

## COSC 522 – Machine Learning

### Lecture 16 – From Machine Learning to Deep Learning

Hairong Qi, Gonzalez Family Professor  
Electrical Engineering and Computer Science  
University of Tennessee, Knoxville

<https://www.eecs.utk.edu/people/hairong-qi/>

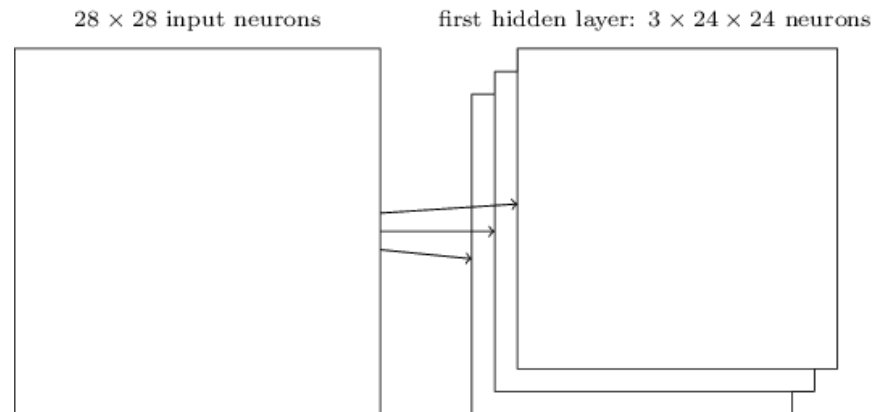
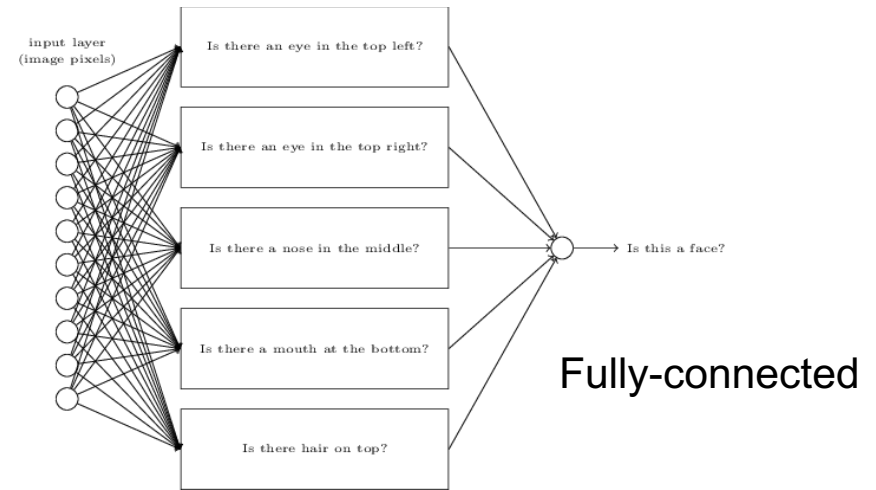
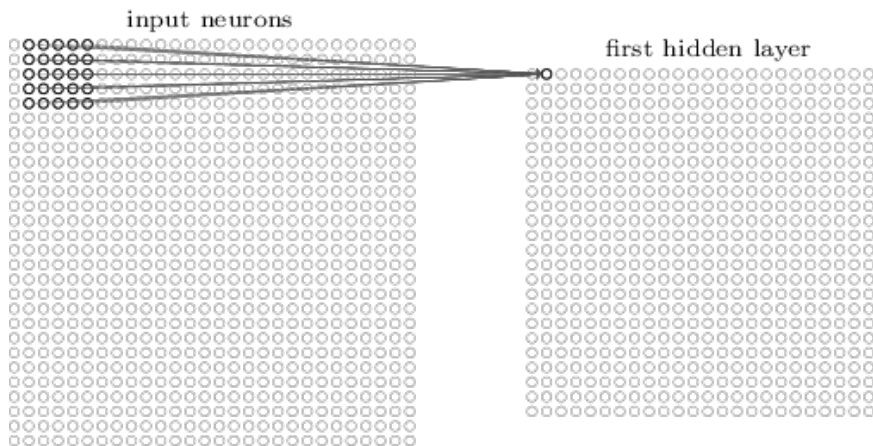
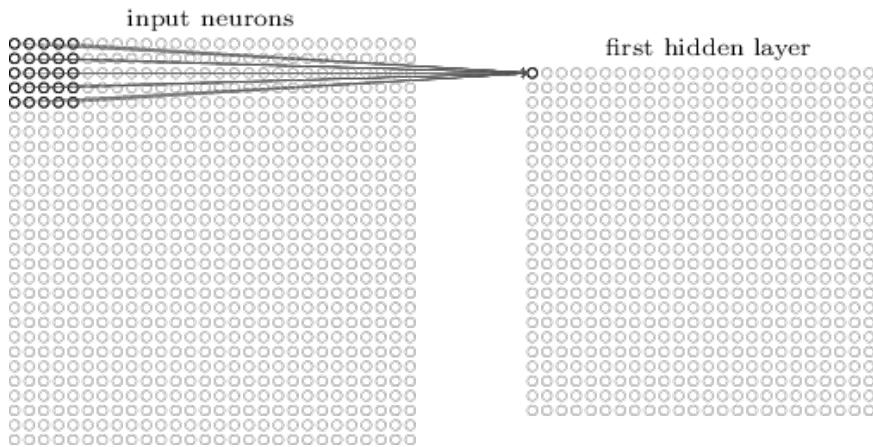
Email: [hqi@utk.edu](mailto:hqi@utk.edu)

Course Website: <http://web.eecs.utk.edu/~hqi/cosc522/>

# A list of misconceptions

- Is deep learning merely deeper?
  - The two unique features of convolutional neural network (CNN)
- Is deep learning a classifier?
  - Engineered features vs. automatic features
- Supervised vs. Unsupervised
- Model-based approach vs. Data-driven approach
  - the two extremes?
- The world beyond CNN
  - GAN, AE, RNN, RL
- Implementation

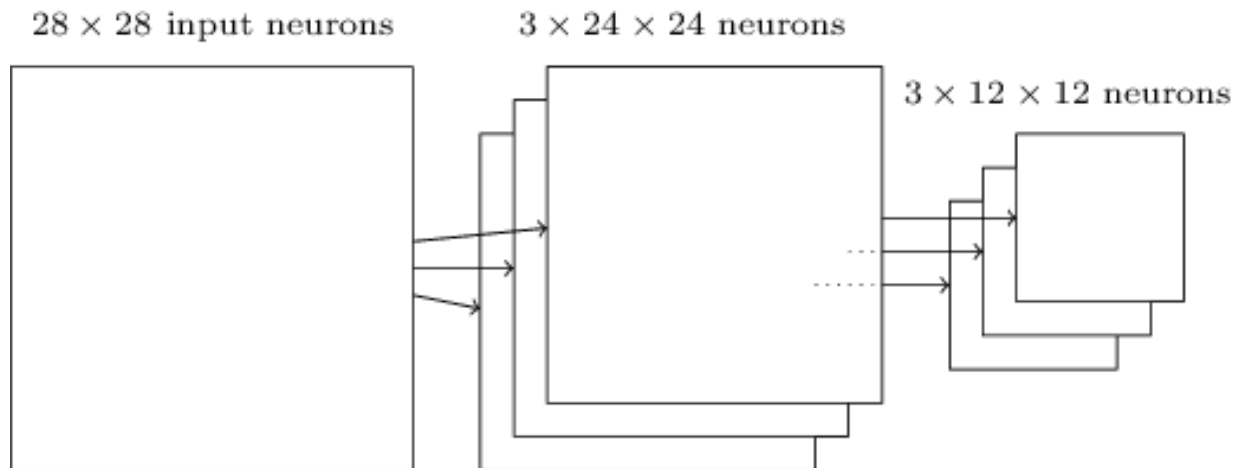
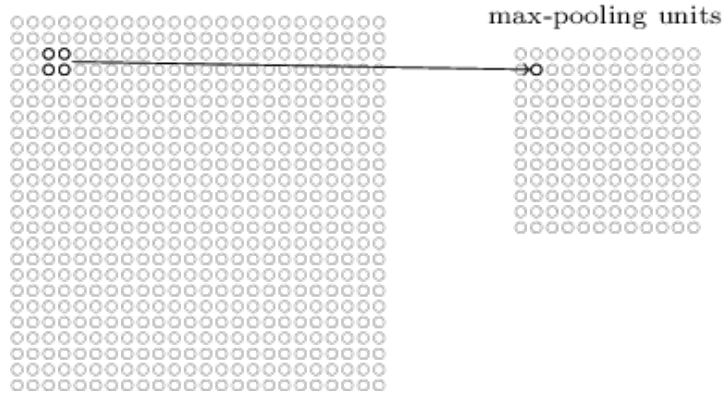
# Core idea 1: Receptive field (RF) and shared weight



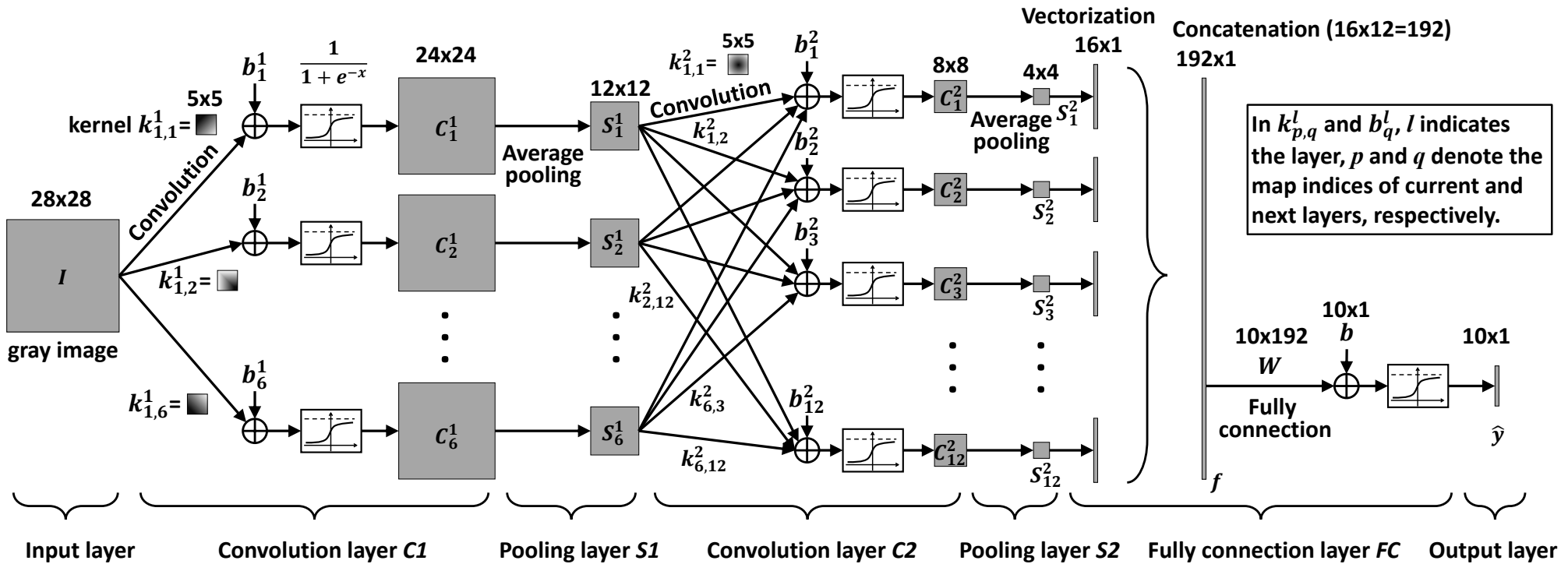
Feature maps

# Core idea 2: Hierarchical vision - Max pooling

hidden neurons (output from feature map)



# A simple CNN framework

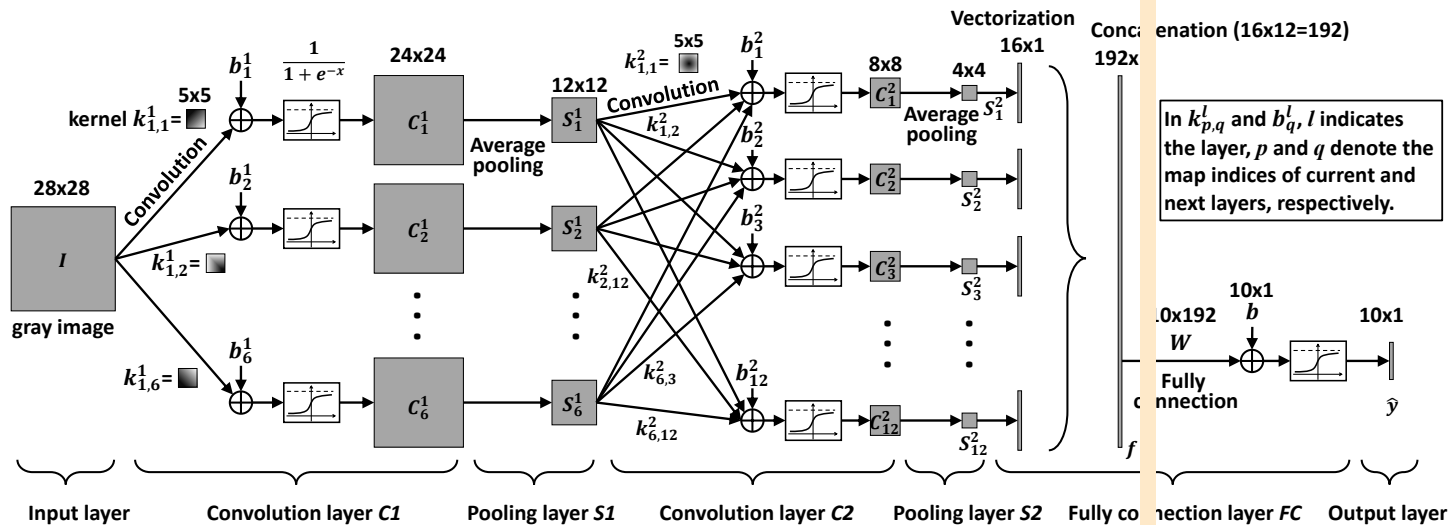
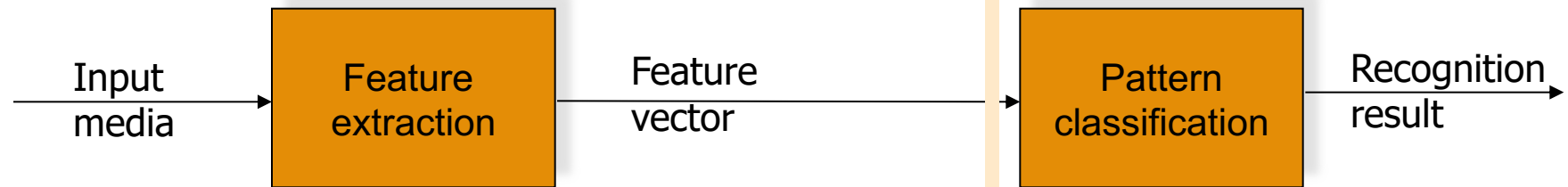


# A list of misconceptions

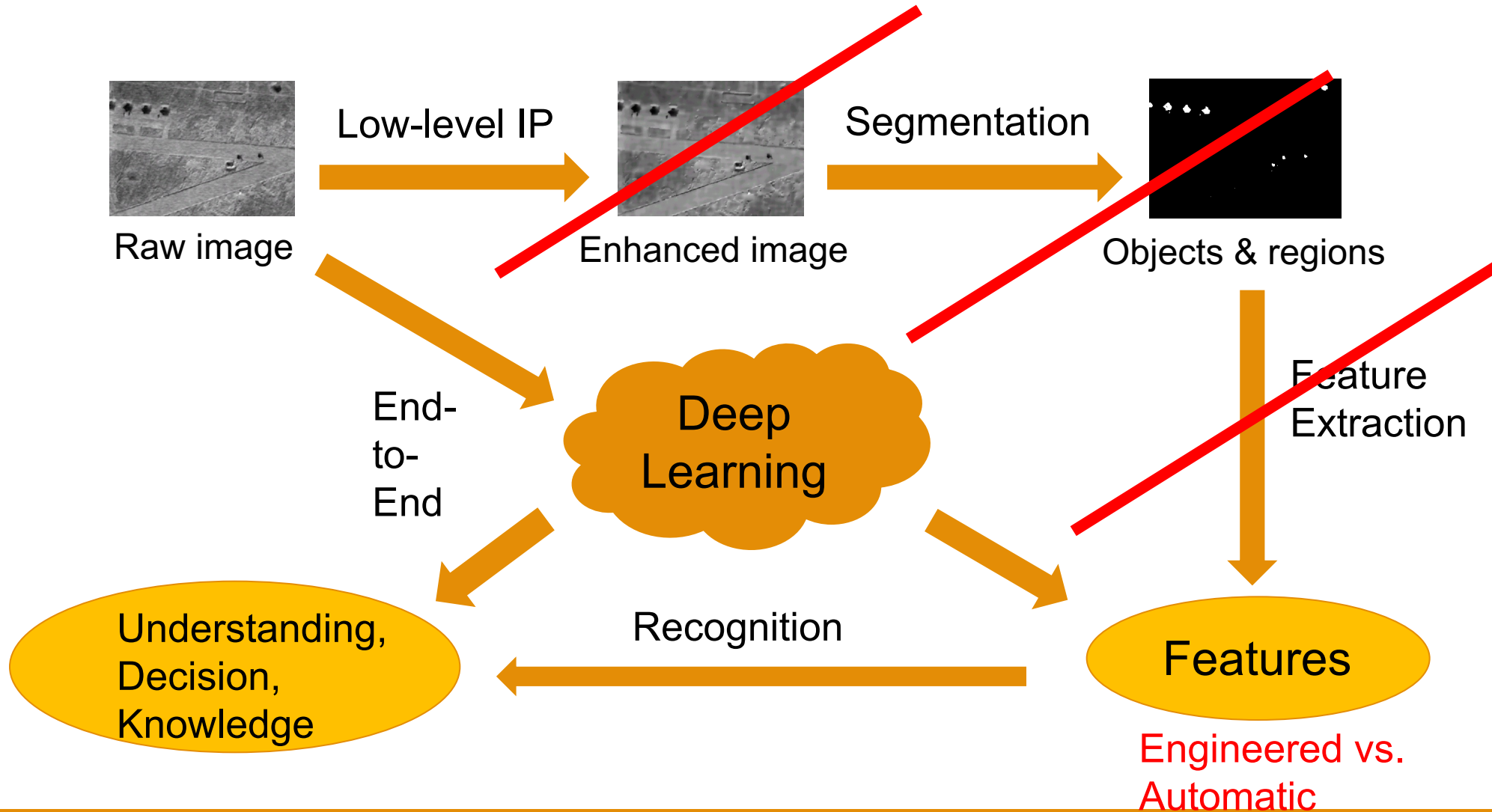
- Is deep learning merely deeper?
  - The two unique features of convolutional neural network (CNN)
- **Is deep learning a classifier?**
  - **Engineered features vs. automatic features**
- Supervised vs. Unsupervised
- Model-based approach vs. Data-driven approach
  - the two extremes?
- The world beyond CNN
  - GAN, AE, RNN, RL
- Implementation

# Engineered features vs. automatic features

Need domain knowledge



# The flowchart comparison





# A list of misconceptions

- Is deep learning merely deeper?
  - The two unique features of convolutional neural network (CNN)
- Is deep learning a classifier?
  - Engineered features vs. automatic features
- **Supervised vs. Unsupervised**
- Model-based approach vs. Data-driven approach
  - the two extremes?
- The world beyond CNN
  - GAN, AE, RNN, RL
- Implementation

# Revisit: A bit of history

- 1956-1976
  - 1956, The Dartmouth Summer Research Project on Artificial Intelligence, organized by John McCarthy, Marvin Minsky, Nathaniel Rochester, and Claude Shannon

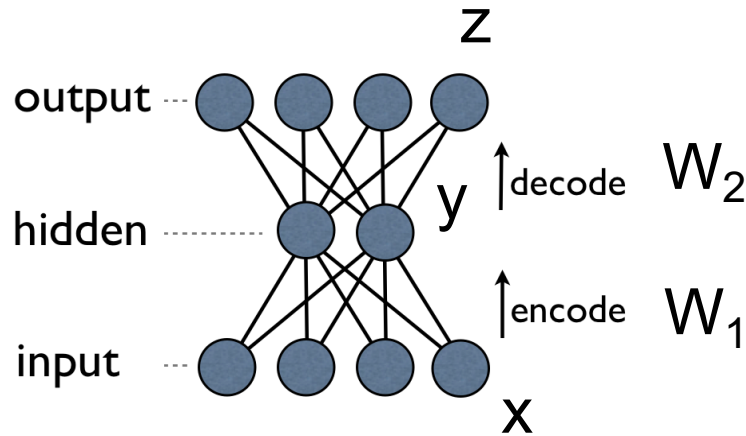
We propose that a 2 month, 10 man study of artificial intelligence be carried out during the summer of 1956 at Dartmouth College ... The study is to proceed on the basis of the conjecture that every aspect of learning or any other feature of intelligence can in principle be so precisely described that a machine can be made to simulate it. An attempt will be made to find how to make machines use language, form abstractions and concepts, solve kinds of problems now reserved for humans, and improve themselves. We think that a significant advance can be made in one or more of these problems if a carefully selected group of scientists work on it together for a summer.

- The rise of symbolic methods, systems focused on limited domains, deductive vs. inductive systems
- 1973, the Lighthill report by James Lighthill, “Artificial Intelligence: A General Survey” - automata, robotics, neural network
- 1976, the AI Winter
- 1976-2006
  - 1986, BP algorithm
  - ~1995, The Fifth Generation Computer
- 2006-???
- 2006, Hinton (U. of Toronto), Bingio (U. of Montreal), LeCun (NYU)
- 2012, ImageNet by Fei-Fei Li (2010-2017) and AlexNet

[https://en.wikipedia.org/wiki/Dartmouth\\_workshop](https://en.wikipedia.org/wiki/Dartmouth_workshop)

[https://en.wikipedia.org/wiki/Lighthill\\_report](https://en.wikipedia.org/wiki/Lighthill_report)

# Unsupervised learning – Autoencoder (AE)



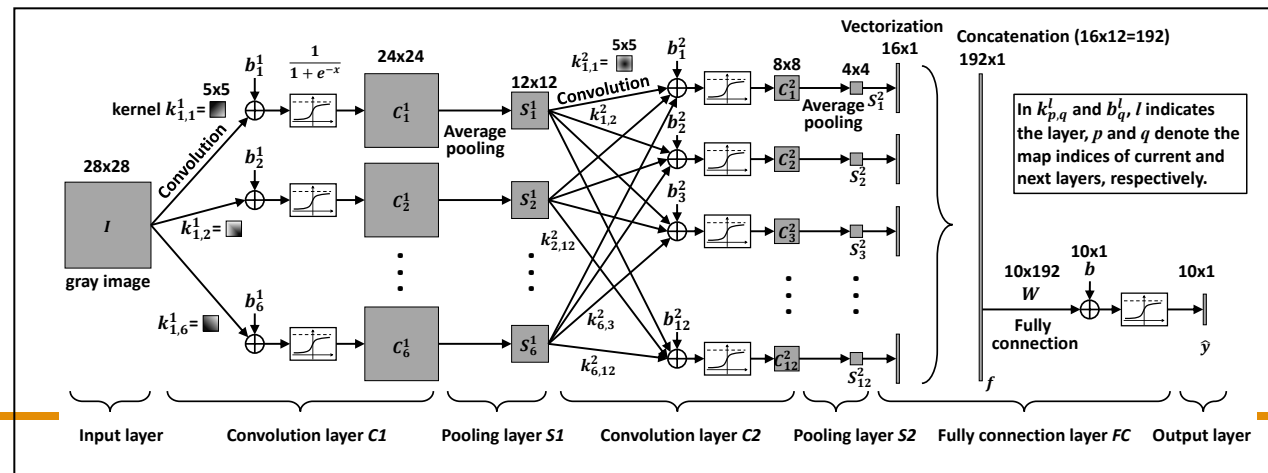
$$y = f_{\theta_1}(\mathbf{W}_1 \mathbf{x} + \mathbf{b}_1)$$

$$z = g_{\theta_2}(\mathbf{W}_2 \mathbf{y} + \mathbf{b}_2)$$

$$\theta_1 = \{\mathbf{W}_1, \mathbf{b}_1\}, \theta_2 = \{\mathbf{W}_2, \mathbf{b}_2\}$$

$$\theta_1^*, \theta_2^* = \arg \min_{\theta_1, \theta_2} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}^{(i)}, \mathbf{z}^{(i)}) = \arg \min_{\theta_1, \theta_2} \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}^{(i)}, g_{\theta_2}(f_{\theta_1}(\mathbf{x}^{(i)})))$$

$$L(\mathbf{x}, \mathbf{z}) = \frac{1}{2} \|\mathbf{x} - \mathbf{z}\|_2^2$$



# Discrimination vs. Representation of Data

- Best discriminating the data

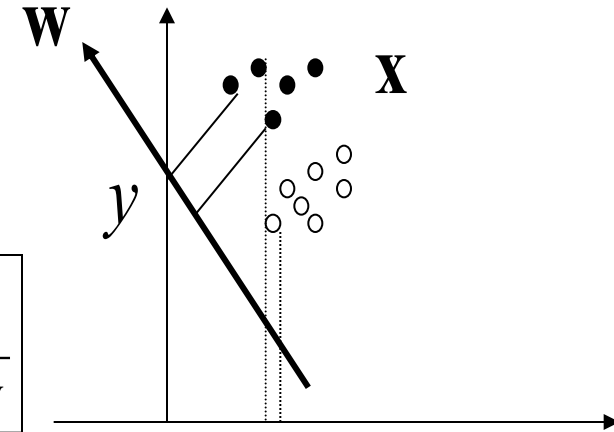
- Fisher's linear discriminant

(FLD)

- NN

- CNN

$$J(\mathbf{w}) = \frac{|\tilde{m}_1 - \tilde{m}_2|^2}{\tilde{s}_1^2 + \tilde{s}_2^2} = \frac{|\mathbf{w}^T (\mathbf{m}_1 - \mathbf{m}_2)|^2}{\mathbf{w}^T \mathbf{S}_1 \mathbf{w} + \mathbf{w}^T \mathbf{S}_2 \mathbf{w}} = \frac{\mathbf{w}^T \mathbf{S}_B \mathbf{w}}{\mathbf{w}^T \mathbf{S}_W \mathbf{w}}$$

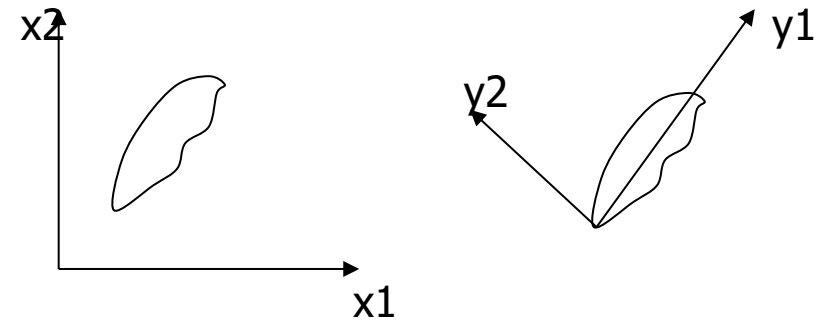


- Best representing the data

- Principal component analysis

(PCA)

$$\mathbf{x} = \sum_{i=1}^m y_i \mathbf{b}_i + \sum_{i=m+1}^d y_i \mathbf{b}_i \approx \sum_{i=1}^m y_i \mathbf{b}_i + \sum_{i=m+1}^d \alpha_i \mathbf{b}_i$$



Error: 
$$\Delta \mathbf{x} = \sum_{i=m+1}^d (y_i - \alpha_i) \mathbf{b}_i$$

# PCA as Linear Autoencoder

Raw data ( $X_{n \times d}$ )  $\rightarrow$  covariance matrix ( $\Sigma_X$ )  $\rightarrow$   
eigenvalue decomposition ( $\lambda_{d \times 1}$  and  $E_{d \times d}$ )  $\rightarrow$   
principal component ( $P_{d \times m}$ )  $\rightarrow Y_{n \times m} = X_{n \times d} * P_{d \times m}$

# The two papers in 2006

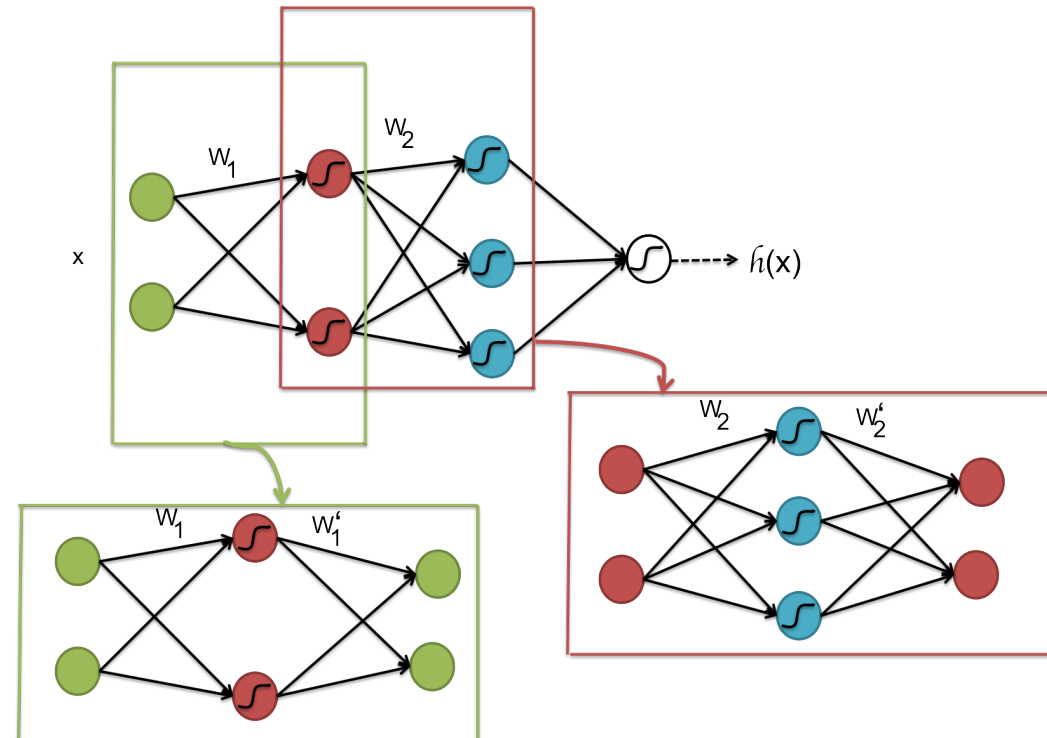
- [Hinton:2006a] G.E. Hinton, S. Osindero, Y.W. Teh, “A fast learning algorithm for deep belief nets,” *Neural Computation*, 18(7):1527-1554, 2006.
- [Hinton:2006b] G.E. Hinton, R.R. Salakhutdinov, “Reducing the dimensionality of data with neural networks,” *Science*, 313:504-507, July 2006.

# Techniques to avoid overfitting

- Regularization
  - Weight decay or L1/L2 normalization
  - Use dropout
  - Data augmentation
- Use unlabeled data to train a different network and then use the weight to initialize our network
  - Deep belief networks (based on restricted Boltzmann Machine or RBM)
  - **Deep autoencoders (based on autoencoder)**

# AE as pretraining methods

- Pretraining step
  - Train a sequence of shallow autoencoders, greedily one layer at a time, using unsupervised data
- Fine-tuning step 1
  - Train the last layer using supervised data
- Fine-tuning step 2
  - Use backpropagation to fine-tune the entire network using supervised data

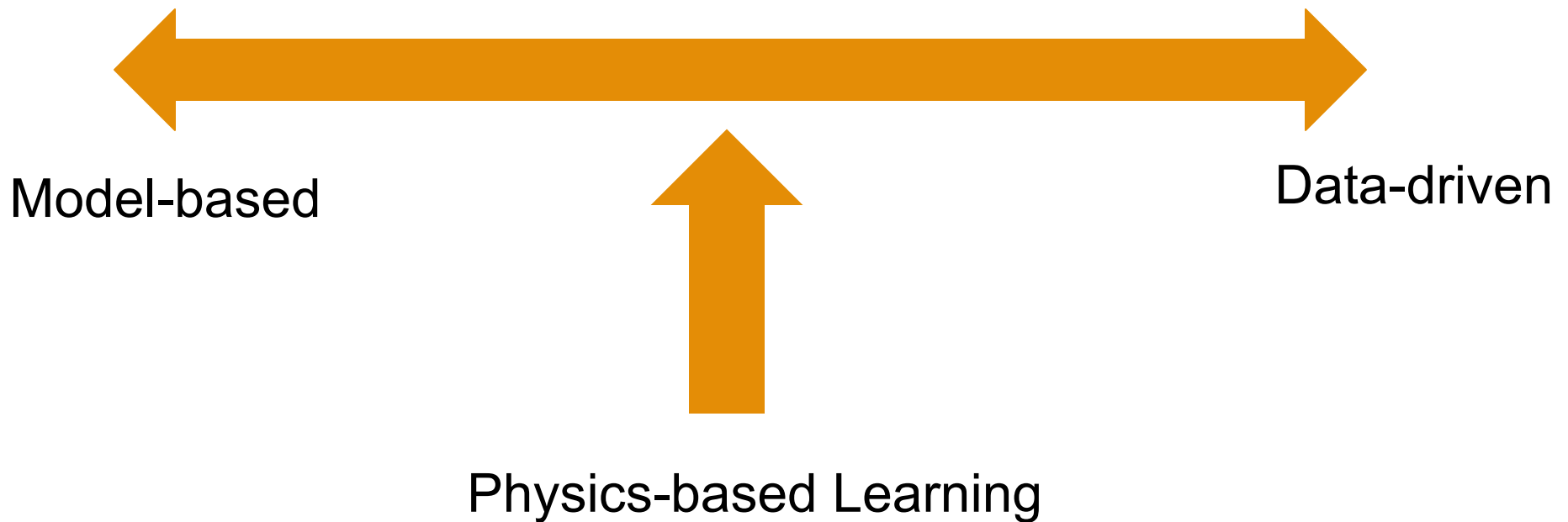




# A list of misconceptions

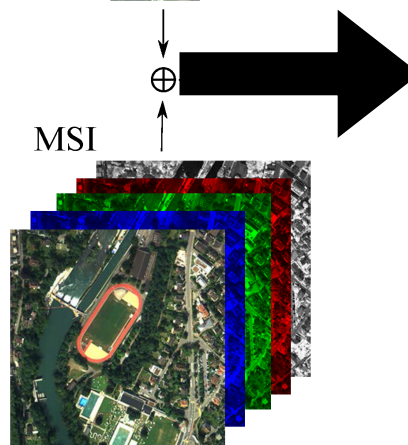
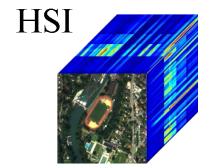
- Is deep learning merely deeper?
  - The two unique features of convolutional neural network (CNN)
- Is deep learning a classifier?
  - Engineered features vs. automatic features
- Supervised vs. Unsupervised
- **Model-based approach vs. Data-driven approach**
  - **the two extremes?**
- The world beyond CNN
  - GAN, AE, RNN, RL
- Implementation

# From model-based to data-driven



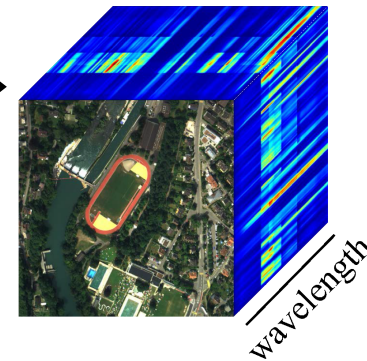
# Case study: Hyperspectral Image (HSI) Super-Resolution (SR)

Hyperspectral images (HSI):  
Low spatial but high spectral  
resolution



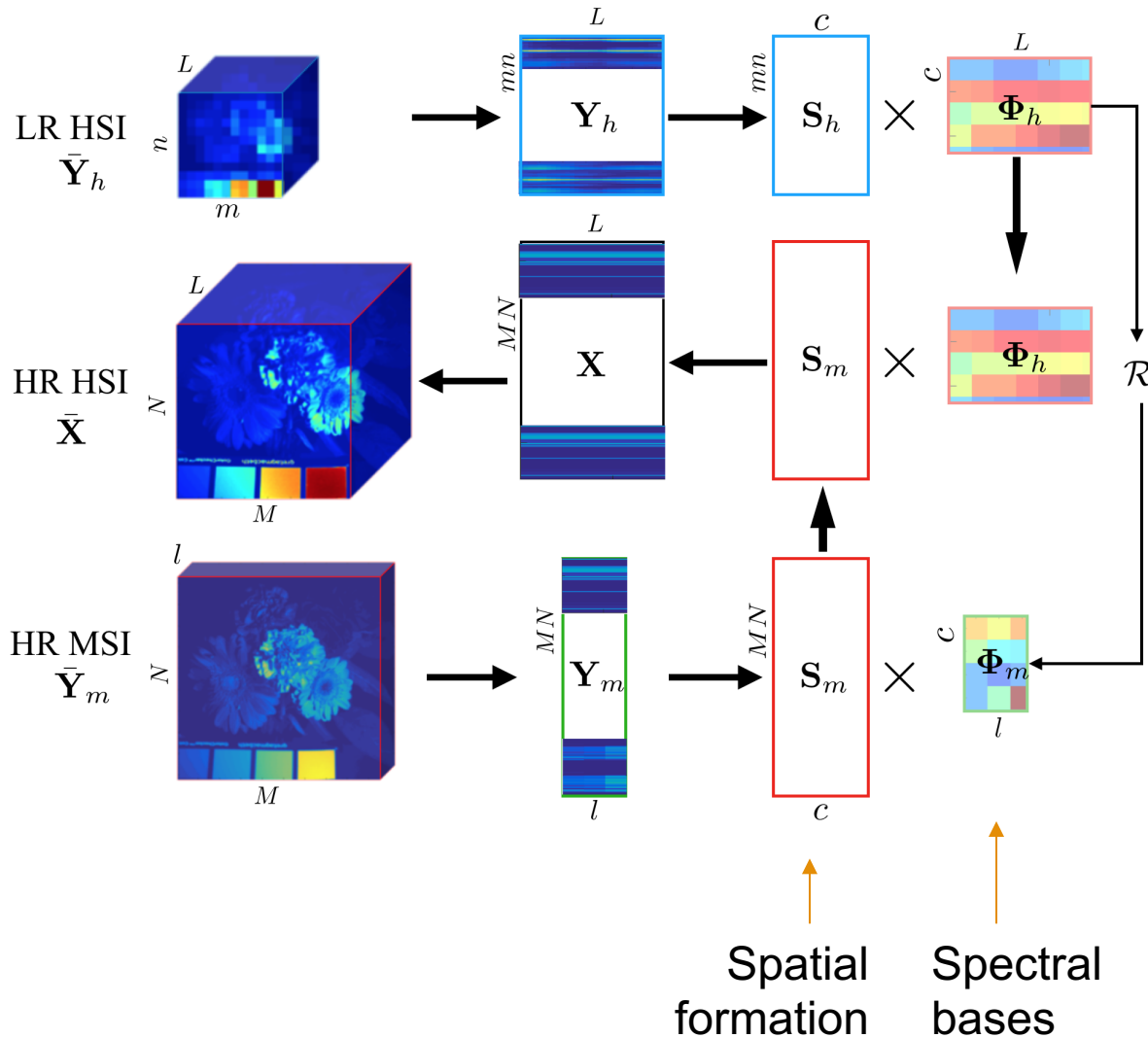
Multispectral images (MSI):  
High spatial but low spectral  
resolution

Super-Resolution



HSI-SR:  
High spatial and  
High spectral  
resolution

# The traditional formulation



$$Y_h = S_h \Phi_h,$$

$$X = S_m \Phi_h.$$

$$Y_m = S_m \Phi_m,$$

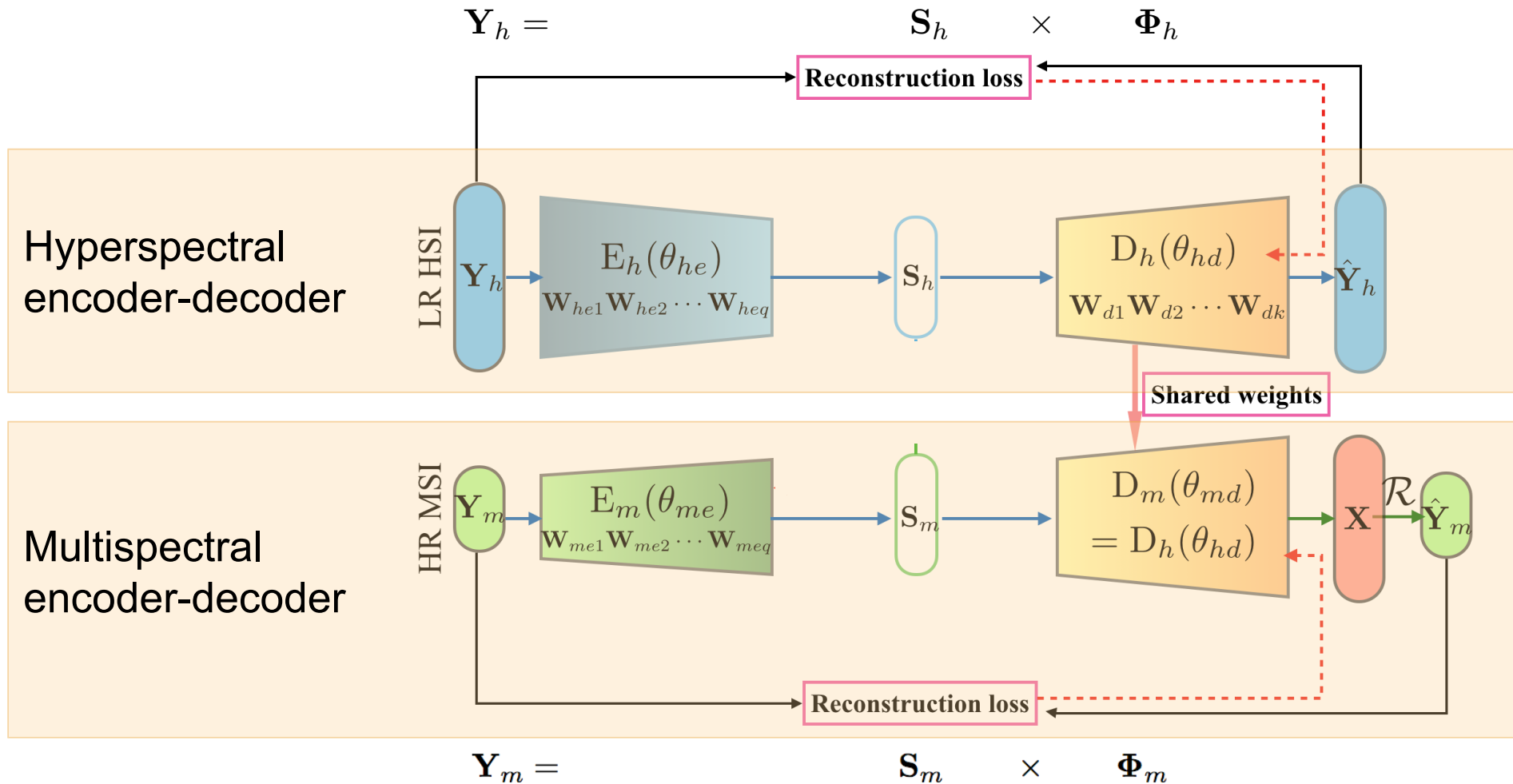
The objective function:

$$P(X|Y_h, Y_m)$$

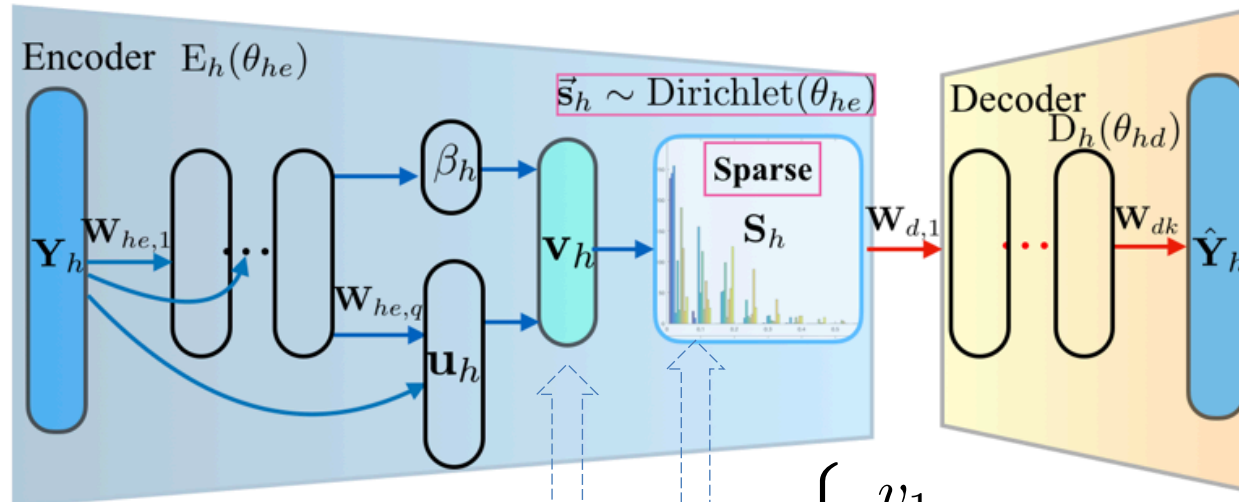
The constraints on S:

- Sum-to-one and non-negative
- Sparse

# The deep-learning approach using unsupervised learning



# The deep-learning approach with two physical constraints on S (Sum-to-one and Non-negativity)



$$s_j = \begin{cases} v_1 & \text{for } j = 1 \\ v_j \prod_{o < j} (1 - v_o) & \text{for } j > 1 \end{cases}$$

Kumaraswamy  $v_o \sim (1 - (1 - u^{\frac{1}{\beta}})^{\frac{1}{\alpha}})$ .

$$v_o \sim \text{Beta}(u, 1, \beta)$$

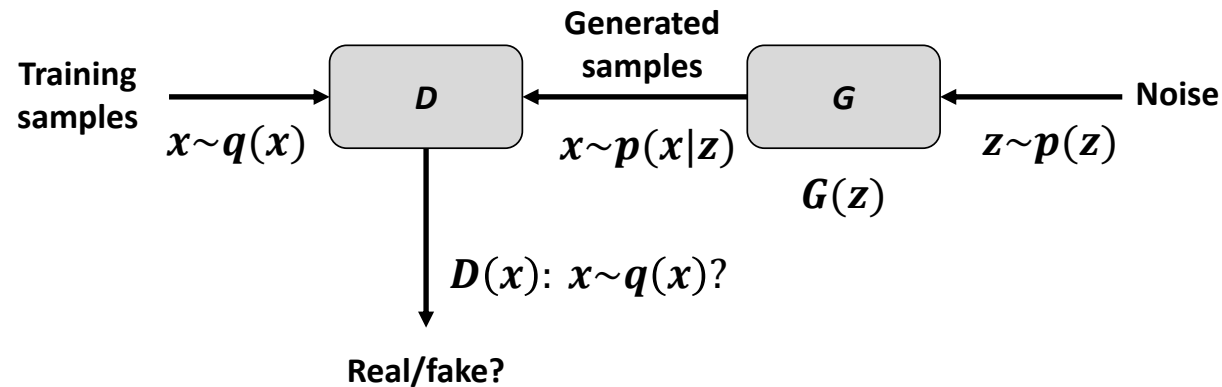
$$\text{kuma}(u, \alpha, \beta) = \alpha \beta u^{\alpha-1} (1 - u^\alpha)^{\beta-1}$$

# A list of misconceptions

- Is deep learning merely deeper?
  - The two unique features of convolutional neural network (CNN)
- Is deep learning a classifier?
  - Engineered features vs. automatic features
- Supervised vs. Unsupervised
- Model-based approach vs. Data-driven approach
  - the two extremes?
- **The world beyond CNN**
  - **GAN, AE, RNN, RL**
- Implementation

# GAN

- Two neural networks compete against each other
  - A **generator** network  $G$ : mimic training samples to fool the discriminator
  - A **discriminator** network  $D$ : discriminate training samples and generated samples



$$\text{For } D: \quad \max_D \mathbb{E}_{x \sim q(x)} [\log(D(x))] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$

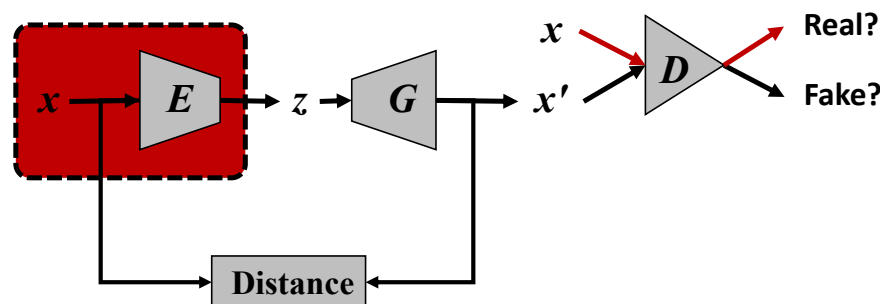
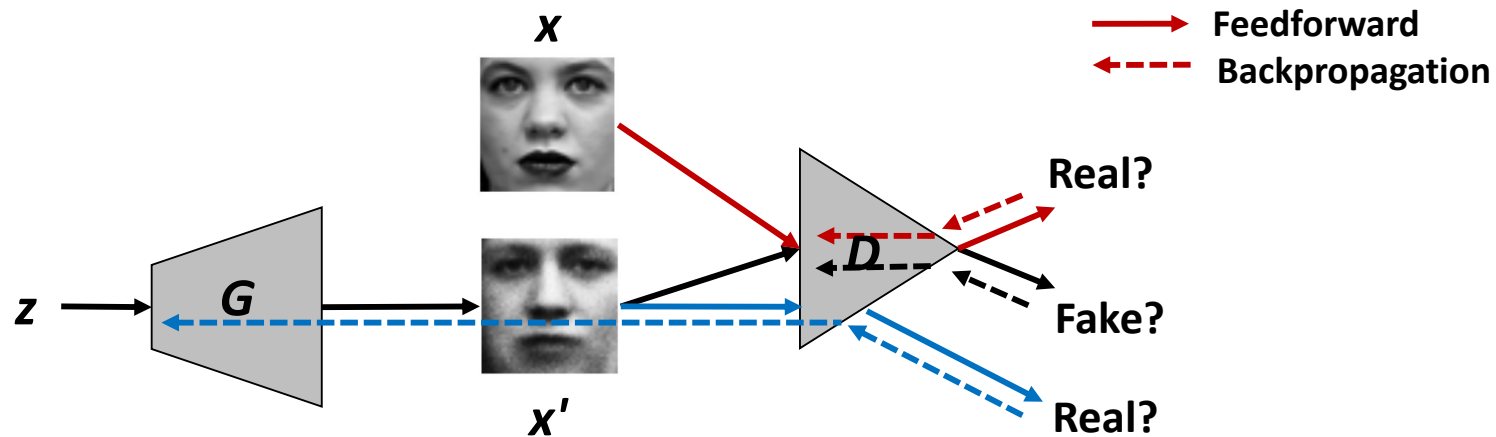
$$\text{For } G: \quad \min_G \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$



# GAN

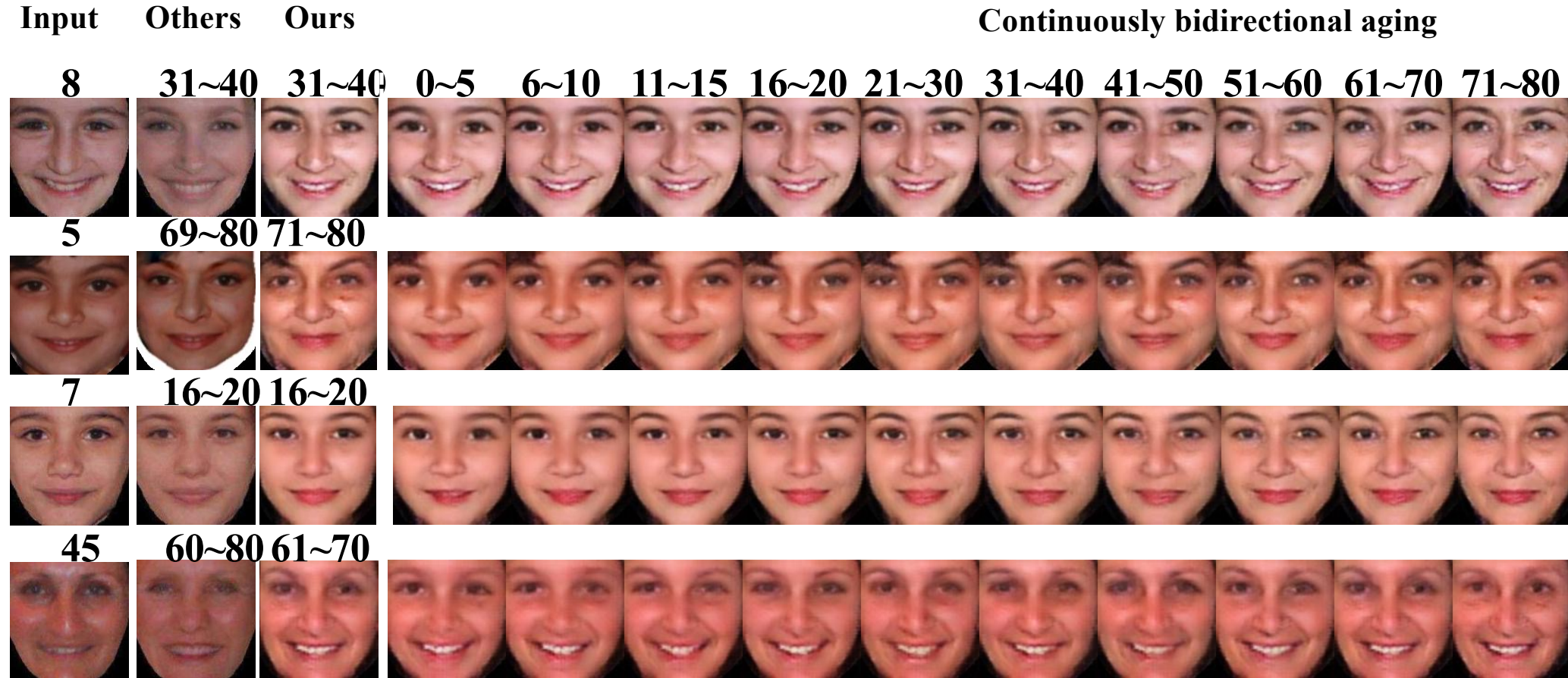
The objective function of GANs:

$$\min_G \max_D \mathbb{E}_{x \sim q(x)} [\log(D(x))] + \mathbb{E}_{z \sim p(z)} [\log(1 - D(G(z)))]$$



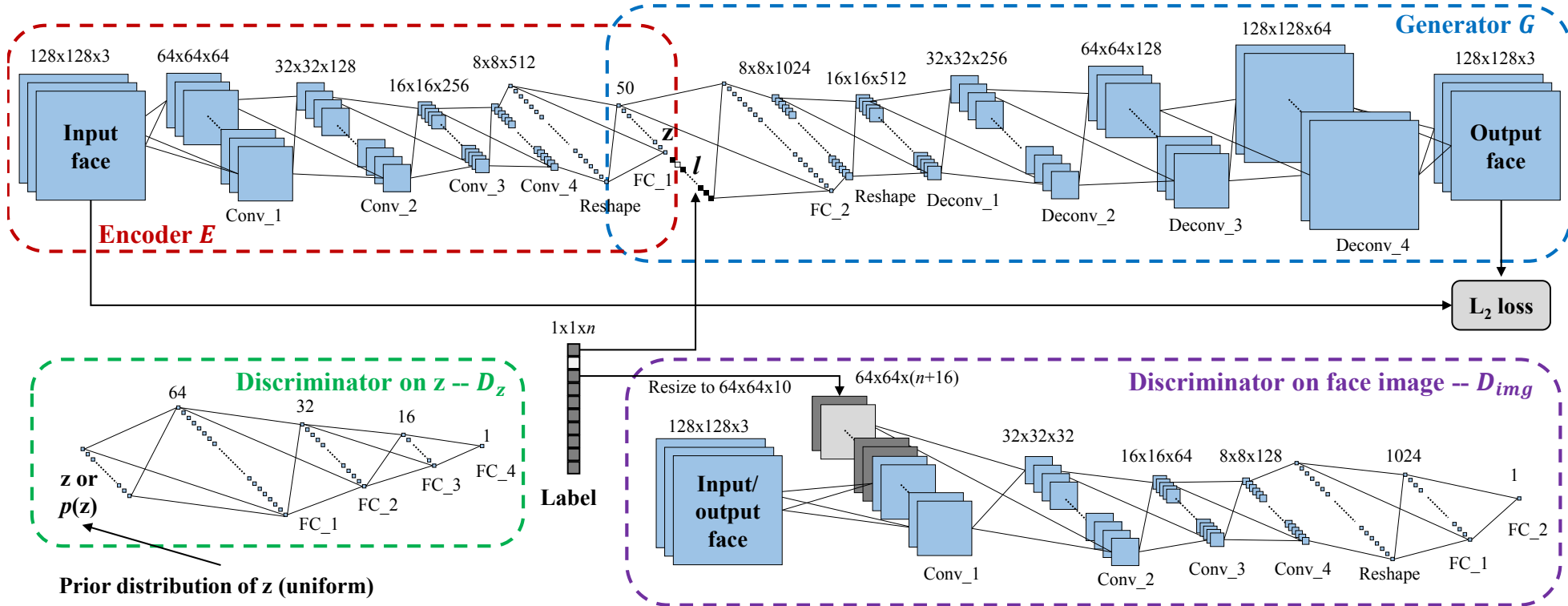
Adding an AE

# Case study: Age progression and regression

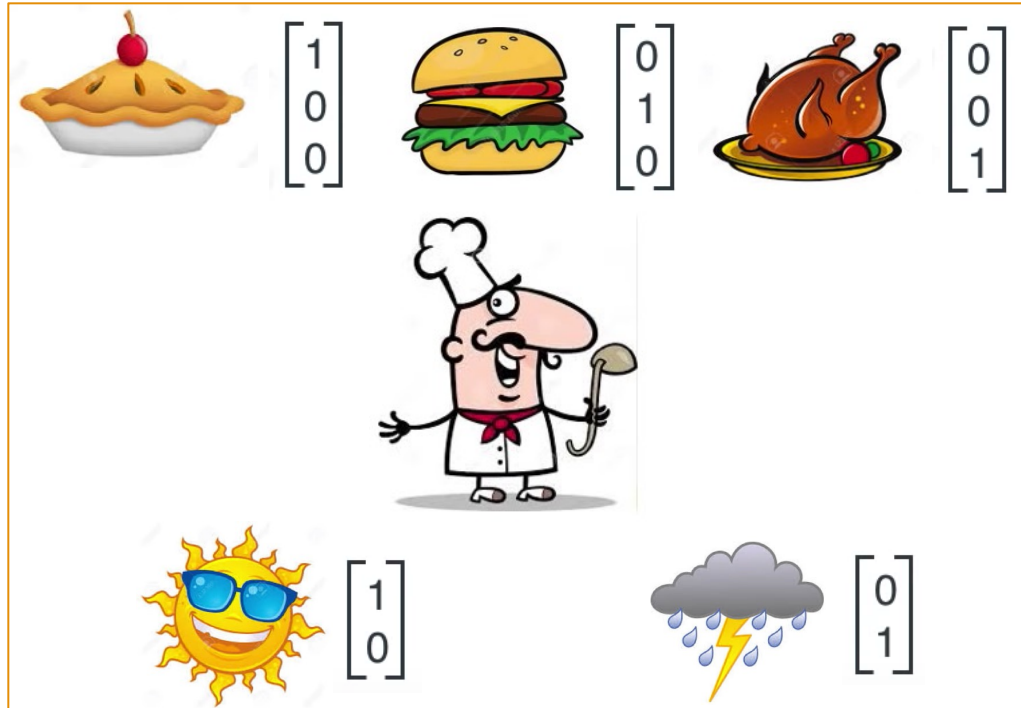


Project page: <https://zzutk.github.io/Face-Aging-CAAE>

# Case study: Conditional Adversarial Autoencoder - CAAE

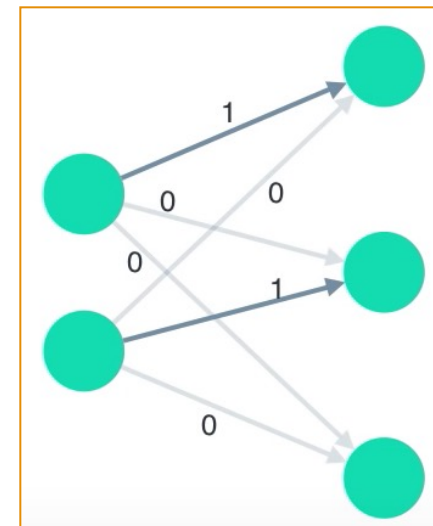
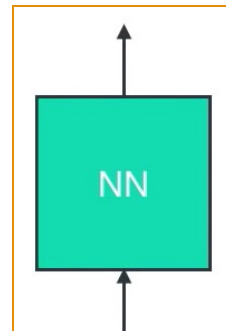


# RNN: A friendly introduction to NN

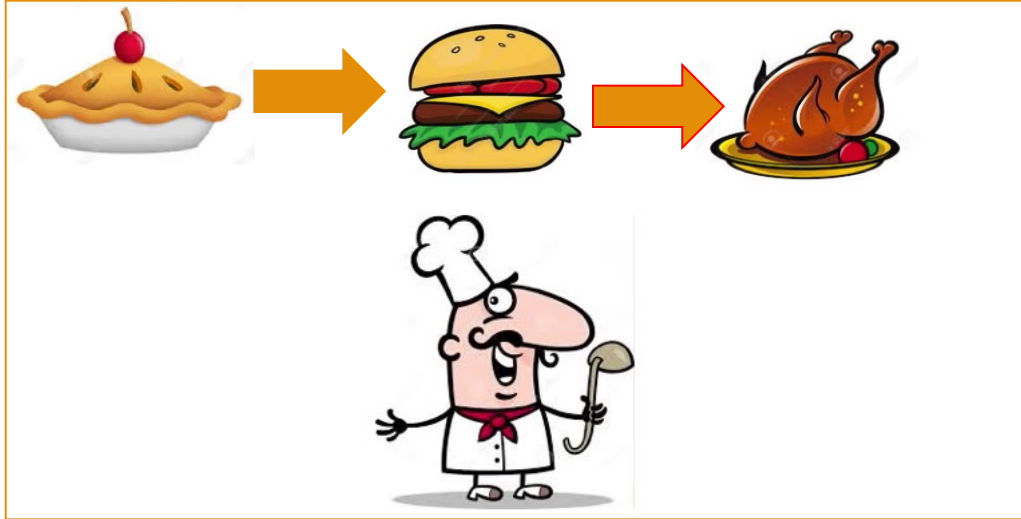


$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix} \text{ (Sun)} = \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \text{ (Pie)}$$

$$\begin{bmatrix} 1 & 0 \\ 0 & 1 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \end{bmatrix} \text{ (Storm)} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \text{ (Burger)}$$

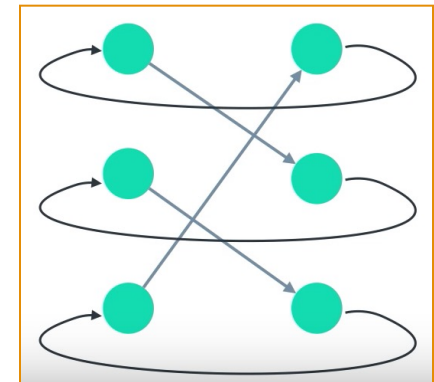
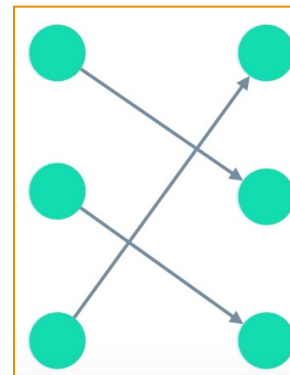


# A friendly introduction to RNN

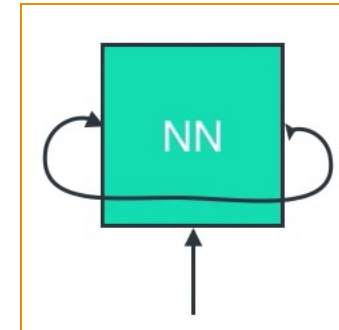
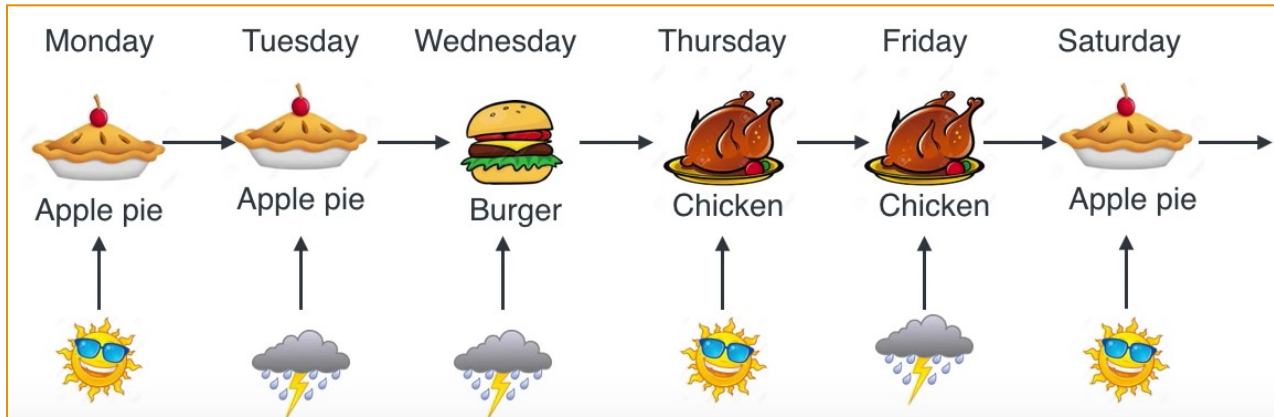


$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \text{ pie} = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \text{ burger}$$

$$\begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 0 \\ 0 & 1 & 0 \end{bmatrix} \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \text{ burger} = \begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix} \text{ turkey}$$



# A more complicated case



$$\begin{array}{c}
 \left[ \begin{array}{ccc|c}
 1 & 0 & 0 & 1 \\
 0 & 1 & 0 & 0 \\
 0 & 0 & 1 & 0 \\
 \hline
 0 & 0 & 1 & \\
 1 & 0 & 0 & \\
 0 & 1 & 0 & 
 \end{array} \right] \begin{array}{c} \\ \\ \\ \\ \\ \\ \end{array} \begin{array}{c} \text{Apple pie} \\ \\ \\ \\ \\ \\ \end{array} = \begin{array}{c} \left[ \begin{array}{c|c}
 1 & \text{Apple pie} \\
 0 & \\
 0 & \\
 0 & \\
 1 & \text{Burger} \\
 0 & 
 \end{array} \right] \begin{array}{c} \\ \\ \\ \\ \\ \\ \end{array} \begin{array}{c} \text{Same} \\ \\ \\ \\ \text{Next day} \\ \\ \end{array}
 \end{array}$$

Food

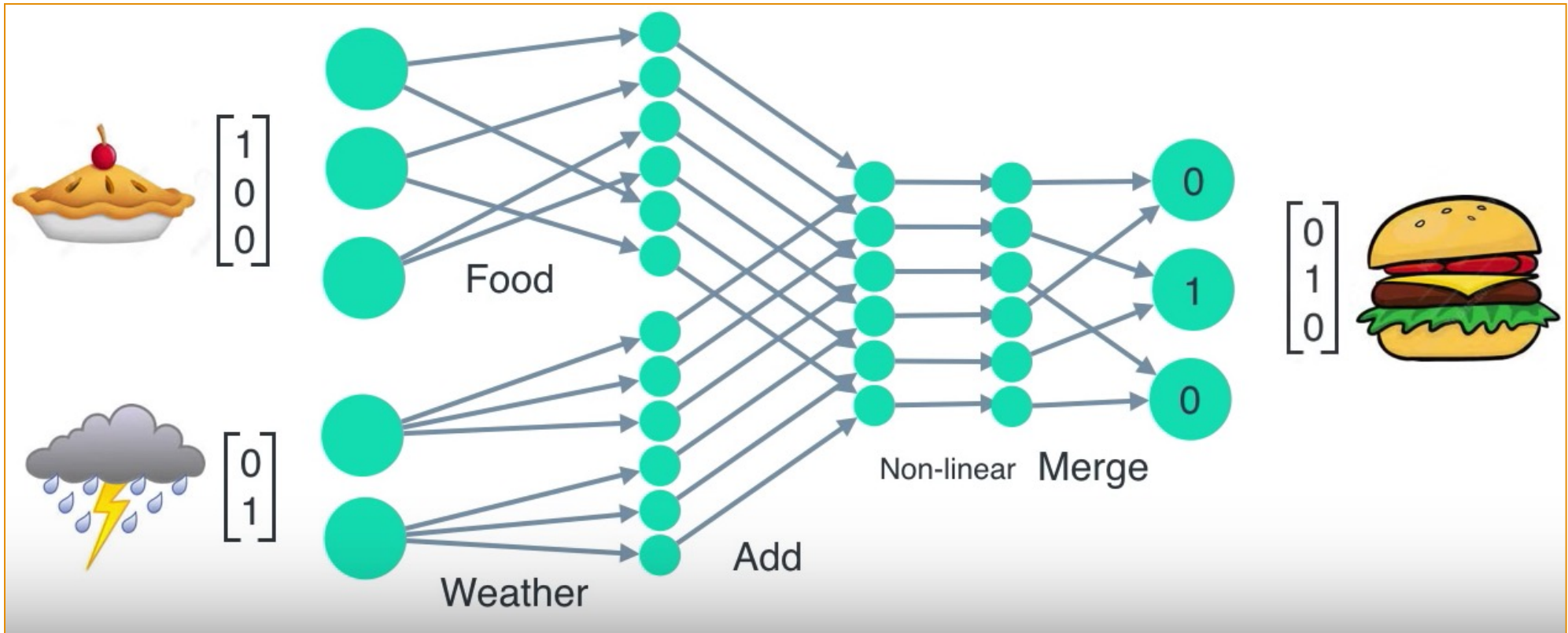
$$\begin{array}{c}
 \left[ \begin{array}{c|c}
 1 & 0 \\
 1 & 0 \\
 1 & 0 \\
 \hline
 0 & 1 \\
 0 & 1 \\
 0 & 1 \\
 \end{array} \right] \begin{array}{c} \\ \\ \\ \\ \\ \\ \end{array} \begin{array}{c} \text{Rainy} \\ \\ \\ \\ \\ \\ \end{array} = \begin{array}{c} \left[ \begin{array}{c|c}
 0 & \text{Sunny} \\
 0 & \\
 0 & \\
 1 & \\
 1 & \\
 1 & 
 \end{array} \right] \begin{array}{c} \\ \\ \\ \\ \\ \\ \end{array} \begin{array}{c} \text{Same} \\ \\ \\ \\ \text{Next day} \\ \\ \end{array}
 \end{array}$$

Weather

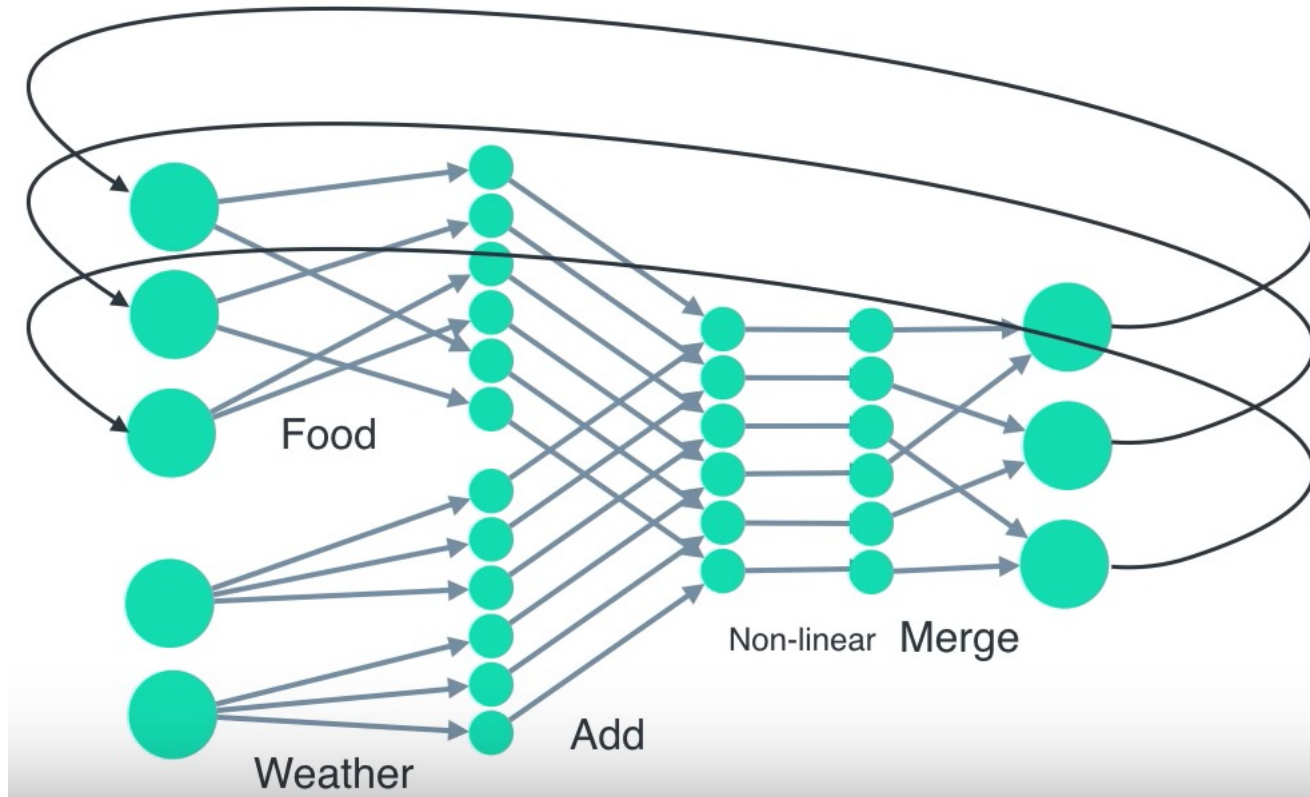


$$\begin{array}{c}
 \left[ \begin{array}{c|c}
 1 & \text{Apple pie} \\
 0 & \\
 0 & \\
 0 & \\
 1 & \text{Burger} \\
 0 & 
 \end{array} \right] \begin{array}{c} \\ \\ \\ \\ \\ \\ \end{array} \begin{array}{c} \text{Same} \\ \\ \\ \\ \text{Next day} \\ \\ \end{array} + \begin{array}{c} \left[ \begin{array}{c|c}
 0 & \text{Sunny} \\
 0 & \\
 0 & \\
 1 & \\
 1 & \\
 1 & 
 \end{array} \right] \begin{array}{c} \\ \\ \\ \\ \\ \\ \end{array} \begin{array}{c} \text{Same} \\ \\ \\ \\ \text{Next day} \\ \\ \end{array} = \begin{array}{c} \left[ \begin{array}{c|c}
 1 & \\
 0 & \\
 0 & \\
 1 & \\
 2 & \\
 1 & 
 \end{array} \right] \begin{array}{c} \\ \\ \\ \\ \\ \\ \end{array}
 \end{array}$$

# A more complicated case (cont'd)

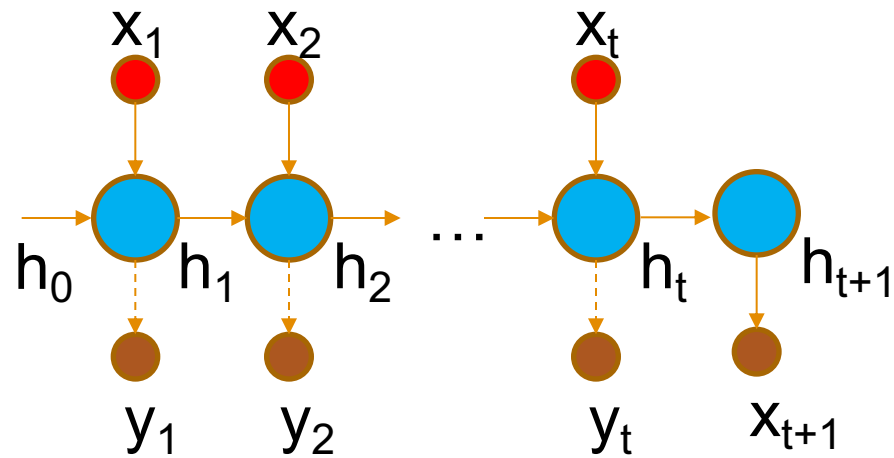


# A more complicated case (cont'd)



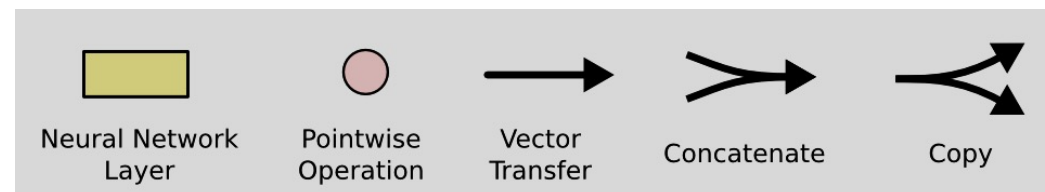
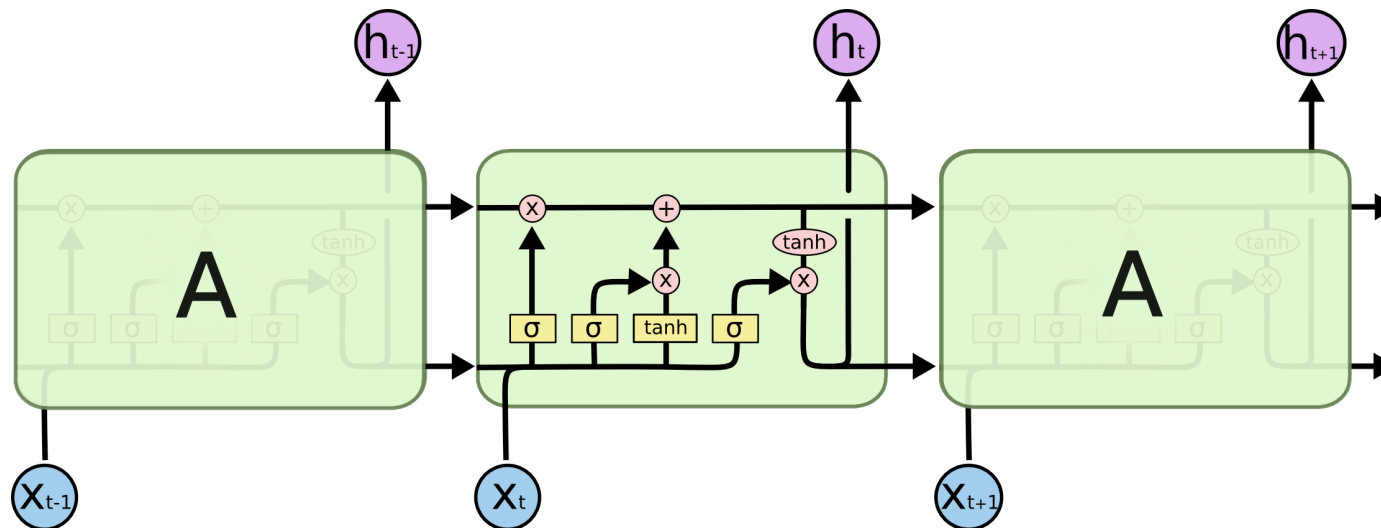


# Recurrent neural network (RNN)



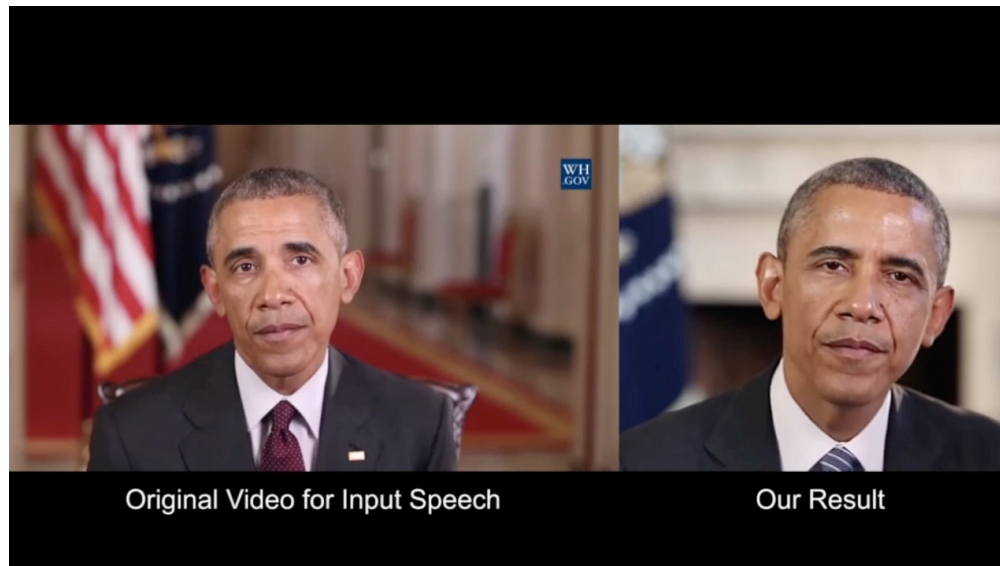
# The long-short term memory (LSTM) module

LSTMs are explicitly designed to avoid the long-term dependency problem.



# Case study: The talking face

Goal: Given an arbitrary audio clip and a face image, automatically generate realistic and smooth face video with accurate lip sync.

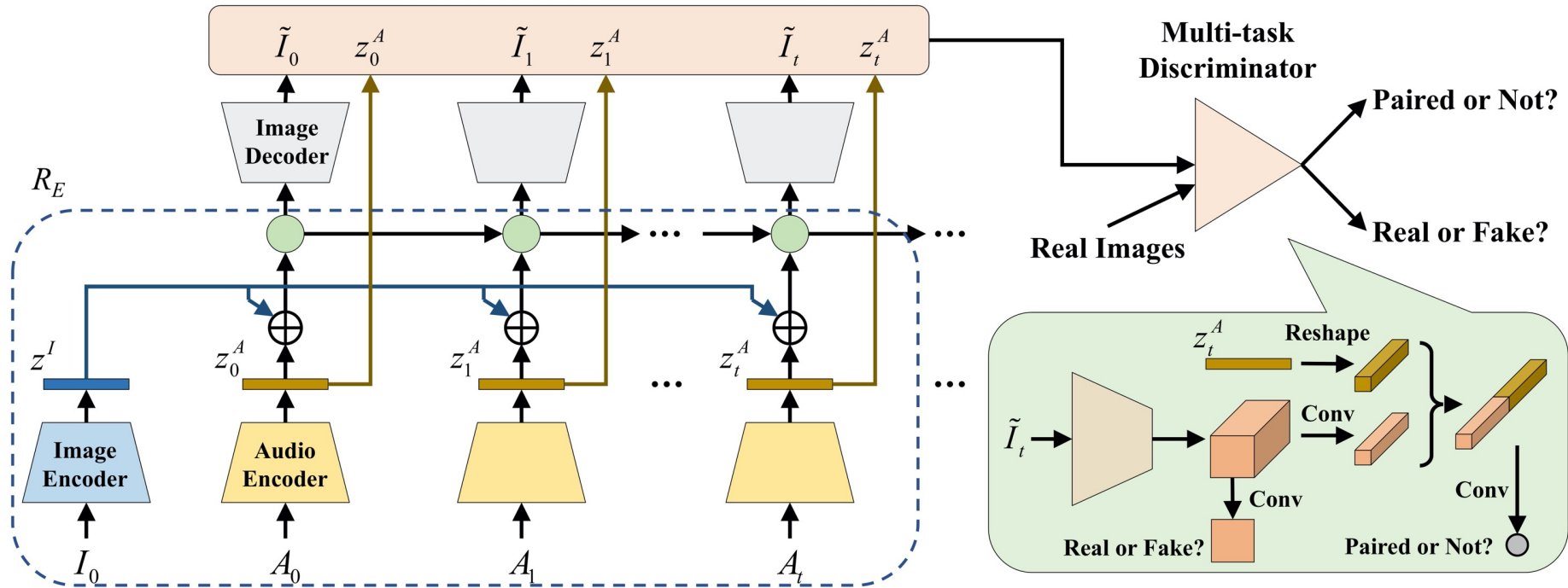


[Suwajanakorn et al., 2017]

Application: Face animation, entertainment, video bandwidth reduction, etc.

# The talking face

The proposed method: conditional video generation



<http://web.eecs.utk.edu/~ysong18/projects/talkingface/talkingface.html>

# A list of misconceptions

- Is deep learning merely deeper?
  - The two unique features of convolutional neural network (CNN)
- Is deep learning a classifier?
  - Engineered features vs. automatic features
- Supervised vs. Unsupervised
- Model-based approach vs. Data-driven approach – the two extremes?
- The world beyond CNN
  - GAN, AE, RNN, RL
- Implementation
  - Matlab
  - TensorFlow
  - PyTorch
  - Keras

Education is what remains  
after one has forgotten  
everything one learned in  
school. -- Albert Einstein