# E9: 309 Advanced Deep Learning
## 7-10-2020

**Instructor**: Sriram Ganapathy
sriramg@iisc.ac.in

**Teaching Assistant** : Jaswanth Reddy
jaswanthk@iisc.ac.in

**Schedule** - MW - 330-5pm (Microsoft Teams)
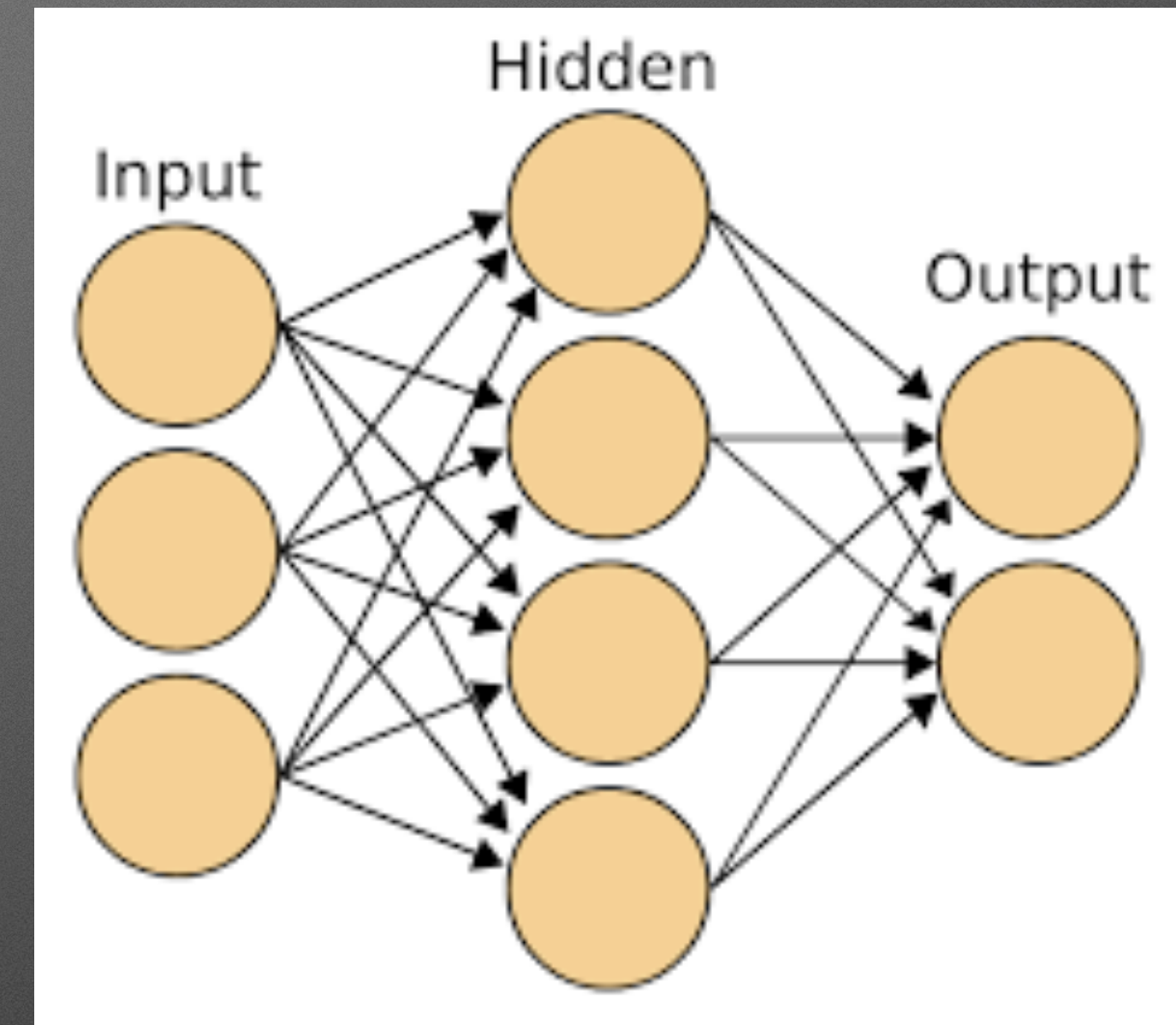http://leap.ee.iisc.ac.in/sriram/teaching/ADL2020/

# Recap of deep learning

# Some notations

* $x \in \mathcal{R}^D$ - input data.

* $y \in \mathcal{R}^C$ - neural network targets.

* $\hat{y} \in \mathcal{B}^C$ - model outputs.



* $e, h \in \mathcal{R}^d$ - hidden model representations or embeddings.

* $\Theta$ - collection of learnable parameters in the model.

* $E(y, \hat{y})$ - error function used in the model training.

# Some notations

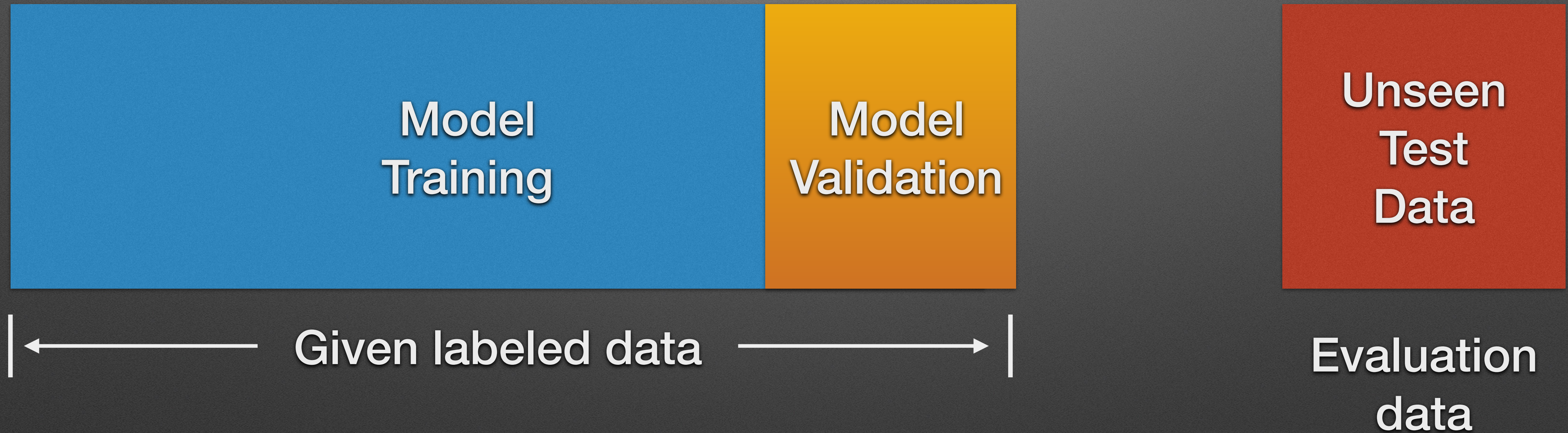* $\{\boldsymbol{x}_1, ..., \boldsymbol{x}_N, \boldsymbol{y}_1, ..., \boldsymbol{y}_N\}$ - labeled training data

* $q = \{1...Q\}$ - iteration index.

* $t = \{1...T\}$ - discrete time index.

* $l = \{1...L\}$ - layer index

* $\eta$ - learning rate (hyper-parameter)

* $N_b$ - mini-batch size and $B$ is the number of mini-batches.

# Premise

* Training data, validation data, test data.

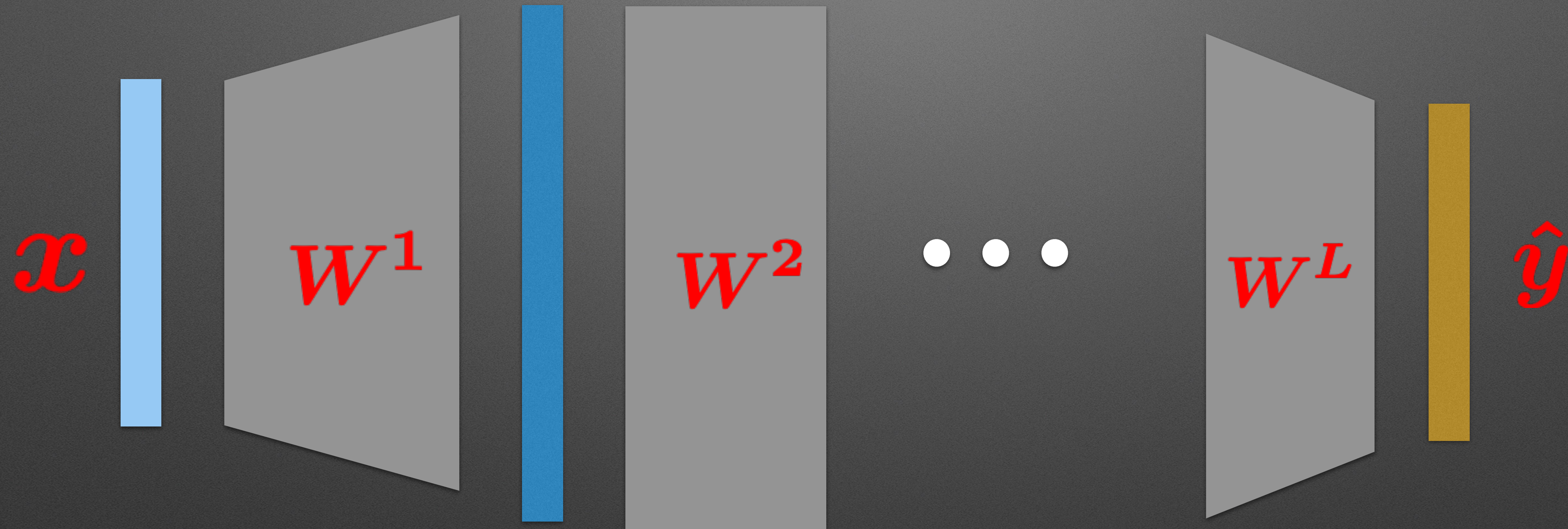| Model Training | Model Validation | | Unseen Test Data |
|:---:|:---:|---|:---:|

|← Given labeled data →| | Evaluation data |

* Model training data - used for parameter learning.

* Validation data - used for hyper-parameter tuning (cross validation CV).

# Feedforward networks



$$x \quad W^1 \quad W^2 \quad \cdots \quad W^L \quad \hat{y}$$

* Dense connections between the input and output - also called fully connected network.

# Learning in feedforward networks

∗ Stochastic gradient descent (SGD) - Initialize the model parameters $\Theta^{0,0}$ (randomly)

for $q = \{1...Q\}$:

    for $b = \{1...B\}$:

$$\Theta^{q,b} = \Theta^{q,b-1} - \eta \frac{\partial E(\{\boldsymbol{y}_{b,1},...\boldsymbol{y}_{b,N_b}\}, \{\hat{\boldsymbol{y}}_{b,1},...\hat{\boldsymbol{y}}_{b,N_b}\})}{\partial \Theta}\bigg|_{\Theta = \Theta^{q,b}}$$

$$\Theta^{q+1,0} = \Theta^{q,B}$$

return $\Theta^* = \Theta^{Q,B}$

# Learning in feedforward networks

* Learning with momentum - accelerate the learning by adding a component of the previous gradient computation.

* RMSprop $\eta' = \dfrac{\eta}{RMS(g)}$

* Adam - adaptive moment estimation

  ➡ combines momentum and RMSprop.

  ➡ empirically shown to be effective in many applications.

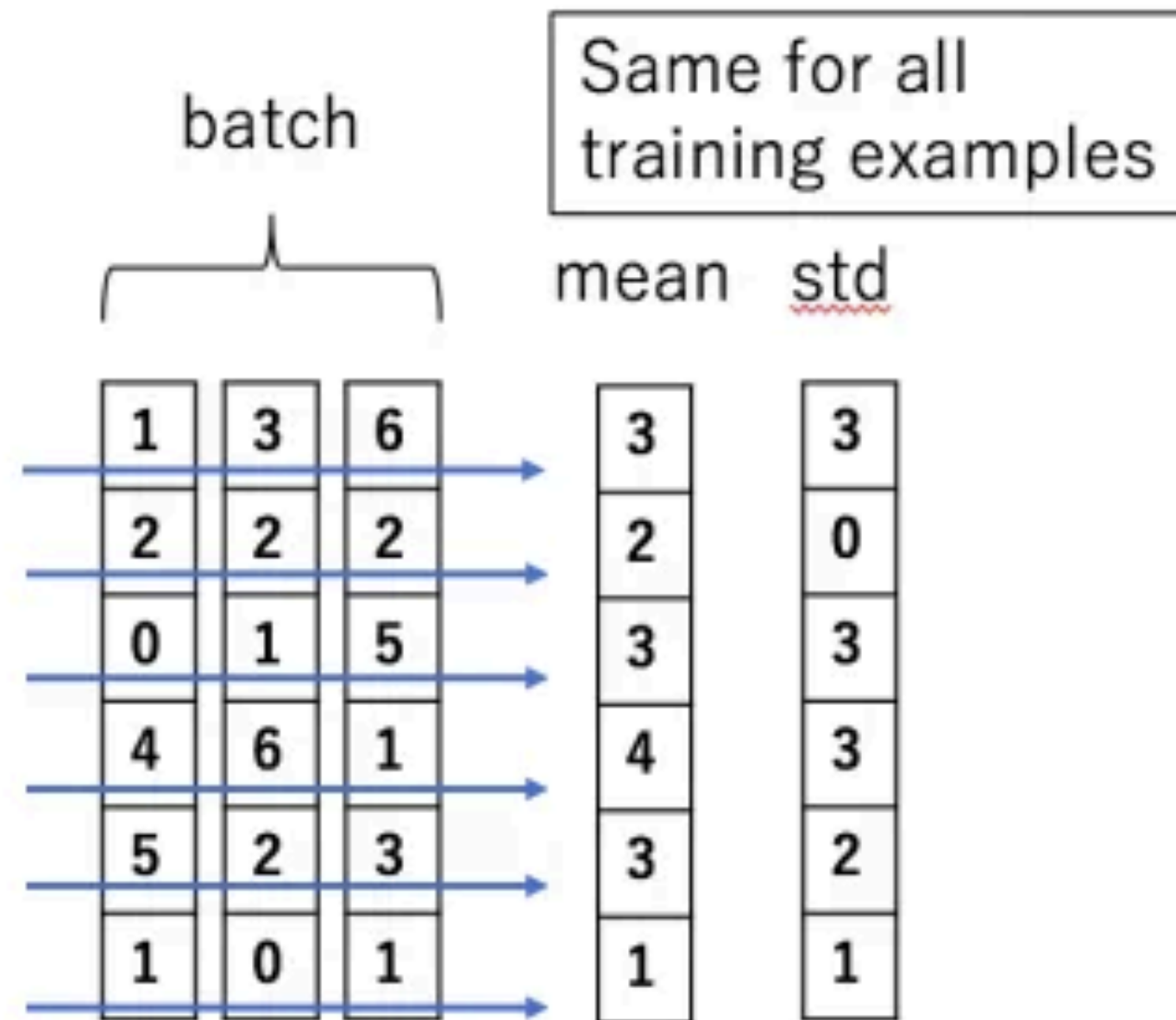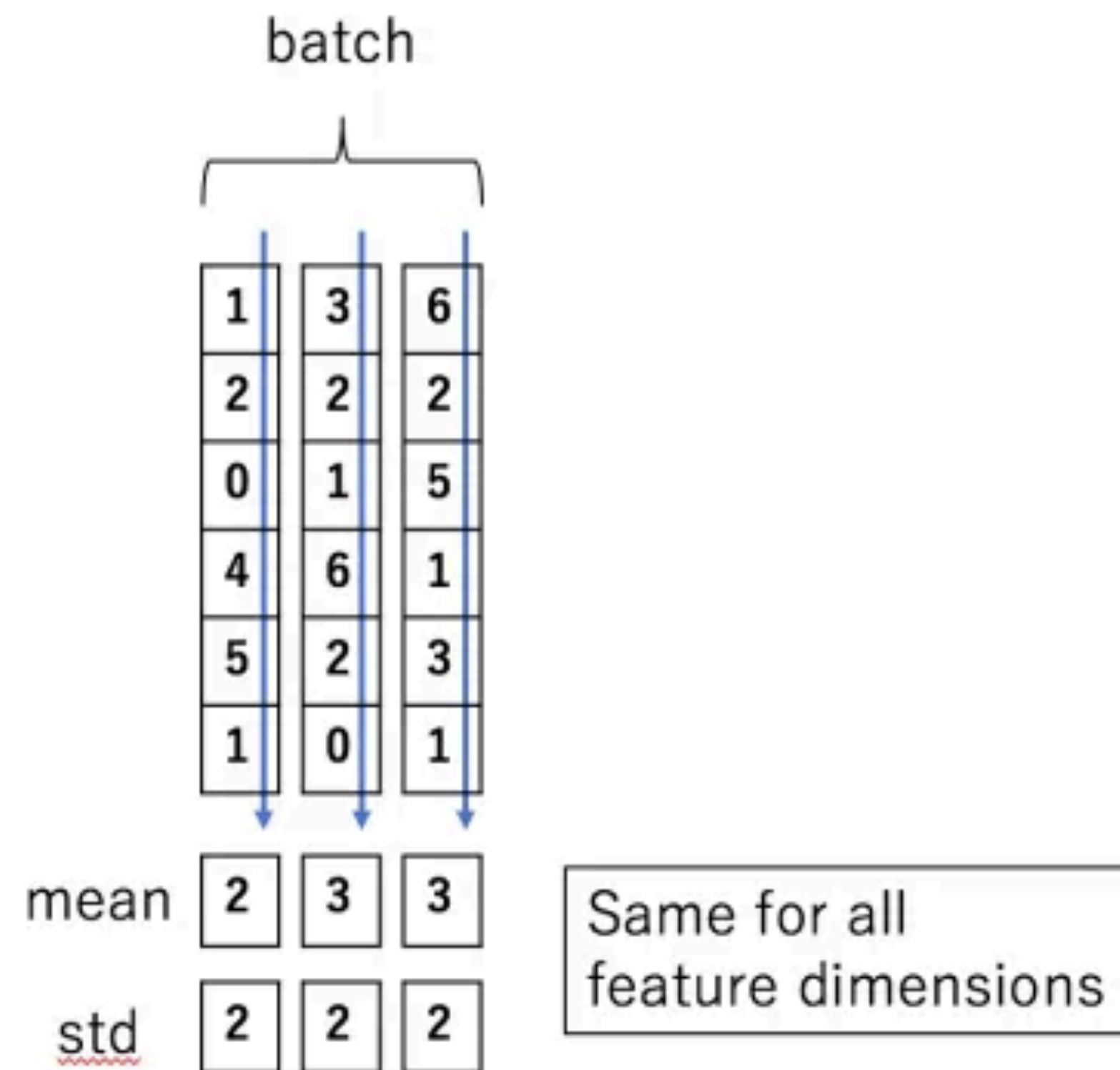Reading assignment - Overview of gradient descent algorithms
https://arxiv.org/pdf/1609.04747.pdf

# Normalization

Reading assignment - How does batchnorm help optimization
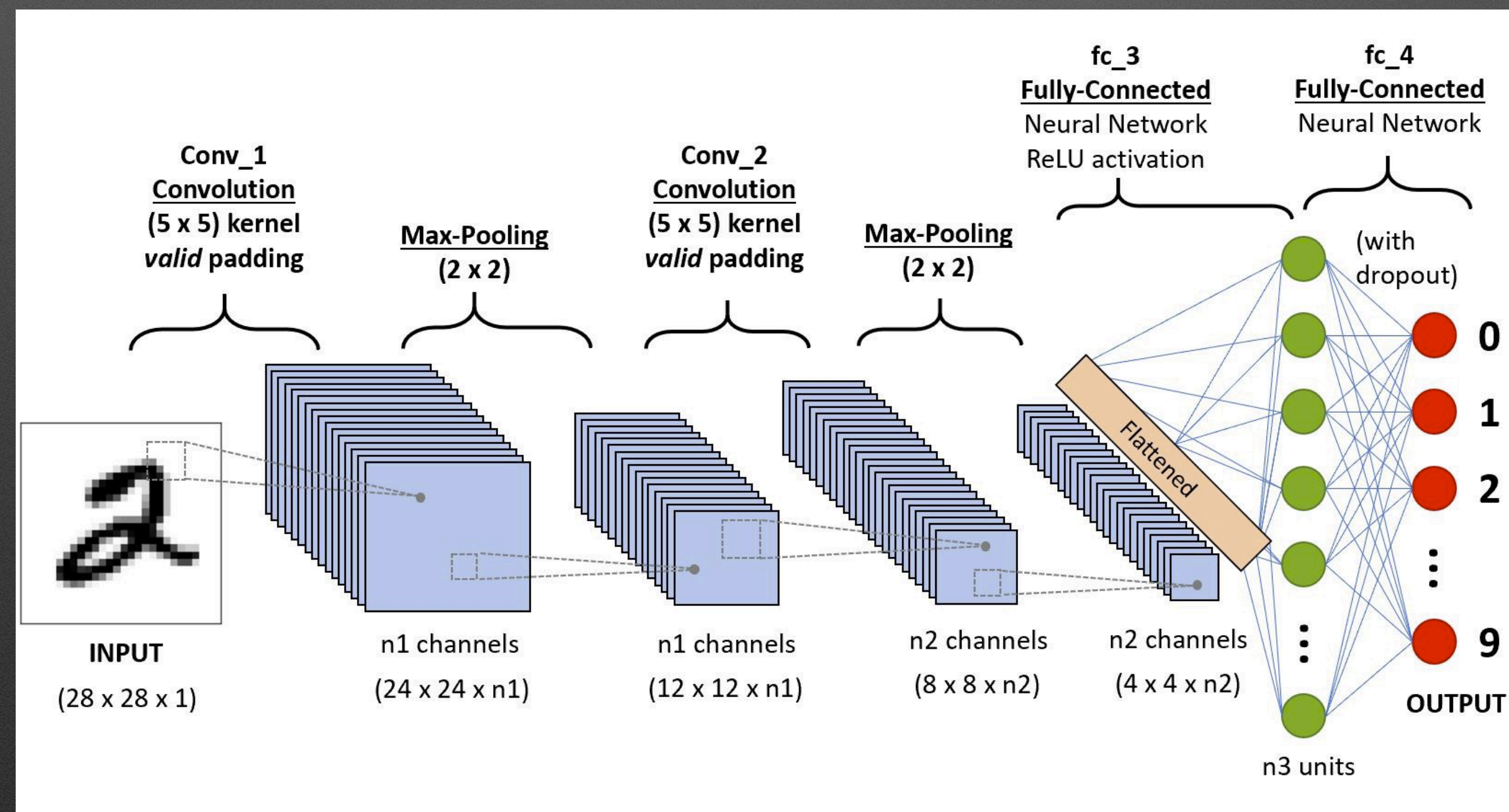https://arxiv.org/pdf/1805.11604.pdf

# Convolutional neural networks

✳ Replacing affine transformations with convolutional operations

$$H(m,n) = X * W(m,n) = \sum_{i,j} X(m+i, n+j)W(i,j)$$

✳ Usually used with max-pooling based sub-sampling

# Module - I Visual and Time Series Modeling

# Why do need recurrent models

* Learn from ordered pairs of $\boldsymbol{x}, \boldsymbol{y}$

  ✓ All the data samples are treated independently.

    ⋆ Data are shuffled before mini-batch formation

* If the input data and output labels are time-series data $\boldsymbol{x}(t), \boldsymbol{y}(t)$

    ⊙ DNNs/CNNs may fail to model the correlation of the data across the time

    ⊙ Question - how can we build models that capture the time evolution of the data and the labels.

# Why do need recurrent models

* An interesting subset of this problem is where the input alone is a time series $x(t), y$ or have different indices $x(t), y(u)$

* Examples

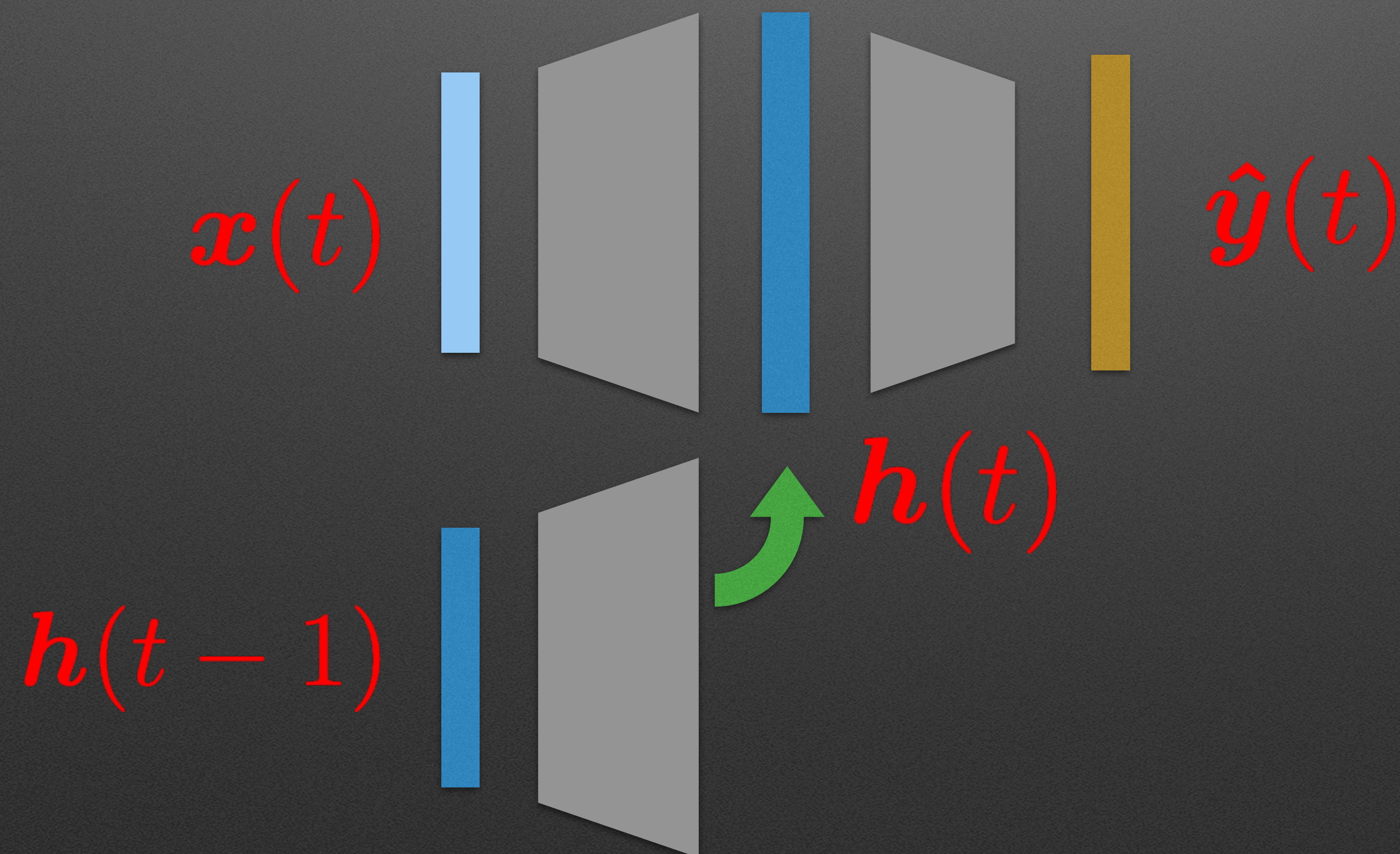  ✓ Text sequences

  ✓ Speech and audio

  ✓ Video sequences

# First order recurrence - hidden layer

∗ Making the hidden layer a function of the previous outputs from the hidden layer along with the input

$$\boldsymbol{h}(t) = f\big(\boldsymbol{h}(t-1), \boldsymbol{x}(t)\big)$$



$\boldsymbol{x}(t)$

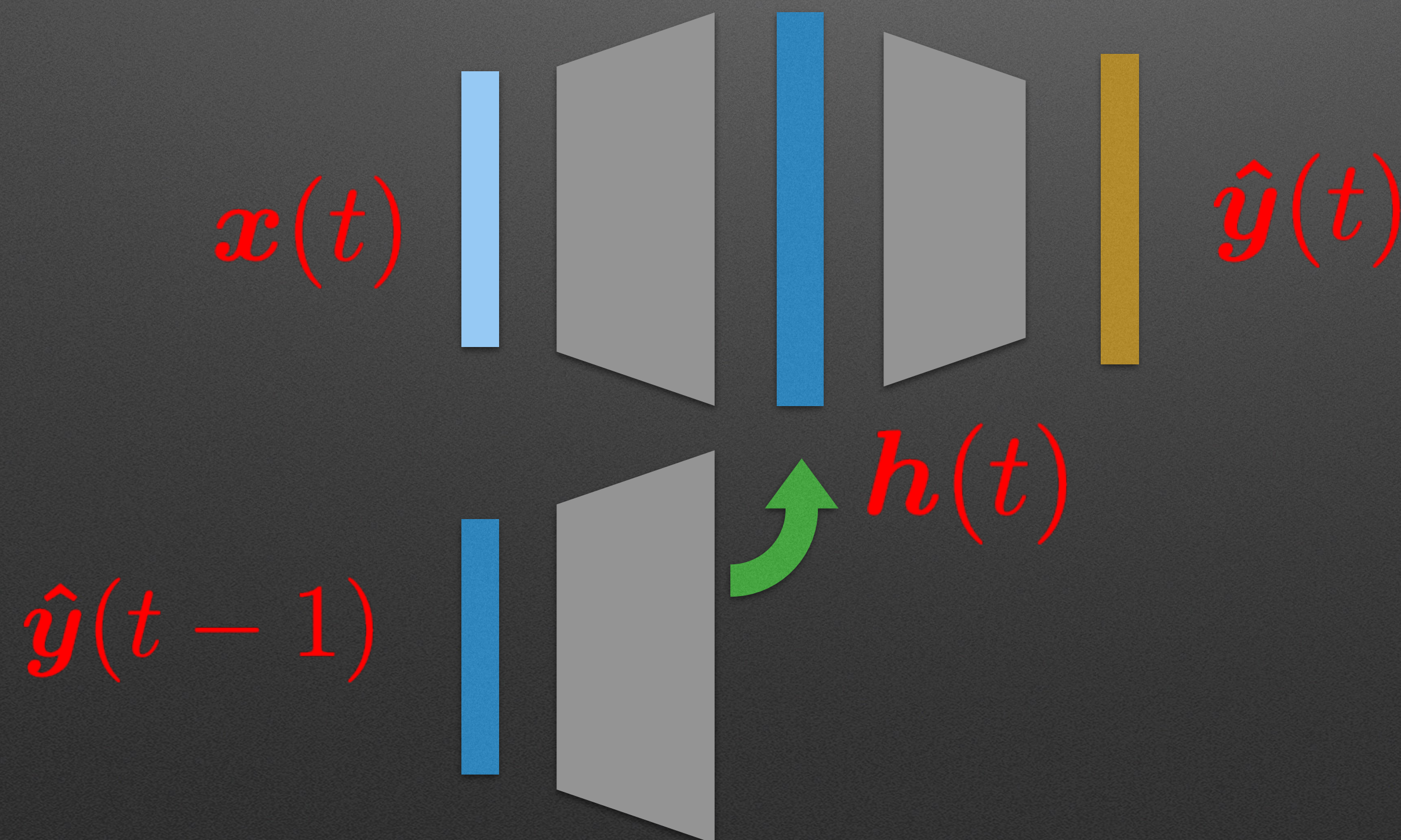$\hat{\boldsymbol{y}}(t)$

$\boldsymbol{h}(t)$

$\boldsymbol{h}(t-1)$

# First order recurrence - output layer

* Making the hidden layer a function of the previous outputs from the hidden layer along with the input

$$\boldsymbol{h}(t) = f\left(\hat{\boldsymbol{y}}(t-1), \boldsymbol{x}(t)\right)$$



$\boldsymbol{x}(t)$

$\hat{\boldsymbol{y}}(t)$

$\boldsymbol{h}(t)$

$\hat{\boldsymbol{y}}(t-1)$

# Looking forward (next lecture)

⁂ Learning in recurrence networks: Back-propagation in time.

⁂ Unsupervised Learning: Issues with forgetting and long-short-term memory networks

# The bell has rung



http://phdcomics.com/