn implicit data distribution (→ outputs are samples & distribution is in the model)
- * Autoregressive models learn an explicit distribution (→ outputs are probabilities) governed by a prior imposed by the model structure

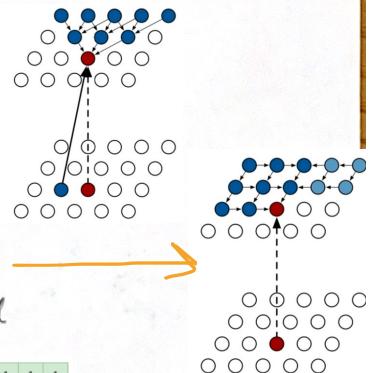
PixelRNN (2016)

- * Interprets image pixels as a sequence → product of conditional distributions
- Next pixel determined by all previous pixels
- ⇒ $p(\vec{x}) = \prod_{i=1}^{n^2} p(x_i | x_1, \dots, x_{i-1})$
- prob. of all pixels (in practice = softmax w/ 256 options) → conditioning on previous pixels over complete image
e.g. prob that a pixel is green



- * To reduce the sequence length: Row LSTM model architecture

- ↳ Each pixel depends on "pyramid" of pixels above
- Image can be processed row by row ⇒ on GPU
- However, incomplete context for each pixel
- ↳ Solution: Diagonal BiLSTM model architecture
⇒ diagonal pixels can be processed in parallel



PixelCNN (2016) → to lower computational cost

- * Masked convolutions to avoid seeing future context → mimics sequence behavior

- * Speciality: Replace ReLUs with Gated Blocks

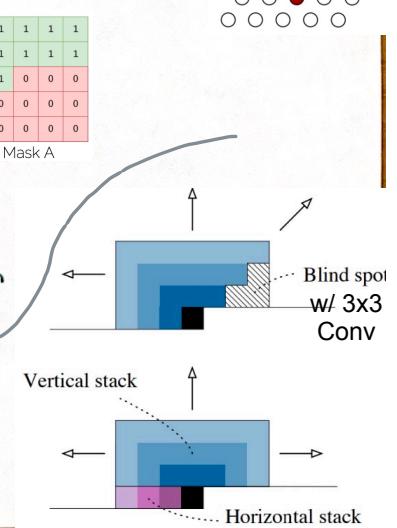
$$y = \tanh(W_{k,f} * x) \odot \sigma(W_{k,g} * x)$$

Weights in k th layer convolution
element-wise product

↳ mimics the LSTM

- * Problem: Unseen context for pixel (Blind spot)

↳ Split convolutions in two stacks



④ Autoregressive models vs. GANs

- ⊕ more stable since probabilistic densities are modelled explicitly
→ easier hyperparam-tuning
- ⊕ applicable to discrete & continuous data
- ⊕ Demonstrations could produce higher quality images
- ⊕ faster to train (But also harder to train)

⑤ Generative Models on Videos

Two options

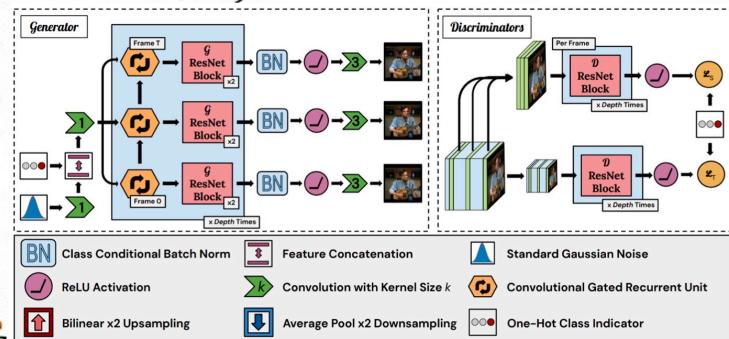
- Single latent code for all frames (⊖ future frames deterministic given past)
- Latent code per-frame (⊖ Need conditioning of future frames on past frames)

⇒ General Issues: Temporal Coherency * Collapse to mean image
(Drift over time, since the most recent frames are also generated and not GT anymore)

→ DVD-GAN (2019) (not discussed in detail)

- * Single random vector seeds a video
- * Per-frame generators
- * Discriminators judge if the sequence is reasonable

→ Conditional GANs work much better on videos



Challenge: cGANs provide high-quality frames, but they are temporally inconsistent.

→ Video2Video Synthesis (2018)

↳ use sequential generator $p(\vec{x}^T | \vec{s}^T) = \prod_{t=1}^T p(\vec{x}_t | \vec{x}_{t-1}, \vec{s}_{t-L})$

Two Discriminators:

past L generated frames past L source frames
in practice: $L=2$

* Conditional img D (is img real?) + conditional vid D (temp. consistency based on optical flow)

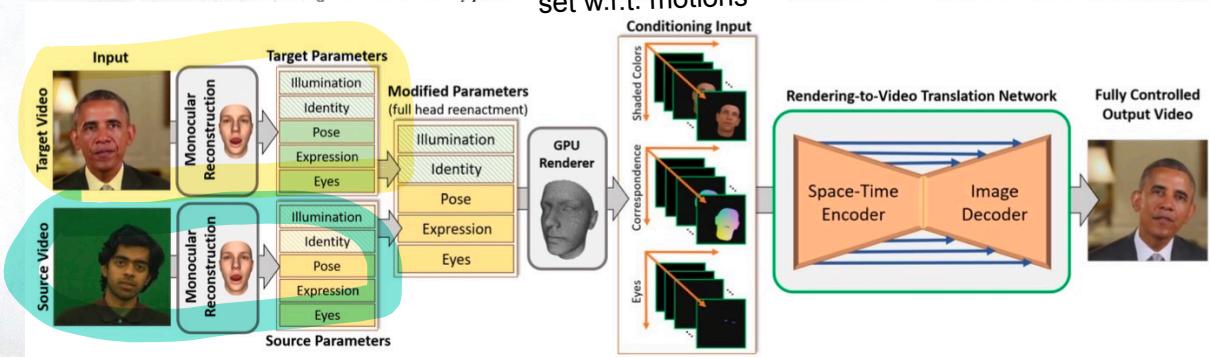
→ Deep Video Portraits (2018)

tracking quality is very important
↳ Train a cond. GAN on a target video w/ face tracking info as guidance (Rerendering)

* At test time use guidance from Source video

→ if the source video is temporally consistent the output video will be temporal consistent

⇒ Synthetic data for tracking is a great anchor → Needs to stay within the training set w.r.t. motions



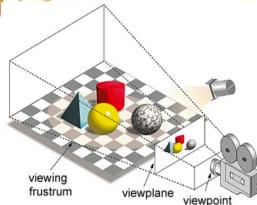
"Everybody Dance Now" has the same concept, but conditions on tracked body poses and generates dance videos

ADL4CV 8: Neural Rendering

(11)

- What is rendering? → Creating a 2D projection of a 3D-scene-representation

- * Rendering according to camera intrinsics + extrinsics
 - ↳ focal length
 - ↳ 6 DOF
 - ↳ principal point



- * Photorealistic renderings of "perfect" models are possible → Reconstructing 3D models from 2D views is difficult

- What is novel viewpoint synthesis?

- * NN encodes the entire scene description (+lighting, materials, etc.)

- * Input is camera pose + viewing direction → Output is novel view image

- Neural Rendering w/ Pix2Pix: (2017)

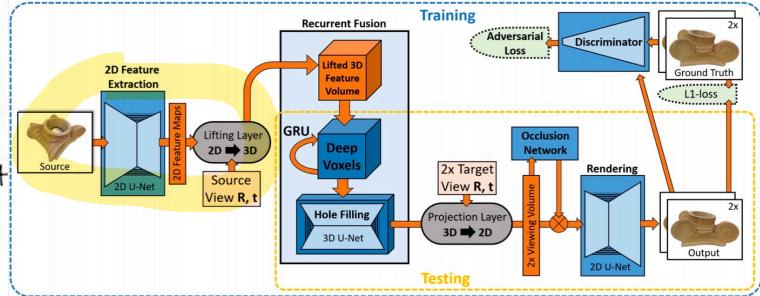
At train time: Pose + ground truth image → Reconstruction loss to train

At test time: Unseen pose → extrapolated, unseen image

- * Problem: Pix2Pix uses conv. → 2D-operators need to learn 3D scene representation
⇒ Artefacts, flickering

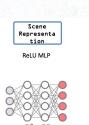
- → w/ Deep Voxels: (2019)

- * Extract features from 2D images and lift them to 3D voxel grid (according to camera pose R_t, t)
- * Deep voxels represented as 3D U-net
- * Projection from 3D to 2D according to target camera pose R_t, t
- * Occlusion network masks out resolved voxels (→ depth estimation network)
But limited resolution due to dense 3D voxel grid
- ⇒ All operations are differentiable



- * Lifting from 2D to 3D works great
 - No need to take specific care for temp. coherency!
 - But: limited resolution due to dense 3D voxel grid

- → w/ an implicit representation:



- ↳ Scene representation networks:

(2018) * Can generalize!

* Scene representation is ReLU MLP

* Renderer is generalized learned sphere tracing
↳ How far until ray hits surface?

- → NeRF: Neural Radiance Fields: (2020)



* Scene repr. is ReLU MLP + pos. encoding

* Renderer is volumetric stratified sampling

* No generalization → MLP just memorizes to color of points in space
⇒ then interpolation

- → w/ Neural Textures: Features on a 3D mesh (2018)

- * Map a neural texture on a given mesh (instead of RGB-texture; features) that live on top of mesh)

"texels"

- * Project 3D mesh + neural texture to 2D

According to camera pose

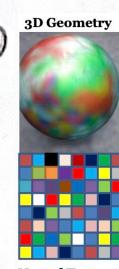
- * 2D U-net renders the realistic image

3D Geometry

ReLU MLP + Positional Encoding

Viewpoints

View R, t

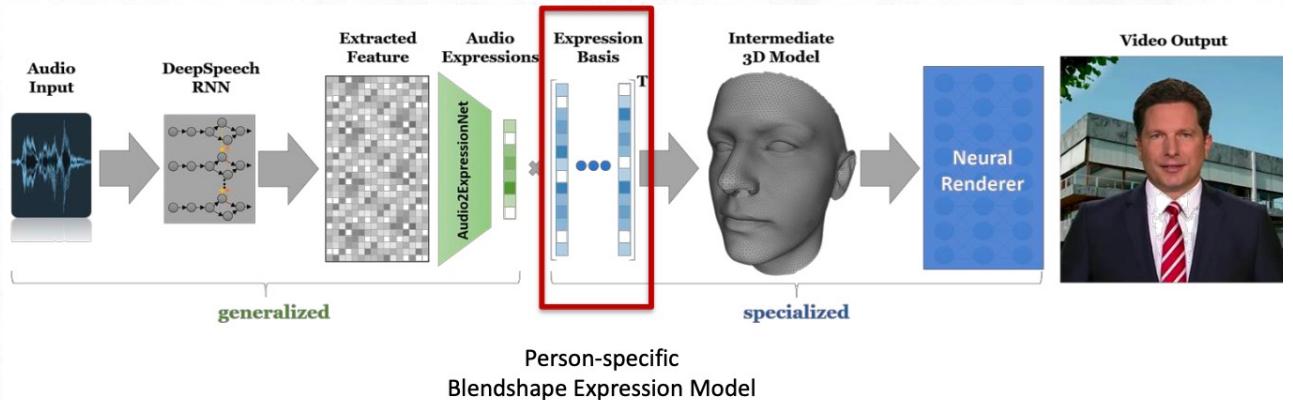


- Comparisons of neural rendering approaches:

Scene Representation					
Renderer	(Alpha) compositing	Volumetric Ray-based	Rasterization	Splatting	Sphere-Tracing Volumetric
Pros	Fast rendering High quality Generalizes	“True 3D” High quality	High quality	High quality	High quality May generalize!
Cons	Only 2.5D Size	No reconstruction priors Memory $O(n^3)$	Requires good SFM No compact representation	Requires good SFM	Expensive rendering, training

- Neural voice puppetry:

- * Use the DeepSpeech model to extract features (logits of alphabet)
 - * Convert to expressions
 - * ↳ drive facial model



ADL4CV 9 : Deep Learning in higher dimensions

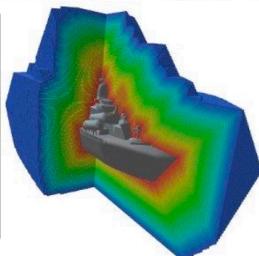
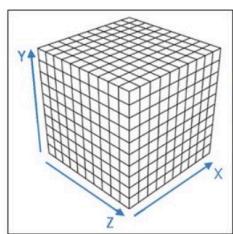
12

- Multi-dim. ConvNets:

- * 1D : Audio, Speech → e.g. WaveNet
- * 2D : Images
- * 3D : Videos, 3D data → 3D-shape/scene classification
⇒ requires more memory, but less data
(as the viewing direction does not need to be learned)
- * 4D : dyn. 3D data, simulations

- Volumetric grids:

"voxel" ≈ pixel in 3D



a field is a continuous representation

→ Volumetric data structures:

surface occupancy gets smaller w/ higher resolution

- Occupancy grids (binary value per voxel)
- Ternary grids (is it empty space? surface point? unknown?)
- Distance field (every voxel encodes distance to closest surface)
↳ at surface point = 0
- Signed distance field (sign encodes if voxel is in front or behind surface)
↳ empty space
↳ unknown space

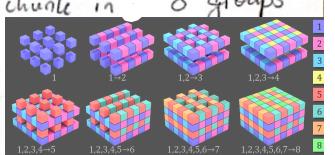
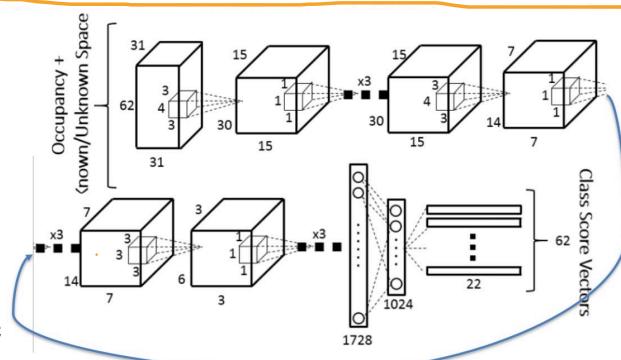
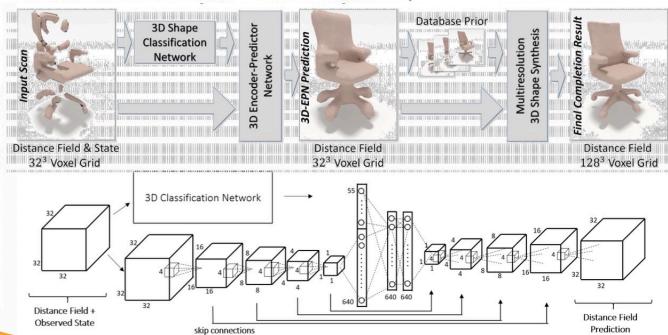
→ 3D shape completion:

- * 3D U-net architecture
- * Input is signed distance field
- * Output is distance field
⇒ network makes prediction

→ Semantic Segmentation in 3D:

- * Process scene volume in chunks (here $62 \times 31 \times 31$)
- * Input is signed distance field chunk
- * Output are classification scores for the 62 central voxels
⇒ Sliding window / chunk since input grid can vary in size
- ⇒ Doing shape completion while semantic segmentation helps the segmentation
- * Also w/ 3D autoregressive neural networks ⇒ 6x6x4 chunk in 8 groups

→ Volumetric grids need a lot of memory + processing time
→ sliding window



- Volumetric Hierarchies: → increase resolution of voxels close to a surface

* Requires the 3D structure to be known beforehand → only for discriminative tasks (e.g. classification)

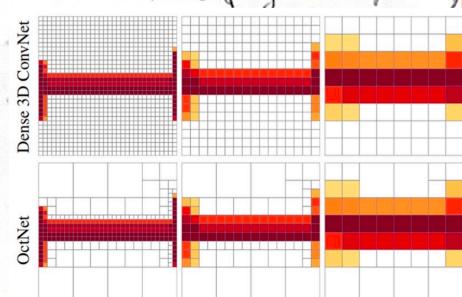
* Voxels in (signed) distance field

* Reduced memory consumption (and runtime)

- ↳ DenseNet: $\mathcal{O}(n^3)$
- ↳ OctNet: $\mathcal{O}(n^2)$

 small performance hit

* In generative models the model needs to predict whether to split voxels



- Multi-view:

do not use 3D information explicitly
→ instead they use 2D rendering and aggregate features in the 2D domain

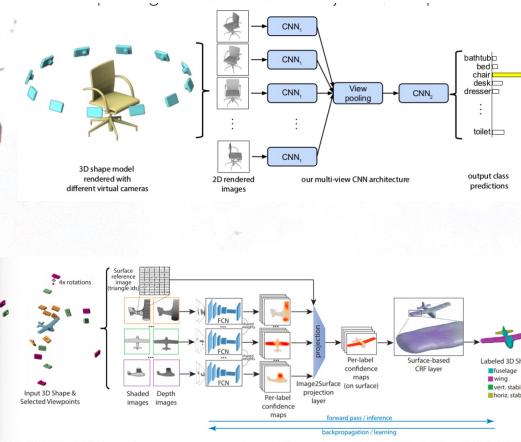
→ Classification:

- * RGB imgs from fixed views around object
- * view pooling (invariant w.r.t. # of views)

* Then classification

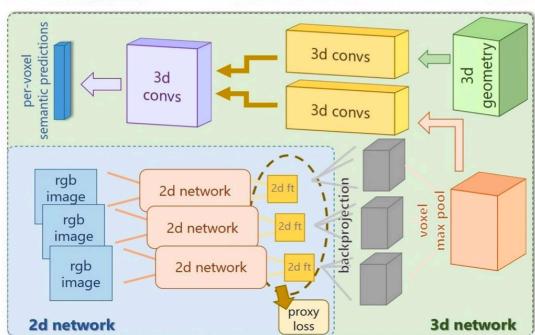
→ Segmentation:

- * 2D-segmentation on every view
→ output are per-label confidence maps
- * Backprojection onto the object surface



- Hybrid: Volumetric + Multi-view:

- * Extracts features from multi-view 2D imgs
- * Backproject color features to voxel grid
- * Add the color voxel grid to the geometry voxel grid
- ⇒ Color and geometry complement each other



- Point Clouds: → Alternative to volumetric voxel grids

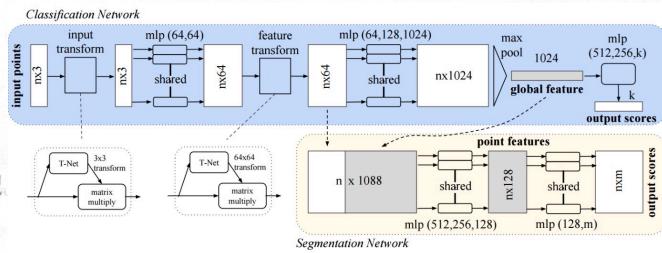
→ PointNet: (2017)

train & test is very fast
But: performance mostly worse than volumetric

- * n input points ($\frac{X}{N}$)
- * extract per-point 1024-feature vector
- * Most indicative features from MaxPool
⇒ output: class scores
- ⇒ concat per-point features with global features for segmentation scores

PointNet ++: (2017)

- * Addresses sampling mismatch (when many samples are in groups far apart)



- Mesh-based: → Polygons represent 3D surfaces

- * Need some differential geometry approx. and perform convolutions in this space prone to noise on real data ⇒ No good results yet

- Sparse Convolutions: → Solve problem where dense convs are expensive in memory and have the dilution problem (set of non-zeros grows quickly)

"Local convolution in the regions where we actually have a surface"

slower than dense convs due to lookup in hash-table
but: higher accuracy

- * Sparse Convs save previously active pixels in a hash-table
⇒ all others are masked out
⇒ less memory usage / allows for larger resolution
- * Submanifold convolution

