

# Visual AI

CPSC 532R/533R – 2019/2020 Term 2

**Lecture 7.** Representation learning I

Helge Rhodin



# Recap: Voxel representations

Idea: A 3D tensor that encodes occupancy

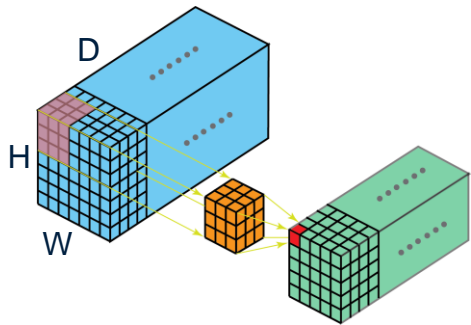
- stores binary values
- occupied or empty cell

Size:  $C \times D \times H \times W$  ( $C$ : channels,  $D$ : depth,  $H$ : height,  $W$ : width)

Batched size:  $N \times D \times H \times W$  ( $N$ : number of elements in mini batch)

Benefits: We can apply 3D convolutions

- A generalization to 2D convolutions with a 3D kernel

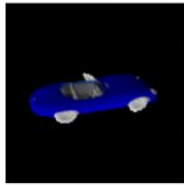
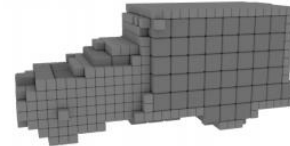


Drawback:

- cubic in memory footprint and computational complexity

Input

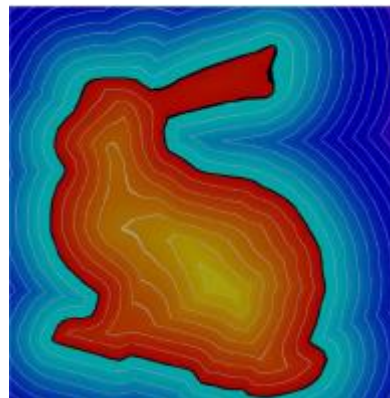
$32^3$



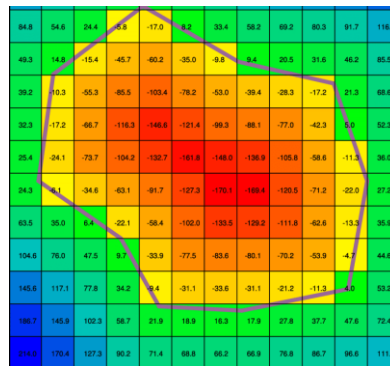
[Octree Generating Networks: Efficient Convolutional Architectures for High-resolution 3D Outputs]

# Signed Distance Field (SDF)

- input domain: dimension equal to the dimension of the space
  - usually two or three-dimensional
- output domain: a scalar
  - negative for inside of the object
  - positive outside
- continuous SDF: defined by a parametric function
  - e.g., sum of Gaussians, neural network
- discrete SDF: defined on a grid
  - e.g. 2D grid or 3D grid
- easy to display SDF in color code  
(red to blue = negative to positive)
- non-trivial to reconstruct the exact shape boundary



Continuous SDF



Discrete SDF

# Implicit functions through NNs

Idea: Train a neural network that takes an image as well as a 3D query point as input and outputs:

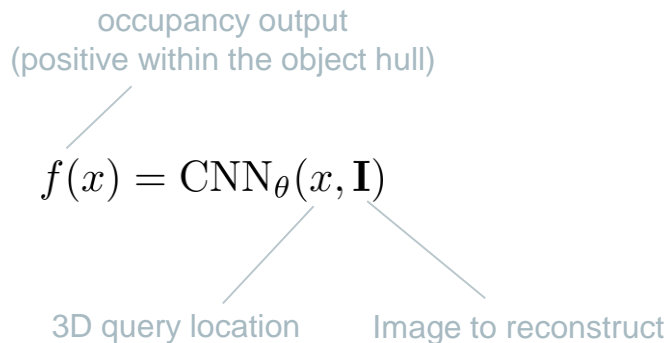
- **negative** for positions **inside** the object
- **positive** **outside** the object
- reconstruct by querying a dense sampling

Advantage:

- No explicit limit on resolution (only limited by NN capacity)

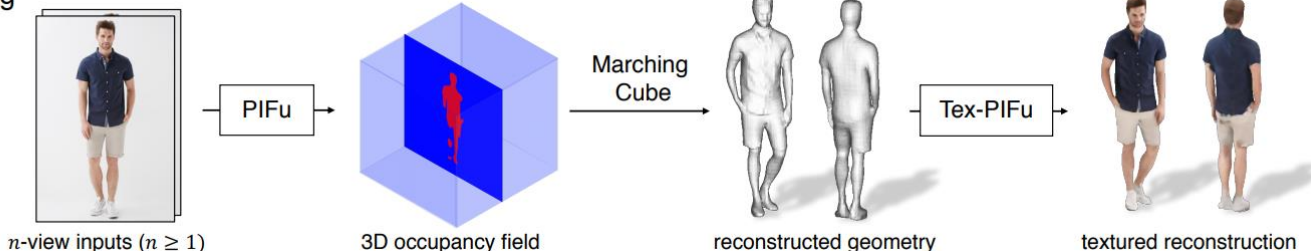
Disadvantage:

- Reconstruction requires many network evaluations, its slow!



*Not straightforward to train...  
wait for the paper presentation*

Testing

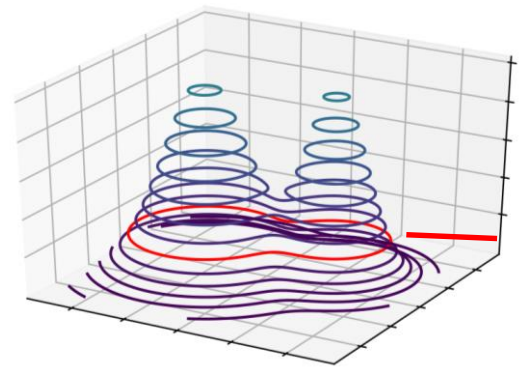


# Recap: Implicit functions

Idea: define complex shapes as the zero-crossing of a function

Size:  $W$  (the number of parameters of the function)

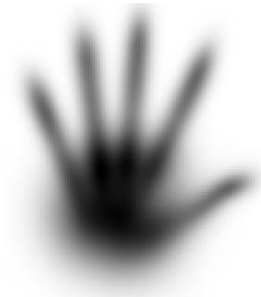
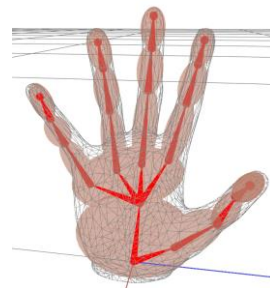
- independent of output space dimension!
- Any parametric function works
  - e.g., mixtures of  $n$  Gaussian distributions with position  $\mu$  and covariance  $\Sigma$



contour line / zero crossing

$$f(x) = \sum_{i=1}^n G(x, \mu_i, \sigma_i)$$

- a neural network?!



[Real-time Hand Tracking Using a Sum of Anisotropic Gaussians Model]

# Recap: Surface mesh

Representation: Vertices connected by edges forming faces (usually triangles)

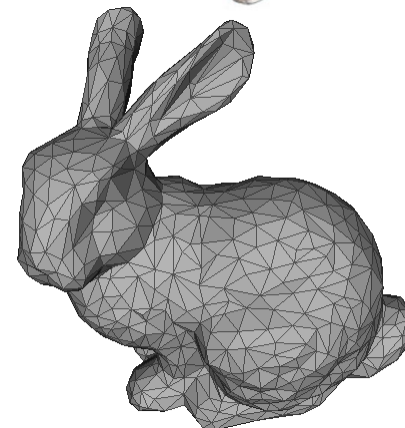
- Size:  $N \times D + F \times 3$  (N: # points, D: space dimension, F: #triangles)
- A 3D surface parametrization (can be higher-dimensional)
  - Piece-wise linear with adaptive detail; triangle faces are usual

## Benefits

- Good for single and multi-view reconstruction
- Provides orientation information (surface normal)
- Graph convolutions possible

## Drawbacks

- Irregular structure (number of neighbors, edge length, face area)
- Difficult to change topology  
(shape changes require to create new vertices and edges)



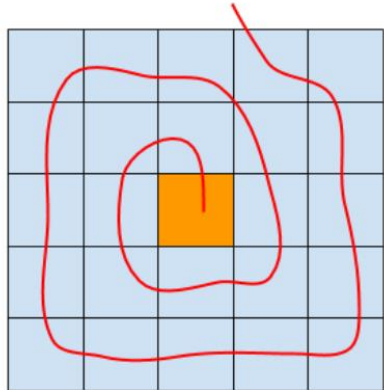
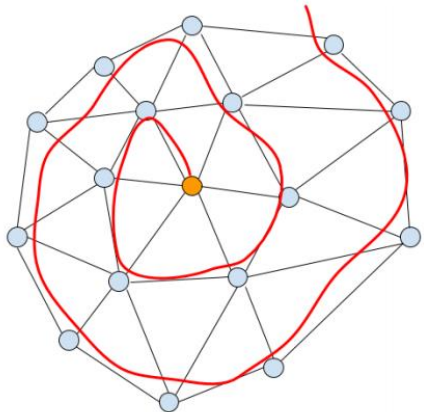
# Spiral convolution

Goal: break the permutation invariance of neighbors

Idea: Order neighbors by simple rules

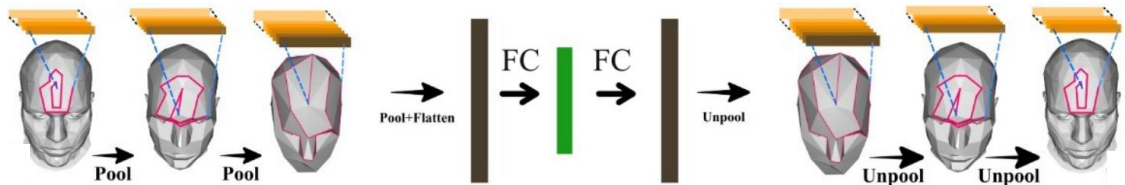
1. collect all neighbors (d hops in the graph)
2. pick the closest one (geodesic distance)
3. continue **counterclockwise** until spiral is of length
4. multiply features h along spiral with weight matrix

$$\mathbf{h}_i^{(l+1)} = \sigma \left( h_{\text{spiral}(\text{neighbors}(i))} W^{(l)} \right)$$



Advantages:

- fixed number of points in each spiral
- efficient to compute
- anisotropic and topology-aware
- easy to optimize



# Details: Mesh Laplacian

**Goal:** A form of 2<sup>nd</sup> order derivative on the mesh

Laplacian for a function in 3D space:

$$\nabla^2 f = \frac{\partial^2 f}{\partial x^2} + \frac{\partial^2 f}{\partial y^2} + \frac{\partial^2 f}{\partial z^2}$$

Difficulty:

- irregularity, where is left / right / up / down?

Solution:

- (weighted) average over all neighboring nodes Ni

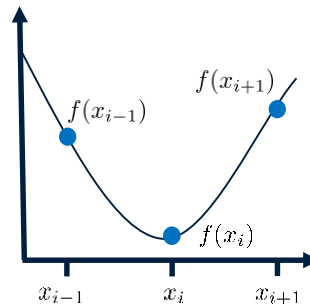
$$\mathcal{L}(\mathbf{v}_i) = \mathbf{v}_i - \frac{1}{d_i} \sum_{j \in \mathcal{N}_i} \mathbf{v}_j$$

- Widely used to encode surface detail and to compare meshes
  - as a loss to compare surfaces

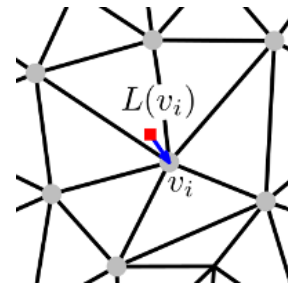
Finite differences approximation in 1D

$$f'(x_i, x_{i+1}) \approx \frac{f(x_{i+1}) - f(x_i)}{h}$$

$$f''(x_{i-1}, x_i, x_{i+1}) \approx \frac{f(x_{i+1}) - 2f(x_i) + f(x_{i-1}))}{h^2}$$



1D Laplacian

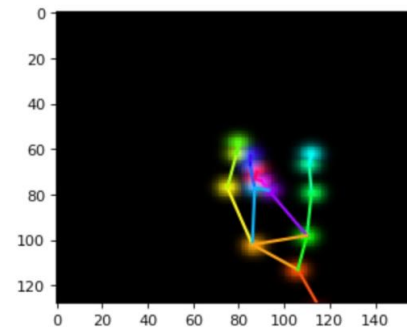
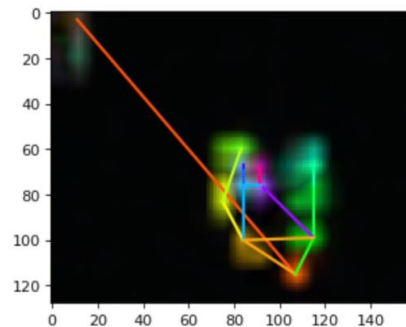
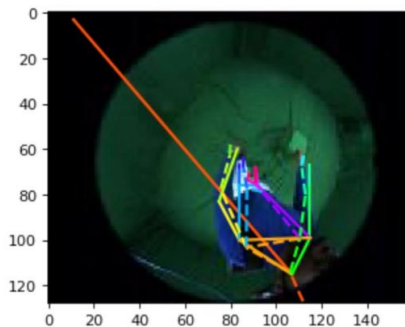


Graph Laplacian



# Assignment 2 discussion

- Issues of heatmap prediction
  - outliers at inference time



- Issues of integral pose regression
  - bias towards the center

For example, if we have the following 1D heatmap

0, 1, 1, 0, 0, 0, 0,

after applying softmax we get the probability map

0.0958, 0.2605, 0.2605, 0.0958, 0.0958, 0.0958, 0.0958,

which leads to predicted position  $0 \times 0.0958 + 1 \times 0.2605 + 2 \times 0.2605 + 3 \times 0.0958 + 4 \times 0.0958 + 5 \times 0.0958 + 6 \times 0.0958 \approx 2.5$

# Assignment 2 discussion II

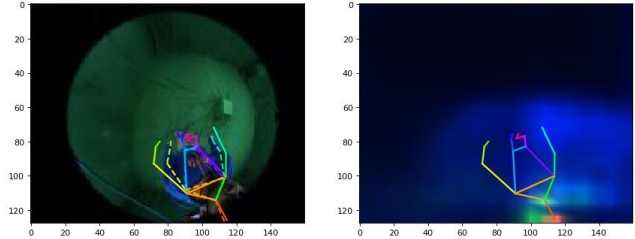
## Mind the numerical stability of soft-max

- a stable implementation was introduced in an earlier lecture

### All black probability map and pose - Task II

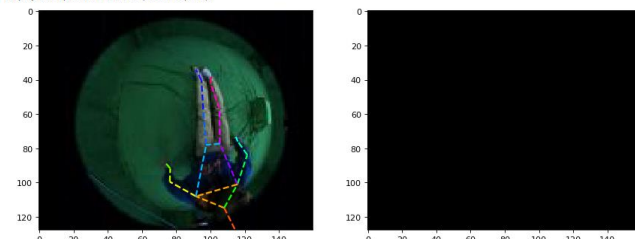
For task II, the training process seems to work fine during the first epoch. However, if I keep training, during the second epoch the predicted probability map and  $\hat{g}$ . Anybody knows what can be the problem?

Displayed output of an iteration (first epoch):



Epoch 0, Iteration 98 of 2823 (3%), loss=0.0017087692394852638

Displayed output of an iteration (second epoch):

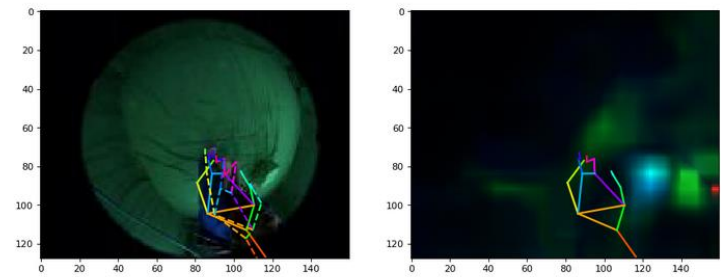


Epoch 1, Iteration 2650 of 2823 (93%), loss=nan

## Dimensions, width and height...

Probability maps look strange!!

For task 2, my predicted poses look quite consistent with the reference, and loss is always less than 0.003. However, the probability maps I



Epoch 0, iteration 2820 of 2823 (99%), loss=0.0017382814548909664

*Any other issues?*

# Debugging best practice!

1. Basic principle: garbage in, garbage out
  - make sure your input has the correct type
    - correct tensor dimension, correct order of dimension, correct values, ...
  - if it is an image or matrix, plot it
  - if you deal with points, plot them

2. How do I determine whether my input/output values are correct?

- read the specification (e.g., assignment)
- if there is no specification, write one
- toy examples where you know the correct behavior
  - e.g., a single object, single color, primitive shape
  - try to separate influence factors, such as scale and shape
    - e.g., two images with the same shape but different scale

3. My input is correct, but the output is wrong, how do I find the bug?

***Wolf fence algorithm*** by Edward Gauss:

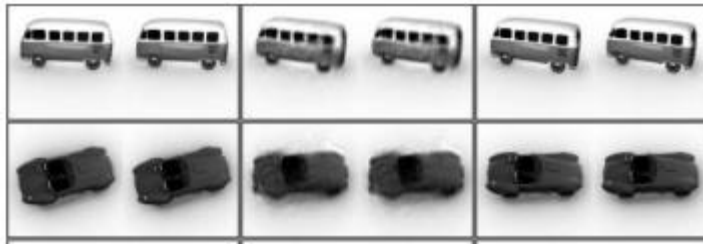
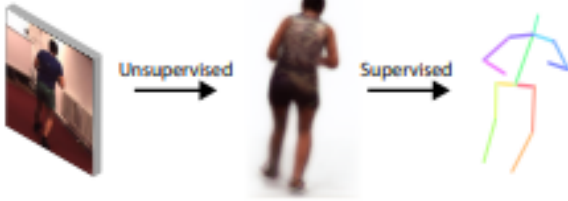
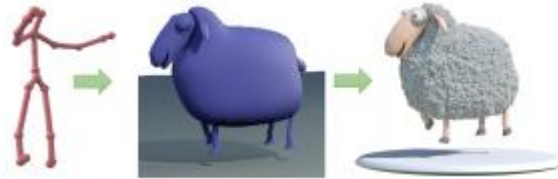
There's one wolf in Alaska; how do you find it?

- build a fence down the middle of the state, wait for the wolf to howl, determine which side of the fence it is on (point 1&2).
- Repeat process on that side only, until you can see the wolf.

# PCA and AEs will be important for the paper reading!

Principal Component Analysis (PCA) and Auto Encoder (AE)

**Interpretable Shape Representation**



**Abstraction I**

Spatial layout (bounding boxes & depth)

**Abstraction II**

Instance segmentation and depth maps

**Abstraction III**

Encoding and novel view decoding



# Principal Component Analysis (PCA)

# Recap: Principal component analysis overview

- The orthogonal linear transformation that transforms the data to a new coordinate system such that the greatest variance by some scalar projection of the data comes to lie on the first coordinate

- First weight vector  $w(1)$

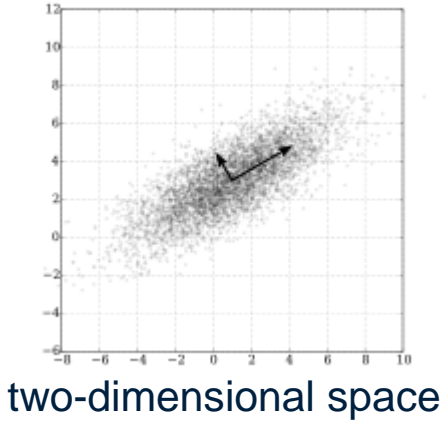
$$\mathbf{w}_{(1)} = \arg \max_{\|\mathbf{w}\|=1} \left\{ \sum_i (\mathbf{x}_{(i)} \cdot \mathbf{w})^2 \right\}$$

computed over all  $x(i)$  in the dataset

- ... continue iteratively in orthogonal directions
- Stacking all weight vectors **as rows** into a matrix  $W$  yields a 'linear auto encoder'

$$\hat{p} = \overbrace{WW^T}^{\text{reconstruction}} p$$

reconstruction (decoding)      projection (encoding)

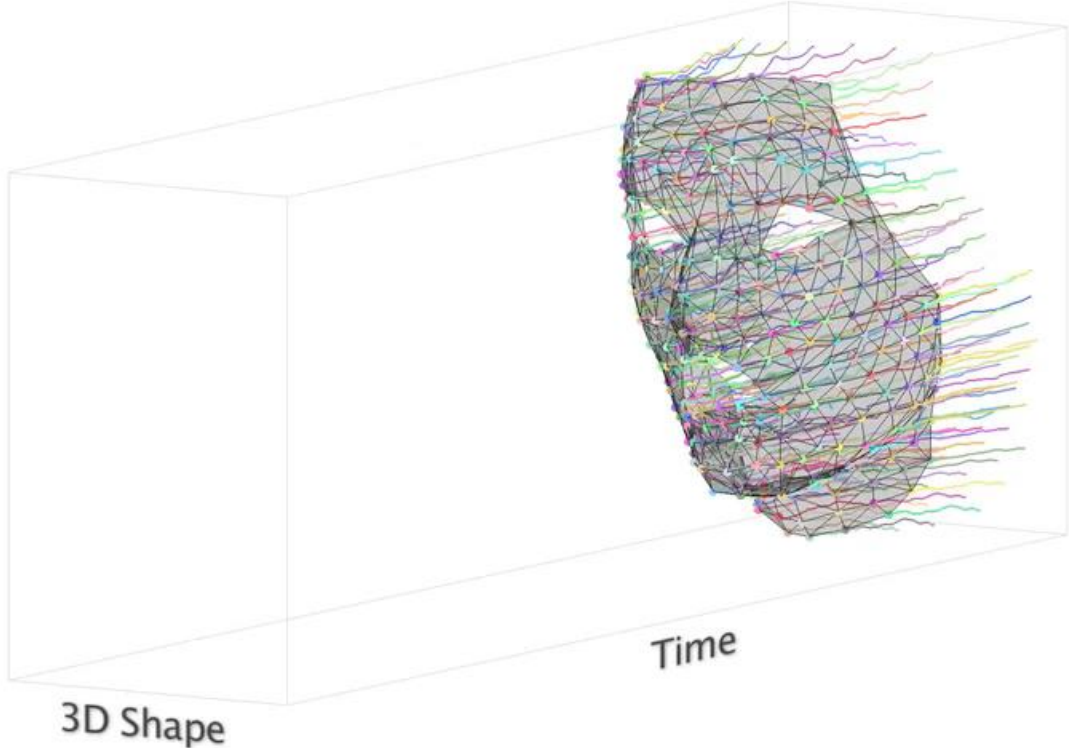


# PCA-like body model



[Movie Reshape]

# PCA space: time or space?





# Data matrix

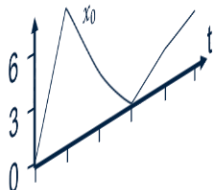
The data matrix  $X$  encodes

- each row represents a new measurement
- each column represents the

**Measurement at time  $t_0$**

	$x_0$	$y_0$	$z_0$	$x_1$	$y_1$	$z_1$
$t_0$	1	7	3	1	2	3
$t_1$	8	2	4	4	5	1
$t_2$	1	9	7	4	5	9
$t_3$	2	9	4	4	5	2
$t_4$	1	4	4	2	6	2
$t_5$	1	9	2	1	5	5

$X =$



**Motion of coordinate  $x_0$**



# The covariance matrix

The empirical sample covariance matrix of  $\mathbf{X}^T$

$$\frac{1}{1-n} \mathbf{X}^T \mathbf{X}$$

We consider its unscaled form

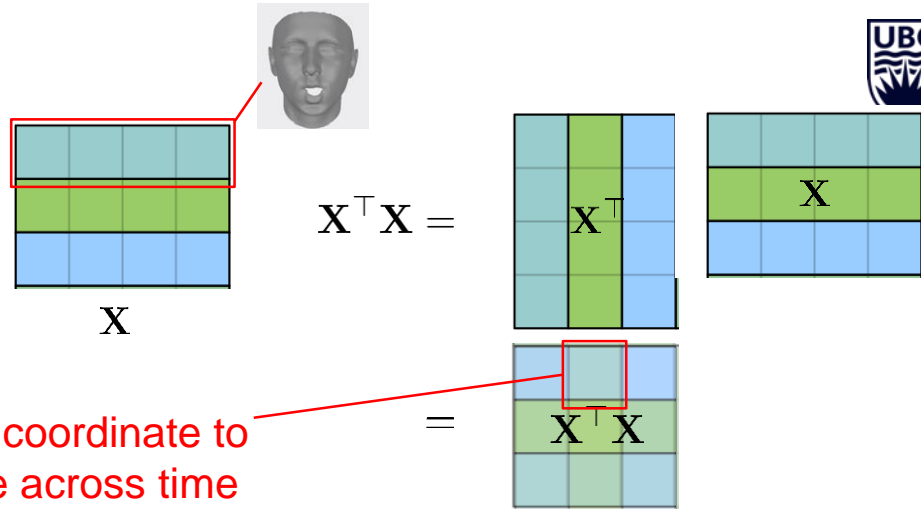
$$\mathbf{X}^T \mathbf{X}$$

This symmetric matrix can be decomposed  
(Eigendecomposition)

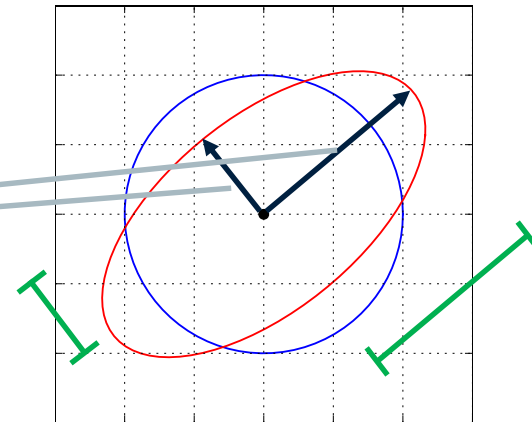
$$\mathbf{X}^T \mathbf{X} = \mathbf{W} \mathbf{\Lambda} \mathbf{W}^T$$

Eigenvalues

Eigenvectors

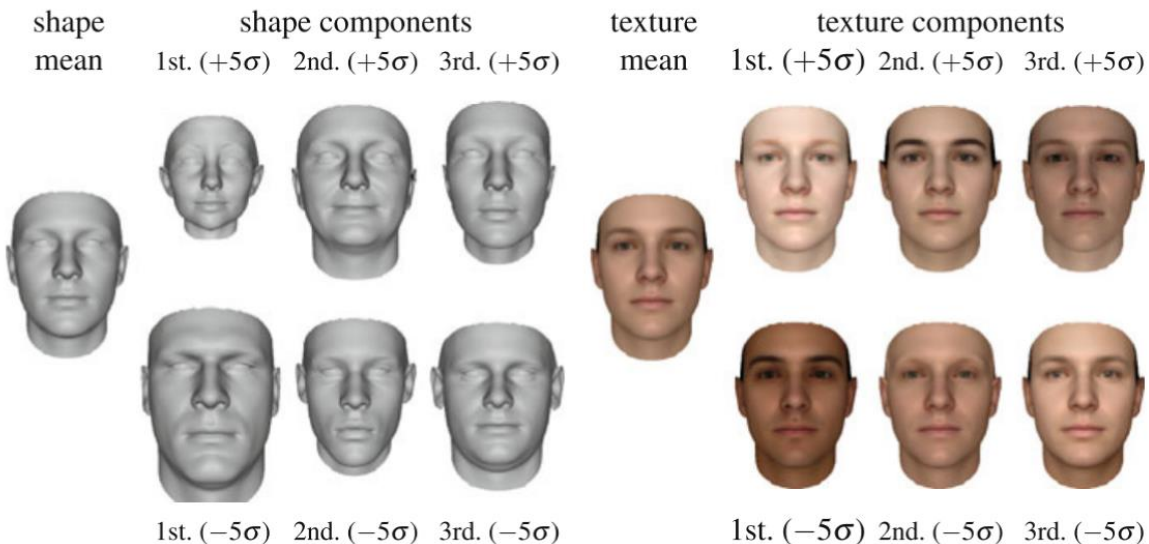


Sphere-ellipse representation  
(how are points on a sphere deformed by  $\mathbf{X}^T \mathbf{X}$ )



# Spatial components

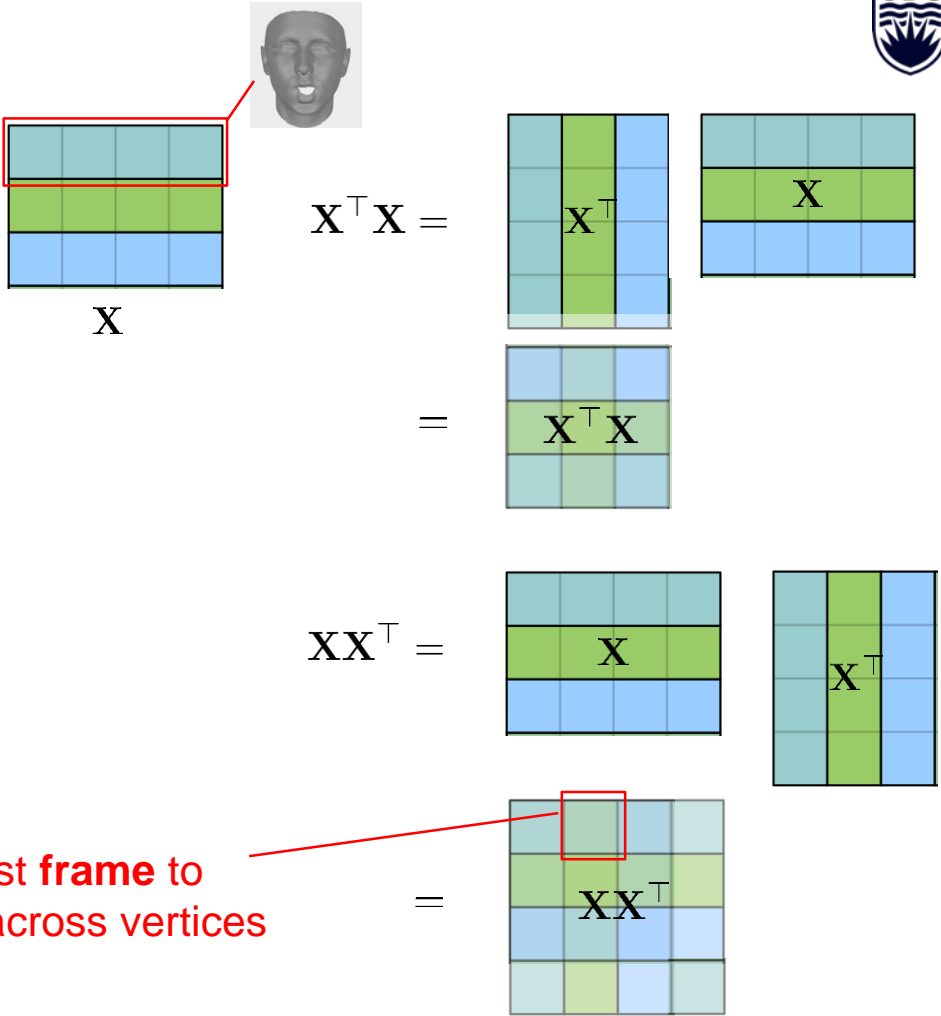
- each component captures
- points that 'move' together
  - together = correlated
  - move = change across different measurements
  - e.g. left and right side
- if the input motion is smooth,
- it will lead to a smooth shape basis
  - global: scale ,male-female
  - forehead wrinkles in one basis
- works also on textures



[Blanz and Fetter, A morphable model for the synthesis of 3D faces. 1999]

# The covariance matrix II

Exchanging the role of rows and columns



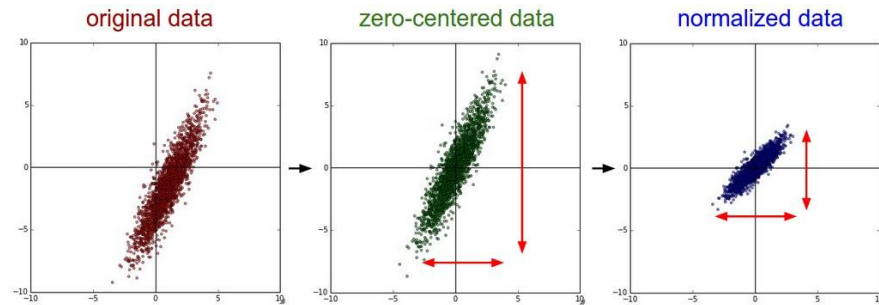
Relation of first **frame** to the second one across vertices

# Recall from lecture 3: Input and output normalization

Goal: Normalize input and output variables to have  $\mu=0$  and  $\sigma=1$

$$\tilde{\mathbf{x}} = \frac{\mathbf{x} - \mu}{\sigma}$$

- For an image, normalize each pixel by the std and mean color (averaged over the **training set**)



<http://cs231n.github.io/neural-networks-2/>

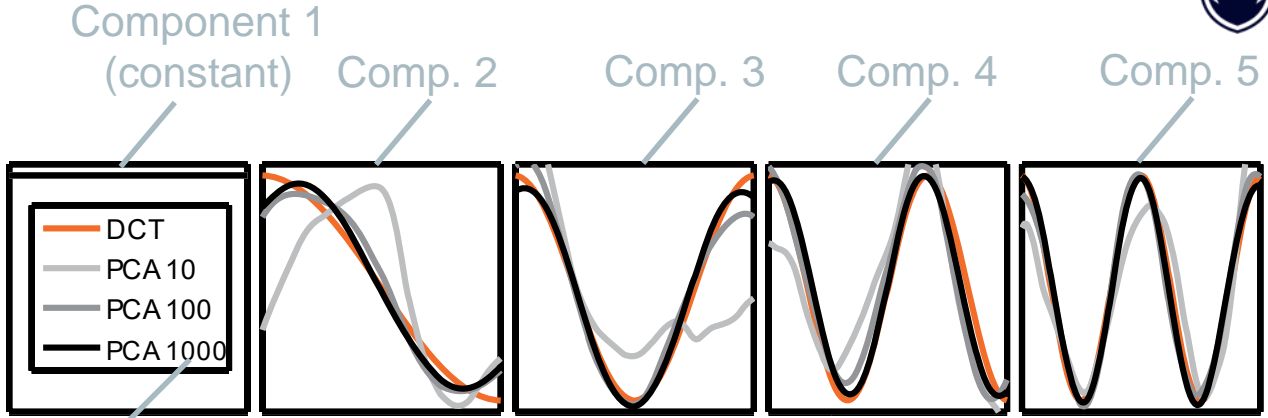
Related to data whitening

- whitening transforms a random vector to have zero mean and unit diagonal covariance
- by contrast, the default normalization for deep learning is element wise, neglecting dependency
  - the resulting covariance is not diagonal!

*This is what we can do with PCA, it's a rotation and scaling of the data*

# Trajectory basis

- smooth input motions lead to a smooth trajectory basis
  - approximates DCT for increasing number of 'training' sequences



PCA with 10,100,1000 training sequences

[Akhter et al., Bilinear Spatiotemporal Basis Models. 2012]

## Fourier transform / Discrete cosine transform (DCT)

- a change of basis
  - orthogonal basis
  - turn a function of time into a function of frequency

# Singular value decomposition (SVD)

Decomposition of the data matrix  $X$  with SVD

$$X = U \Sigma W^T$$

← orthonormal (unit length and linearly independent columns)  
← matrix of eigenvectors  
↑ diagonal matrix of singular values

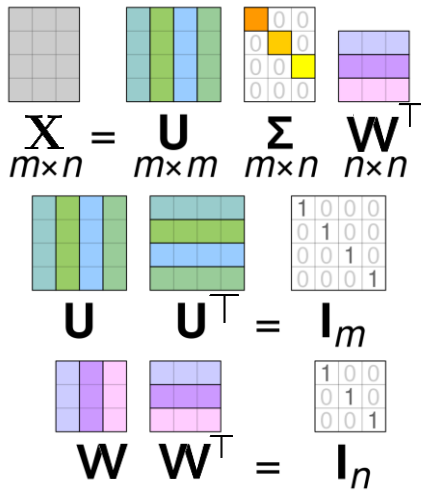
- singular values are arranged in descending order (makes SVD unique)
- closely related to PCA:

$$X^T X = W \Sigma^T U^T U \Sigma W^T = W \Sigma^T \Sigma W^T$$

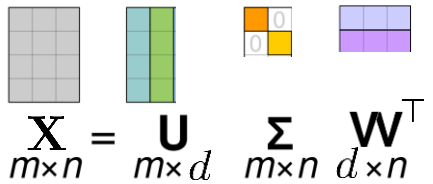
trajectory basis

shape basis

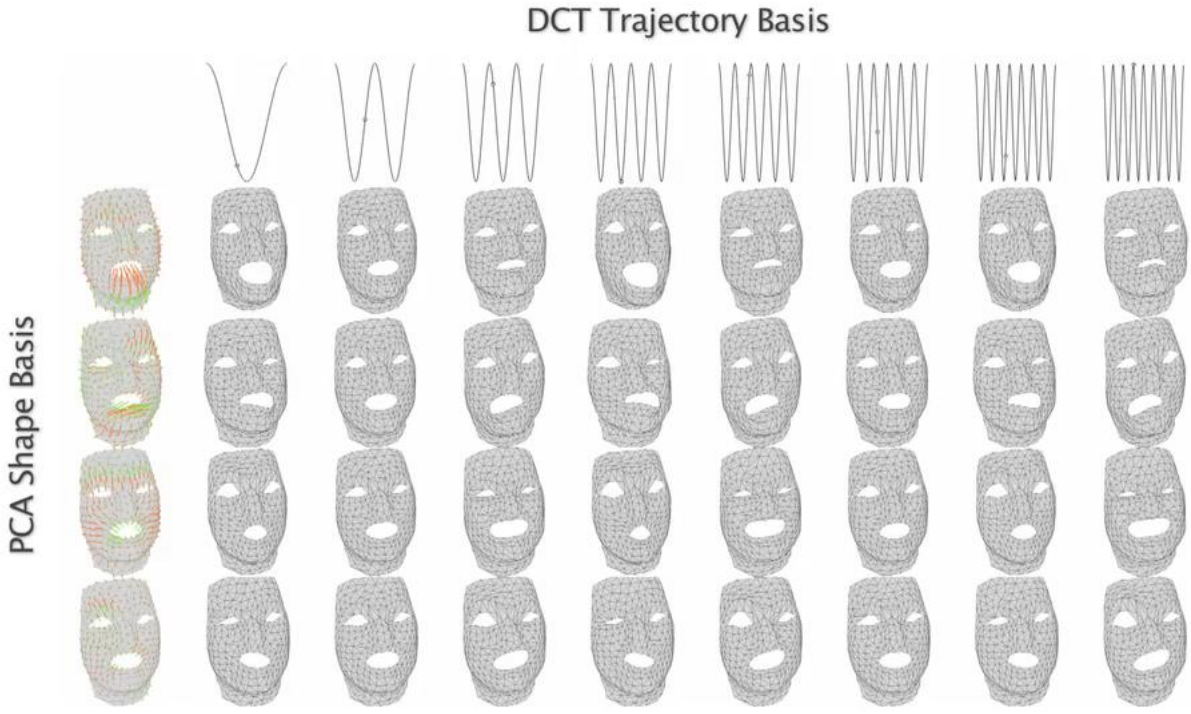
## SVD in matrix form



## Dimensionality reduction



# Extension: Bilinear model



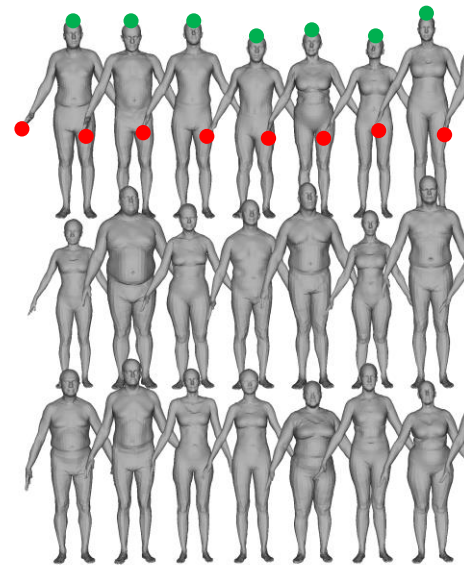
[Akhter et al., Bilinear Spatiotemporal Basis Models. 2012]



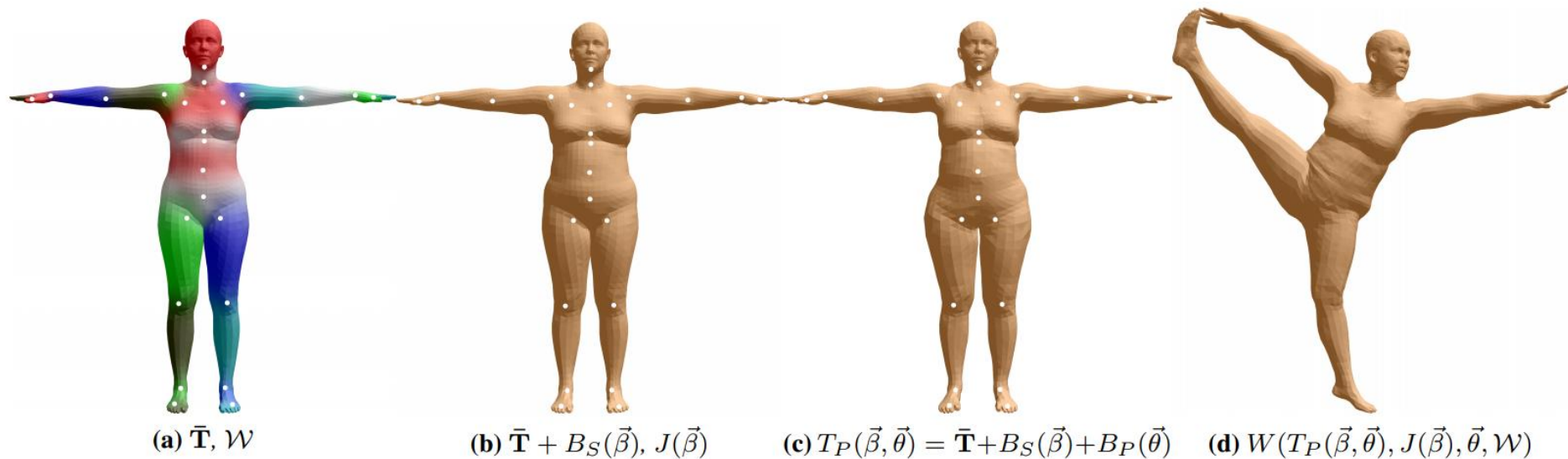
# PCA - correspondence

PCA requires multiple ‘measurements’ of the same quantity

- e.g., for a human mesh model:
  - same number of vertices in mesh
  - the same vertex must correspond to the same semantic position. E.g., vertex 612 is the nose
- holes (missing data) is not supported
  - inappropriate for monocular reconstructions, e.g., where the back of the person is missing
  - generalizations exist to address this case
- scale sensitive
  - estimates those components that maximize variance
    - facial details are outweighed by belly shape
      - for human perception the face is important!

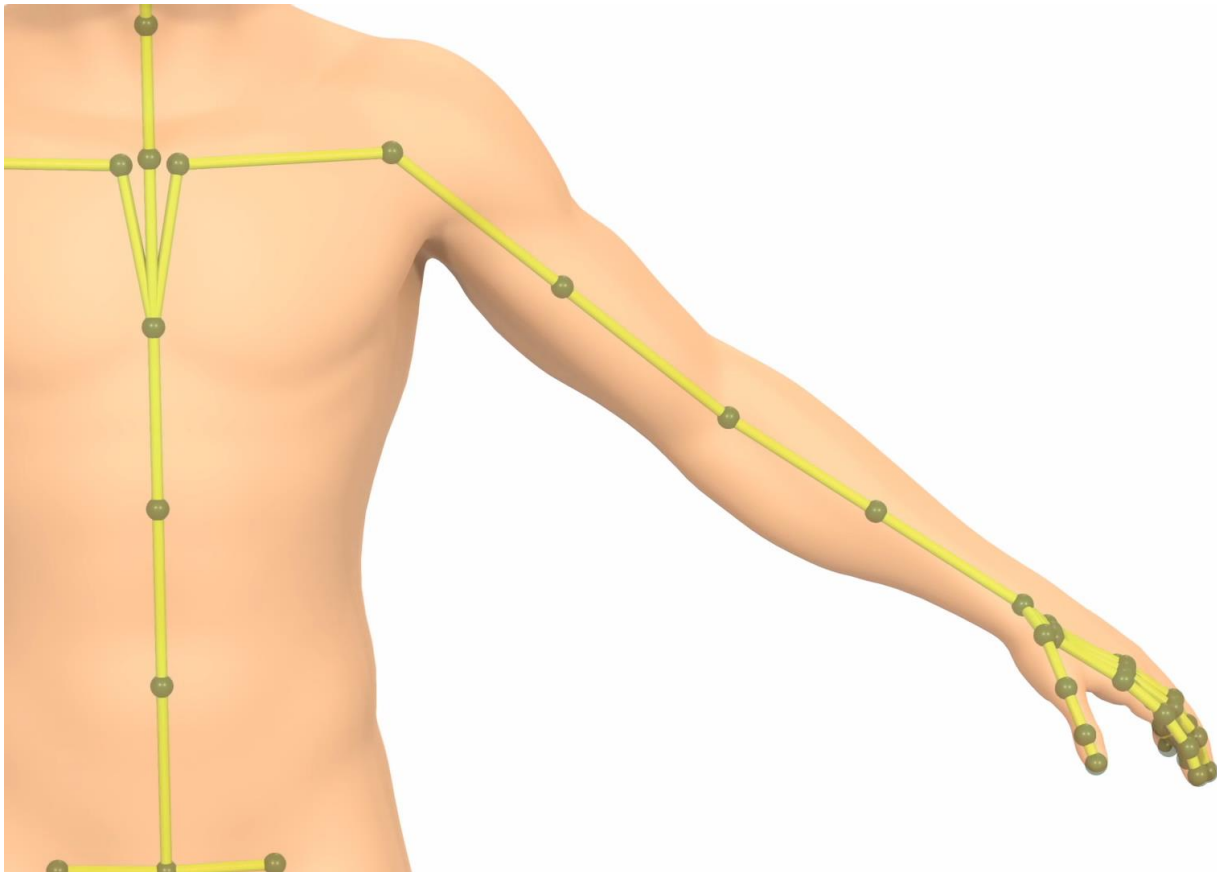


# Recap: SMPL: A Skinned Multi-Person Linear Model



**Figure 3: SMPL model.** (a) Template mesh with blend weights indicated by color and joints shown in white. (b) With identity-driven blendshape contribution only; vertex and joint locations are linear in shape vector  $\vec{\beta}$ . (c) With the addition of pose blend shapes in preparation for the split pose; note the expansion of the hips. (d) Deformed vertices reposed by dual quaternion skinning for the split pose.

# Recap: Skinning



# Auto Encoder (AE)

## General case

$$\mathbf{h} = \text{encoder}_{\theta}(\mathbf{x})$$

$$\mathbf{x}' = \text{decoder}_{\theta}(\mathbf{h})$$

## Simple non-linear case

$$\mathbf{h} = \sigma(\mathbf{W}\mathbf{x} + \mathbf{b})$$

$$\mathbf{x}' = \sigma(\mathbf{W}'\mathbf{h} + \mathbf{b}')$$

## Linear case

$$\mathbf{h} = \mathbf{W}\mathbf{x} + \mathbf{b}$$

$$\mathbf{x}' = \mathbf{W}'\mathbf{h} + \mathbf{b}'$$

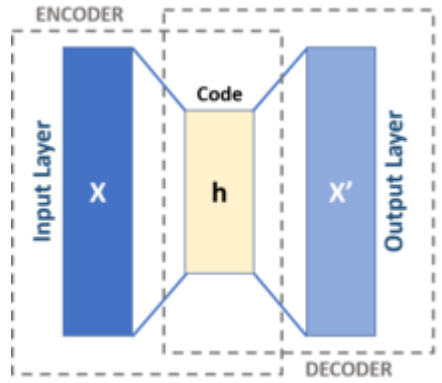
## General reconstruction objective

$$\text{loss}(\mathbf{x}, \mathbf{x}')$$

- e.g., MSE loss

*A two-layer fully-connected neural network*

*Similar to PCA when using squared loss (W spans the same space, but neither forms an ordered nor orthogonal basis)*



<https://en.wikipedia.org/wiki/Autoencoder>

## Linear autoencoder objective

$$\begin{aligned} & \arg \min_{\mathbf{W}} \sum_i \|\mathbf{x} - \mathbf{x}'\|^2 \\ & = \arg \min_{\mathbf{W}} \sum_i \|\mathbf{x}_{(i)} - \mathbf{W}'\mathbf{W}\mathbf{x}_{(i)}\|^2 \end{aligned}$$

## PCA objective

$$\begin{aligned} \mathbf{w}_{(1)} &= \arg \max_{\|\mathbf{w}\|=1} \left\{ \sum_i (\mathbf{x}_{(i)} \cdot \mathbf{w})^2 \right\} \\ &= \arg \max_{\|\mathbf{w}\|=1} \left\{ \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} \right\} \end{aligned}$$

# Autoencoder variants

## Bottleneck autoencoder:

- hidden dimension smaller than input dimension
  - leads to compressed representations
  - like dimensionality reduction with PCA

## Sparse autoencoder:

- hidden dimension larger than input dimension
- hidden activation enforced to be sparse  
(=few activations)

## Denoising autoencoder:

- corrupt the input values, e.g. by additive noise

$$\mathbf{h} = \text{encoder}_{\theta}(\text{noise}(\mathbf{x}))$$

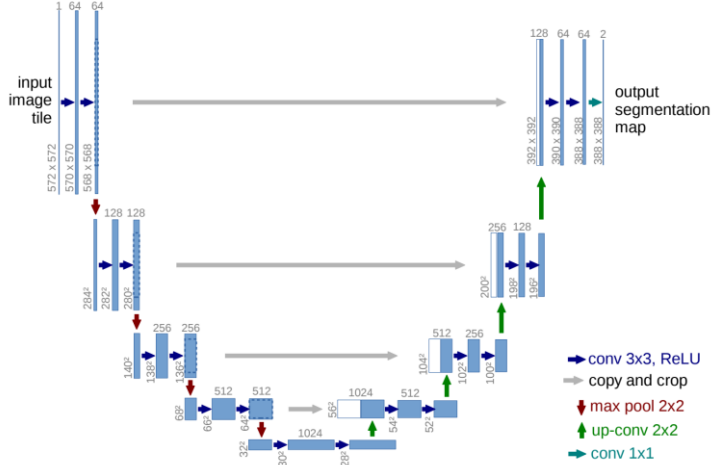
$$\mathbf{x}' = \text{decoder}_{\theta}(\mathbf{h})$$

## Variational Auto Encoder (VAE)

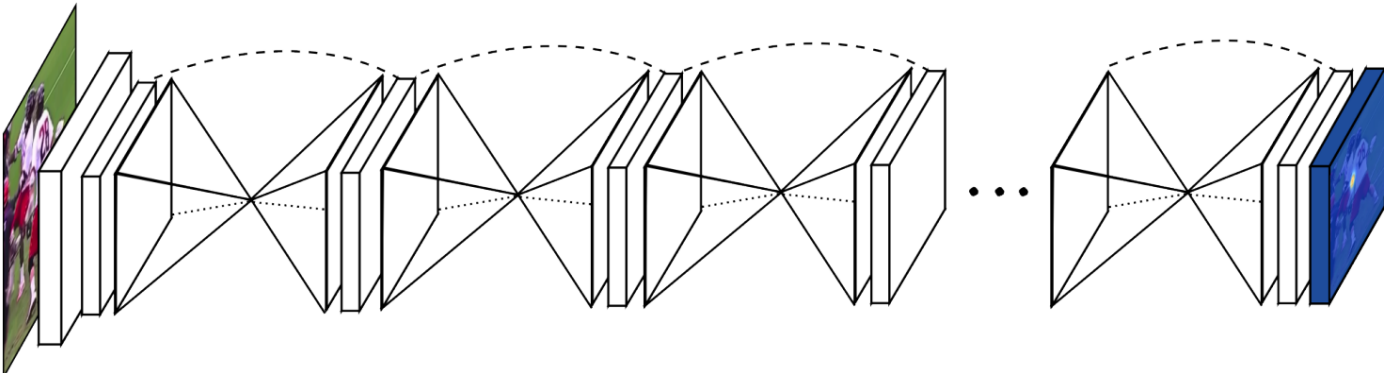
- a probabilistic model
  - *'adding noise on the hidden variables'*
  - more in lecture 8!

# Relation to previous lectures

- The UNet has an encoder-decoder structure



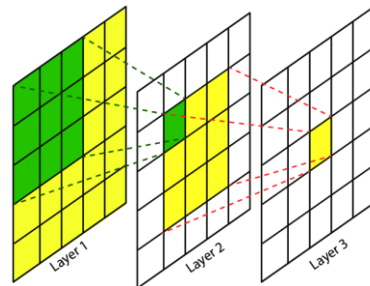
- The stacked hourglass network applies multiple encoders and decoders



## Preparation for Assignment 3

- Will be posted tonight or tomorrow
  - PyTorch issues encountered
    - circular convolutions are broken...
  - the current version is a bit boring

# Feature map size after convolutional kernels



## Transformation of input and output by convolutions

- output size =  $(\text{input size} + 2 * \text{padding} - \text{kernel size} + \text{stride}) / \text{stride}$ 
  - e.g., a 3x3 kernel that preserves size:  $W + 2 * 1 - 3 + 1 = W$
  - e.g., a 4x4 kernel that reduces size by factor two:  $(W + 2 * 1 - 4 + 2) / 2 = W / 2$
- holds per dimension, i.e., 1D, 2D and 3D convolutions

## Transformation of input and output by **transposed** convolutions (aka. deconvolution)

- output size =  $\text{input size} * \text{stride} - \text{stride} + \text{kernel size} - 2 * \text{padding}$ 
  - it has exactly the opposite effect of convolution
  - e.g., a 3x3 kernel that preserves size:  $W - 1 + 3 - 2 * 1 = W$
  - e.g., a 4x4 kernel that **increases** size by **factor** two:  $W * 2 + 2 * 1 - 4 + 2 = W * 2$
  - e.g., a 3x3 kernel that **increases** size by two elements:  $W - 1 + 3 - 2 * 0 = W + 2$





Anyone missing who send their choice?

# Presentation topic assignment ongoing

Summary: 19 votes for 22 papers

- Gives 3 late votes or remaining slots?
- Remaining papers can be presented by auditing students
  - volunteers?

Dingqing

Ege Unlu

Dave Pagurek van Mossel

Shuxian Fan

Shelly C

Peyman Bateni

Tim Straubinger

Shenyi Pan

Michela Minerva - michela

Jerry Yin

Willis Peng

Michelle Appel

stolet

Zicong Fan (Alex)

fjavadi

ssims

Daniele Reda

Shih-Han Chou

Weidong Yin

# Project

## Project proposal

- 3-minute pitch per group
- written plan
  - one page, 11pt font, may include figures
    - not more than one, not less than half a page of text
- the proposal plan must cover
  - the research idea
  - the possible algorithmic contributions
  - and an outline of the evaluation
- get feedback from during office hours
  - Yuchi on Tuesdays
  - me on Wednesdays
  - **Only three weeks left!**

# Hidden questions

1. What is the difference between a strong and a weak hypothesis?

2. What is the difference between a strong and a weak hypothesis?

3. What is the difference between a strong and a weak hypothesis?