

Fundamentals / Foundation

Data Science 9CFU Computer Science 6CFU

of Data Science

Indro Spinelli

Sapienza University of Rome

Composing two networks

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

Network 1:

$$\begin{aligned} h_2 &= a[\theta_{20} + \theta_{21}x] & y &= \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3 \\ h_3 &= a[\theta_{30} + \theta_{31}x] \end{aligned}$$

$$h'_1 = a[\theta'_{10} + \theta'_{11}y]$$

Network 2:

$$\begin{aligned} h'_2 &= a[\theta'_{20} + \theta'_{21}y] & y' &= \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3 \\ h'_3 &= a[\theta'_{30} + \theta'_{31}y] \end{aligned}$$

Composing two networks

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

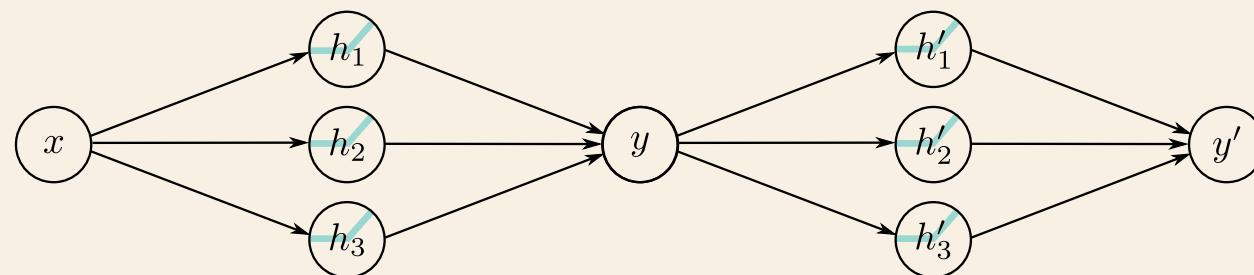
Network 1:

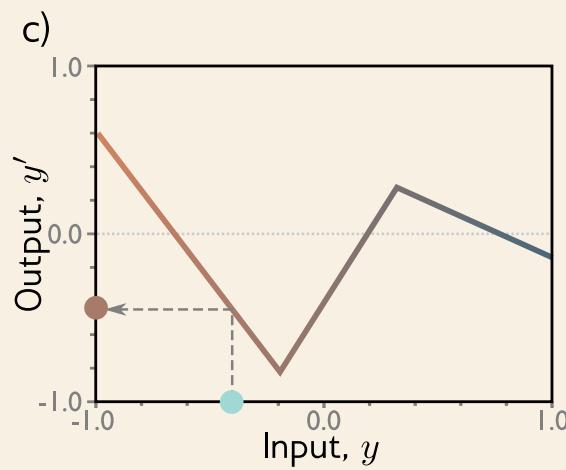
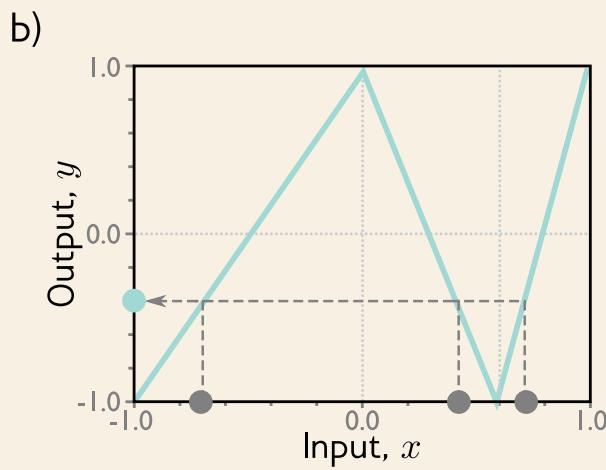
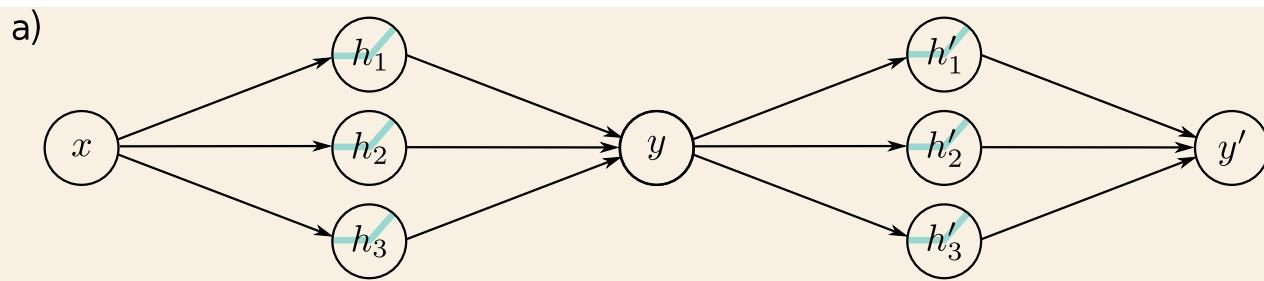
$$\begin{aligned} h_2 &= a[\theta_{20} + \theta_{21}x] & y &= \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3 \\ h_3 &= a[\theta_{30} + \theta_{31}x] \end{aligned}$$

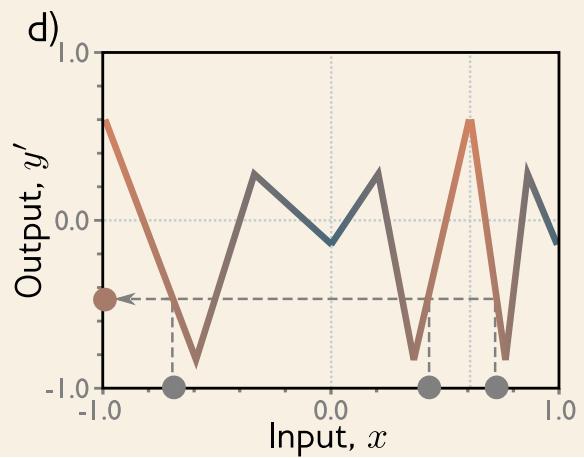
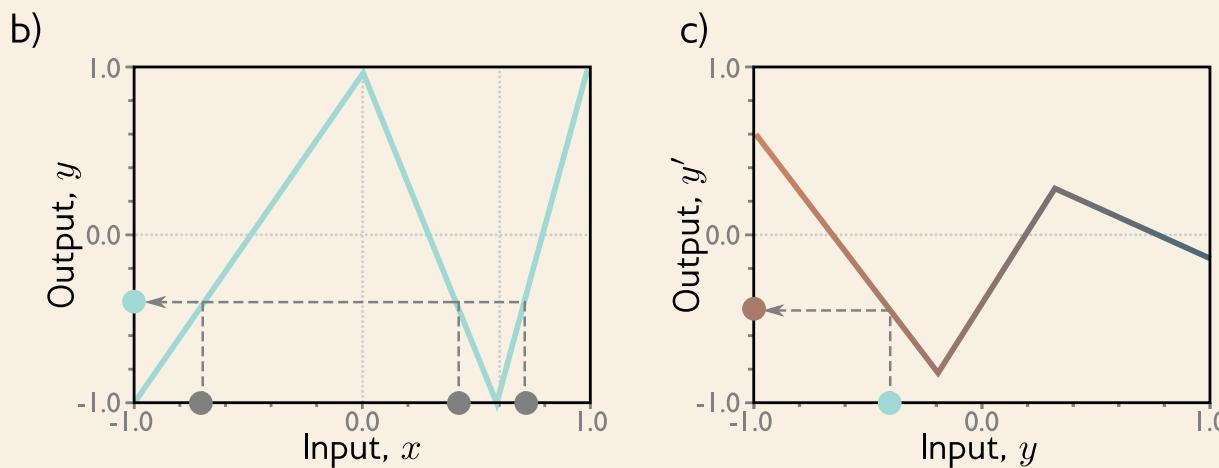
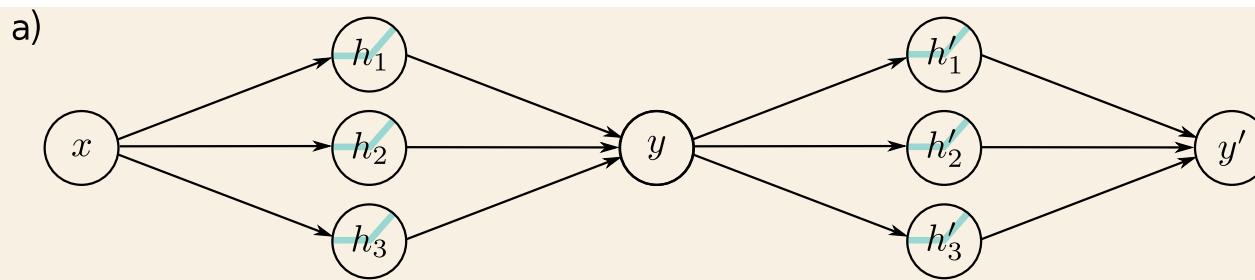
$$h'_1 = a[\theta'_{10} + \theta'_{11}y]$$

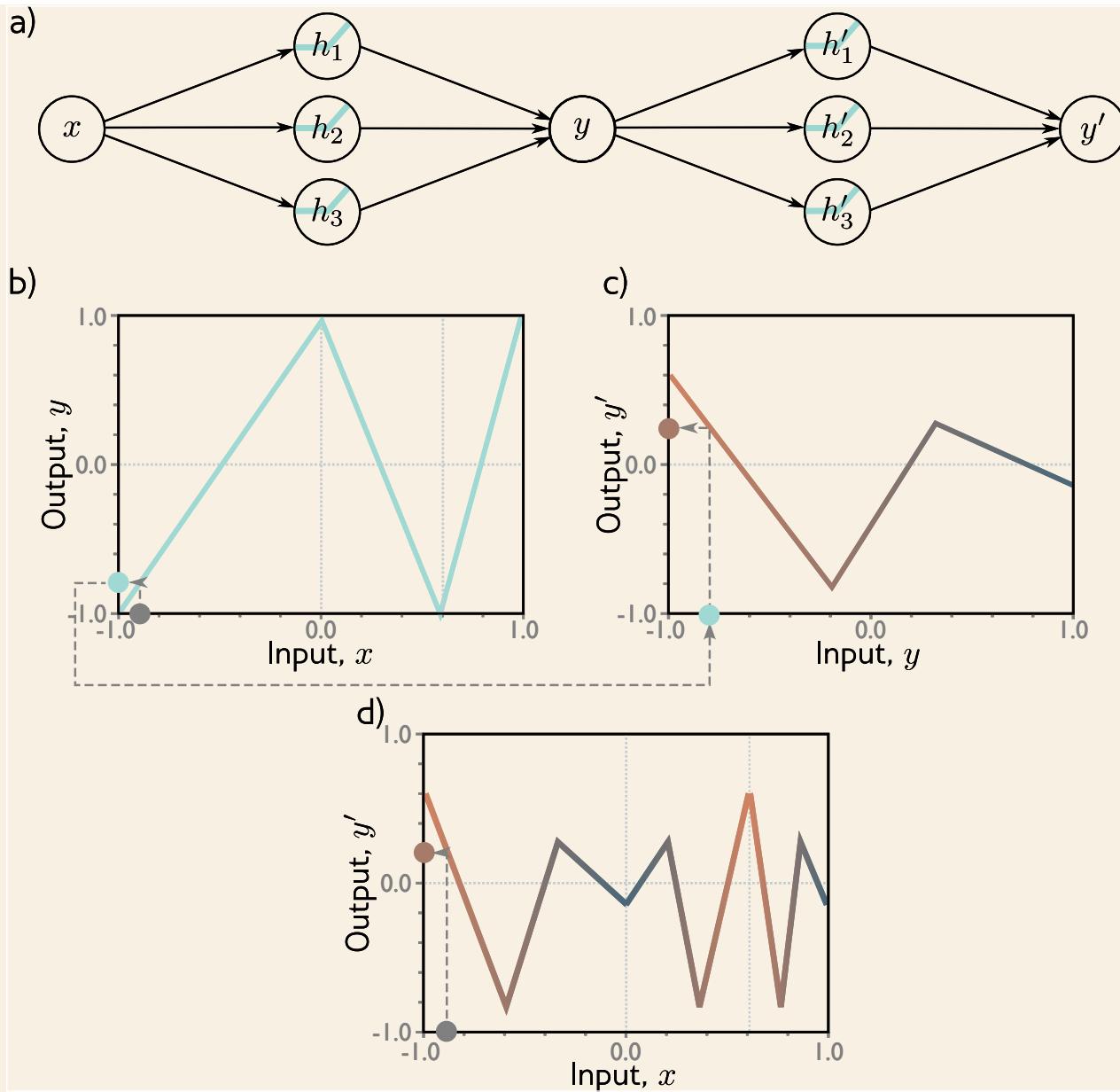
Network 2:

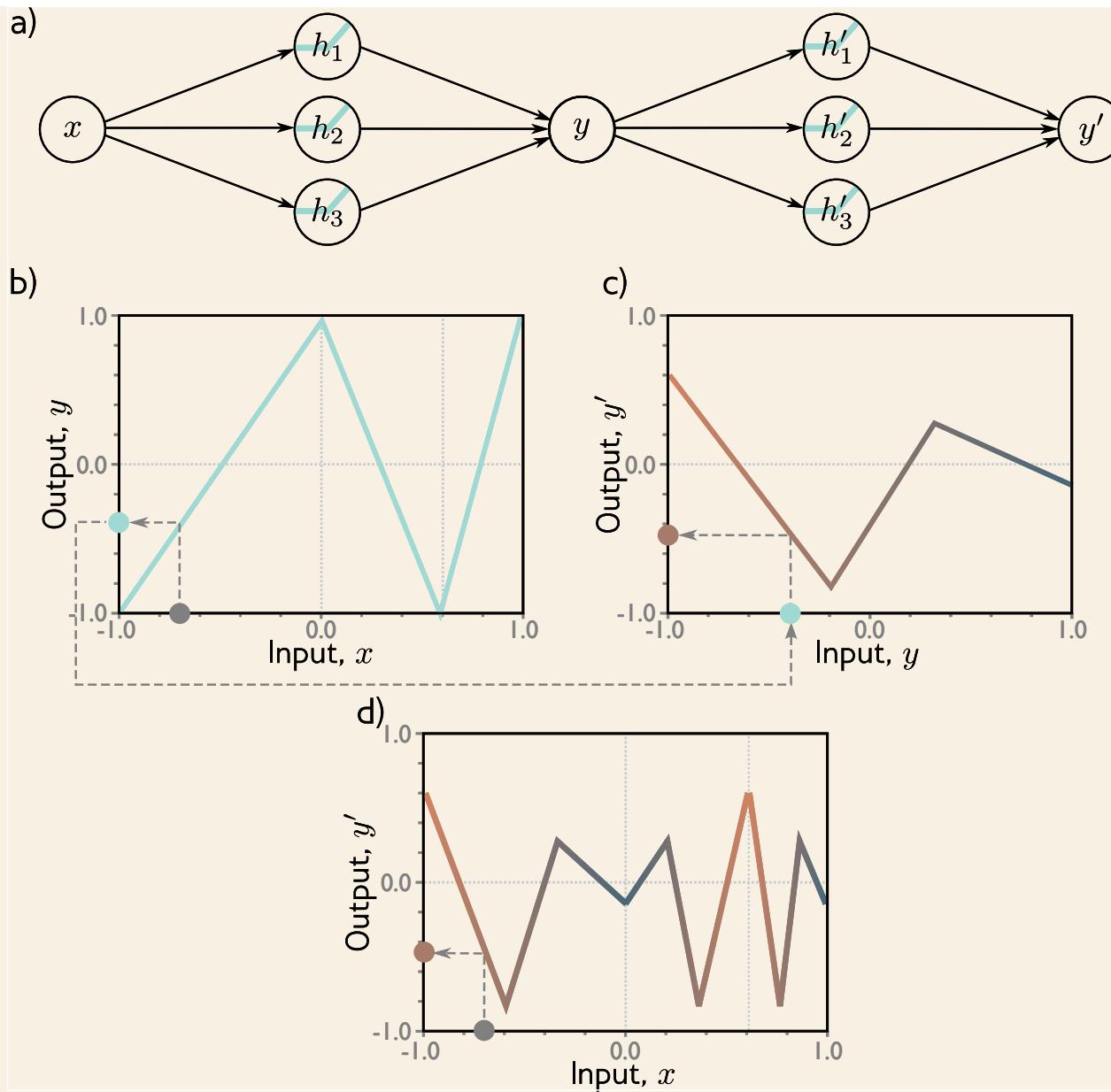
$$\begin{aligned} h'_2 &= a[\theta'_{20} + \theta'_{21}y] & y' &= \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3 \\ h'_3 &= a[\theta'_{30} + \theta'_{31}y] \end{aligned}$$

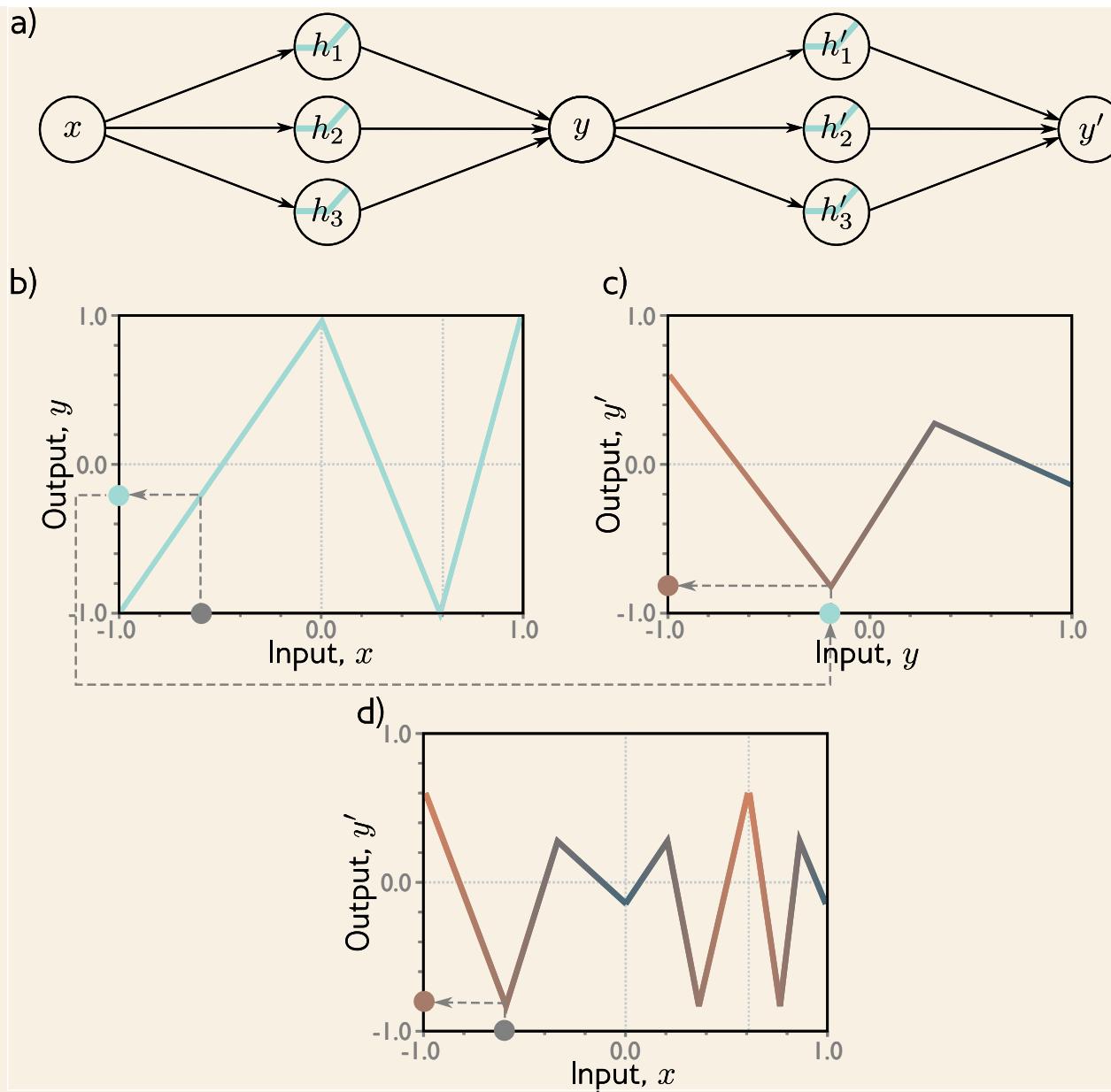


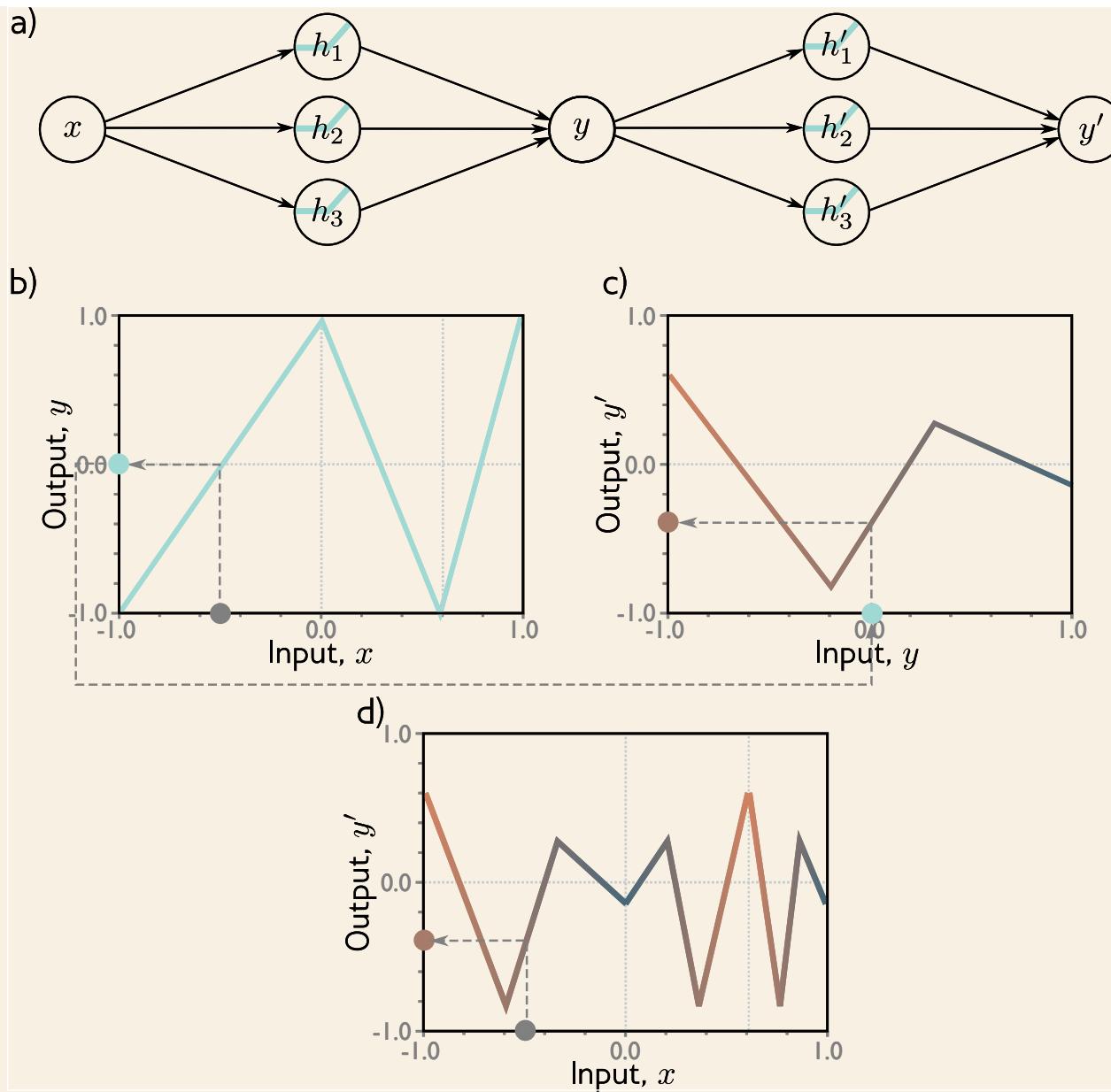


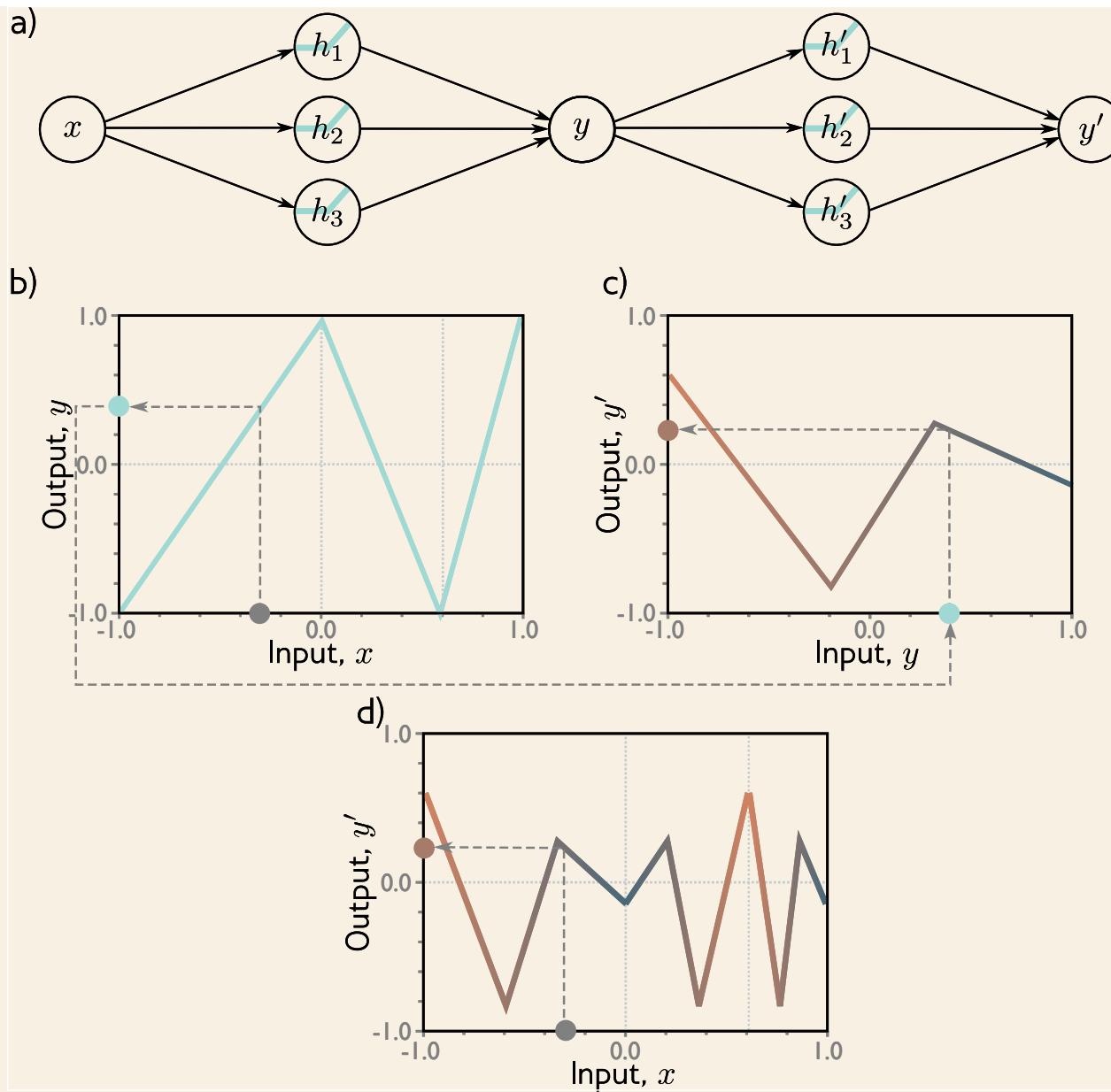


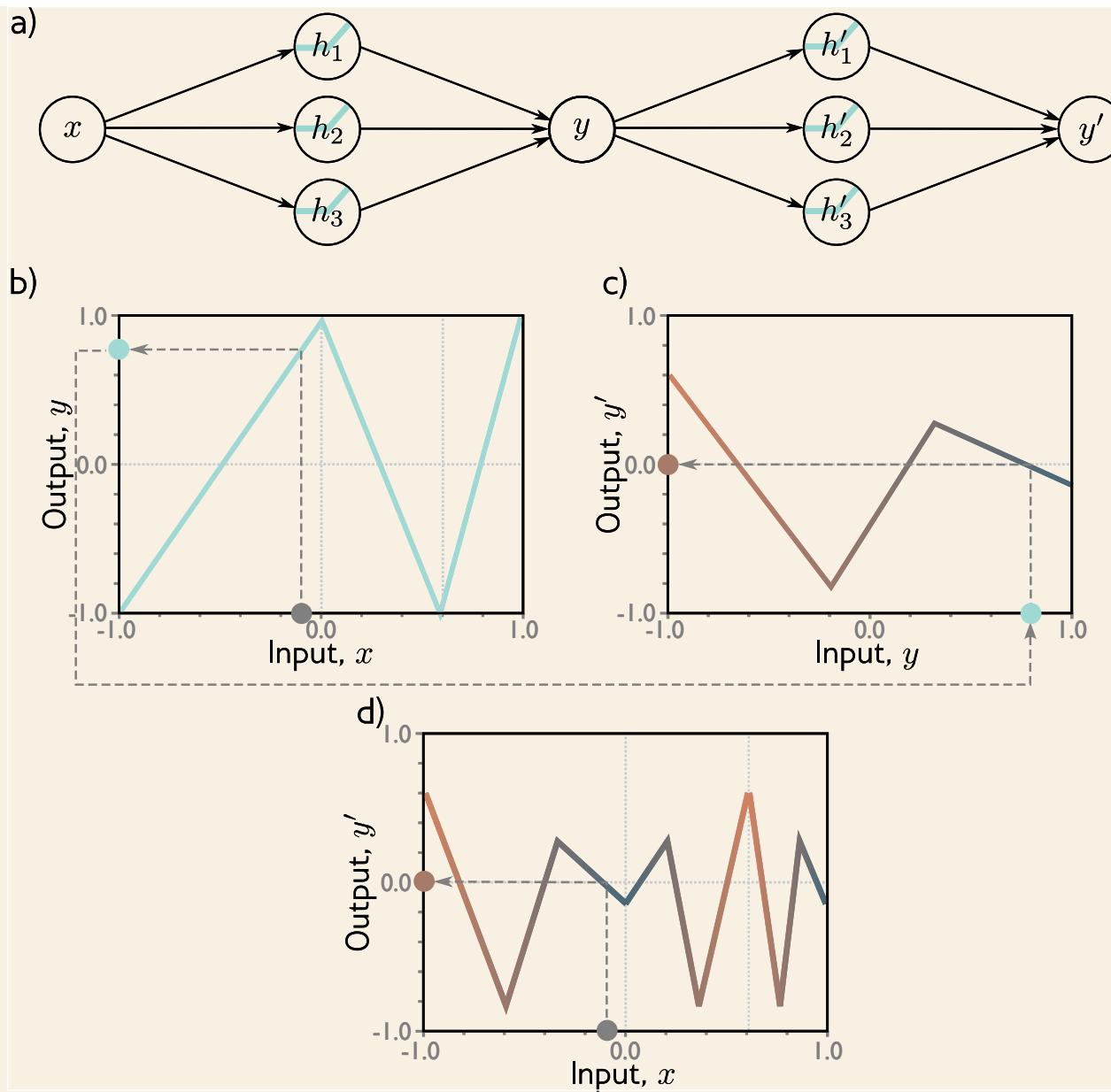


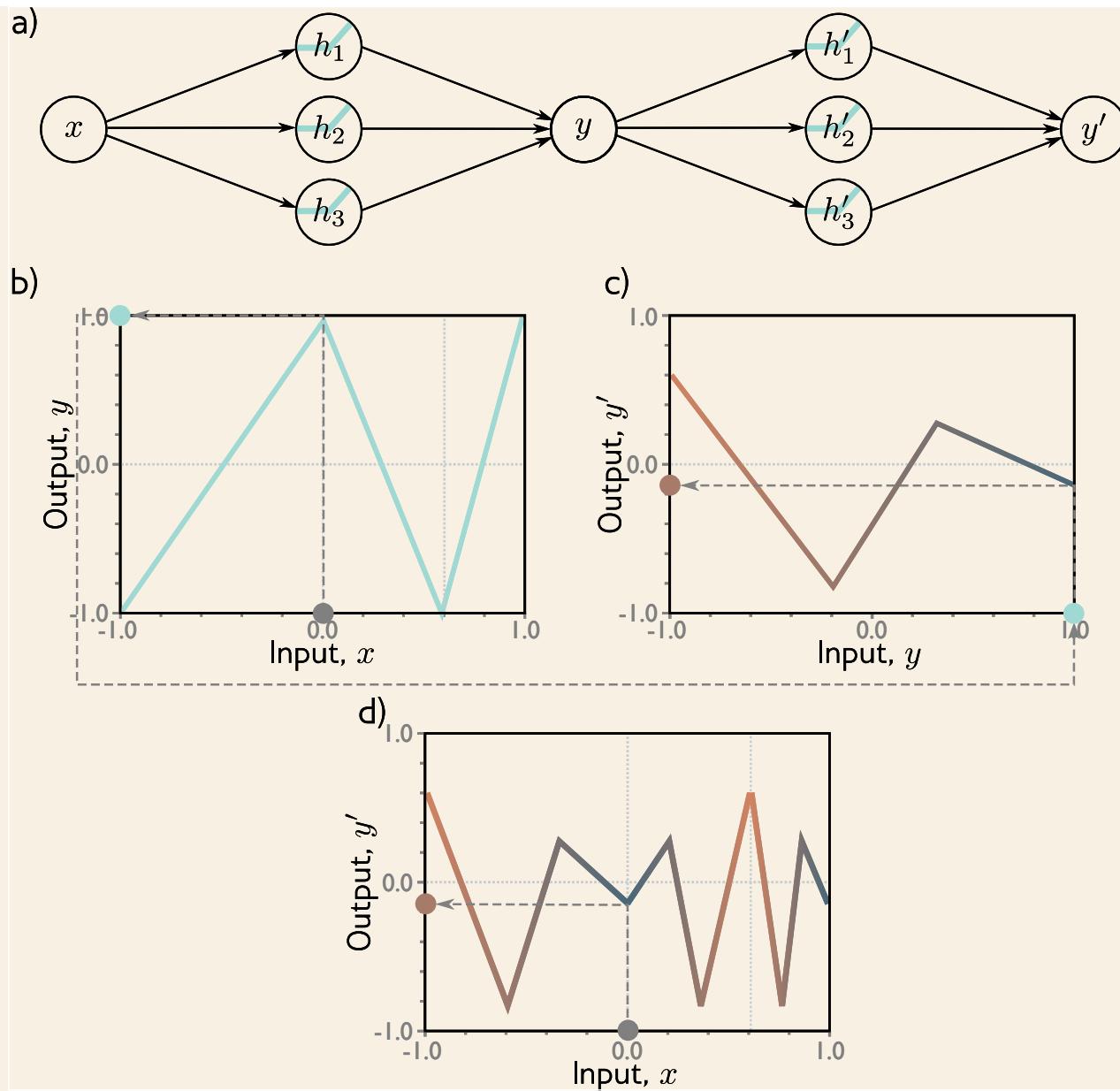


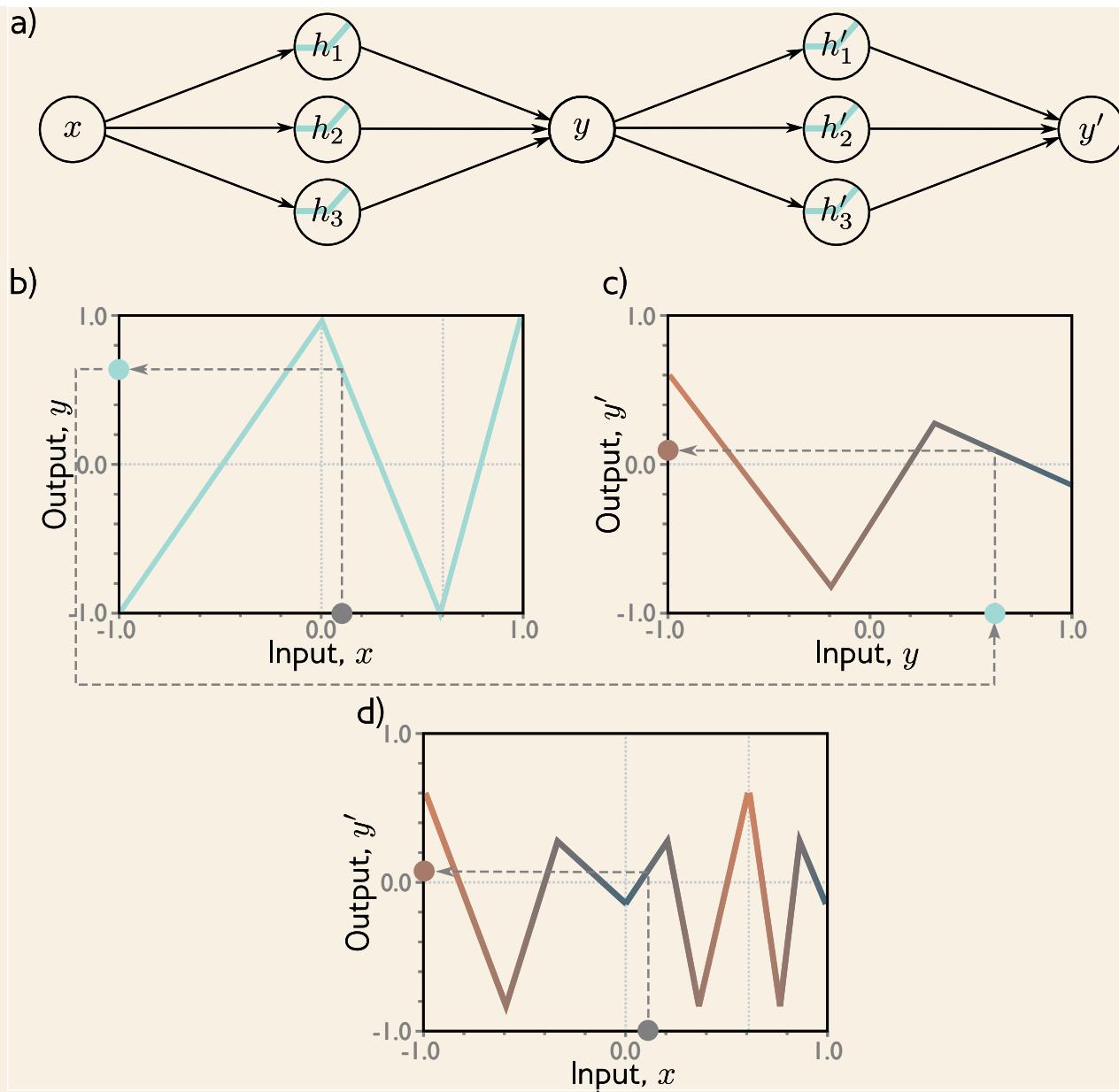


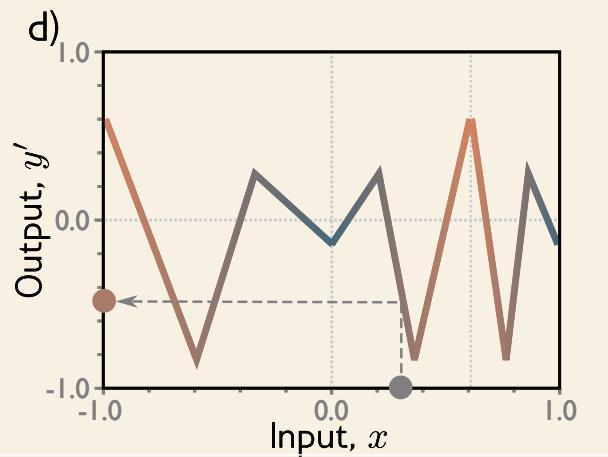
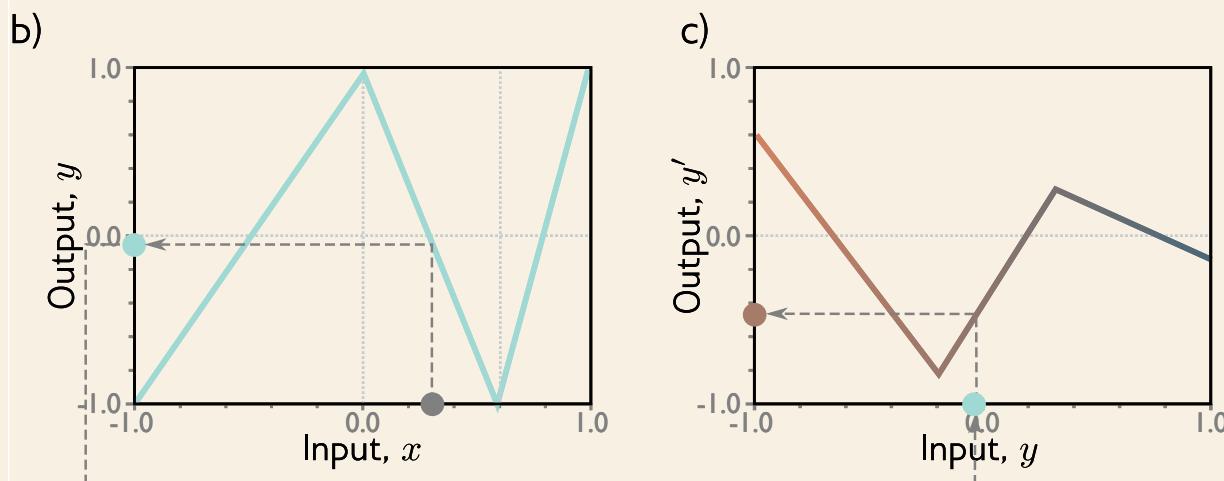
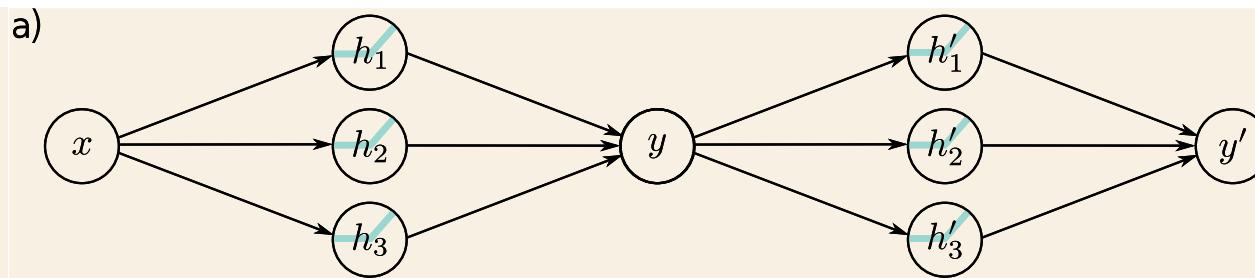


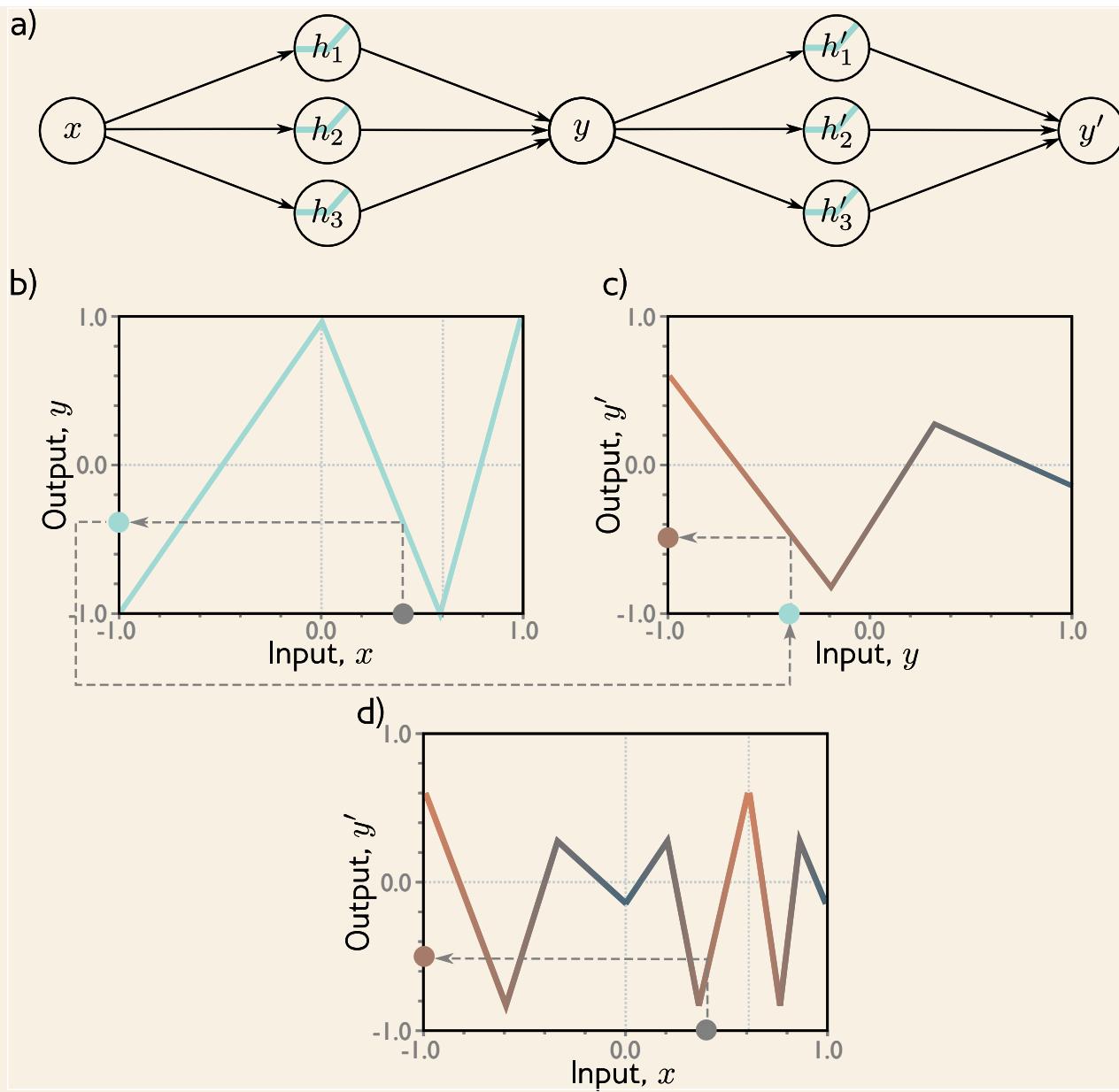


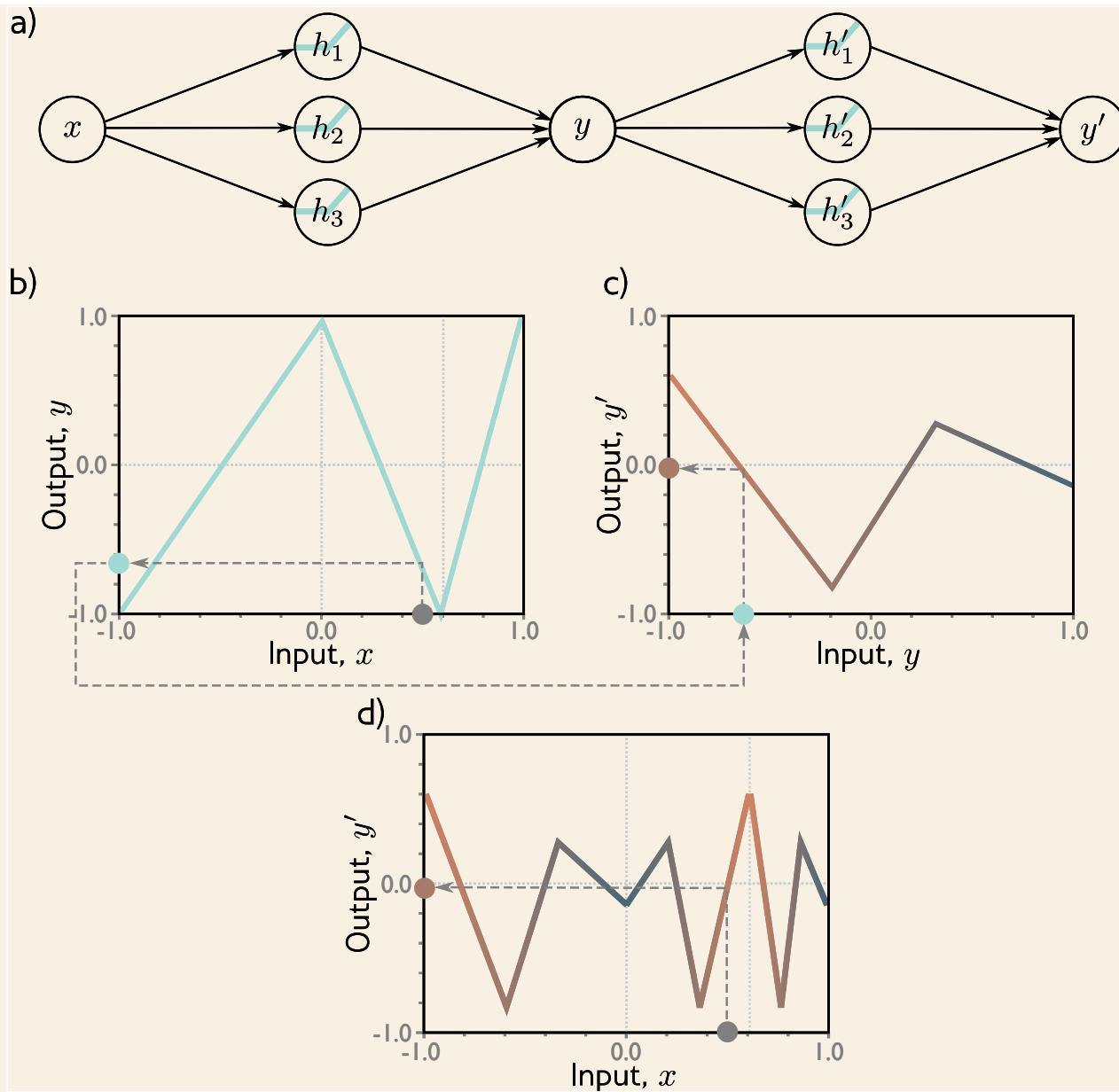


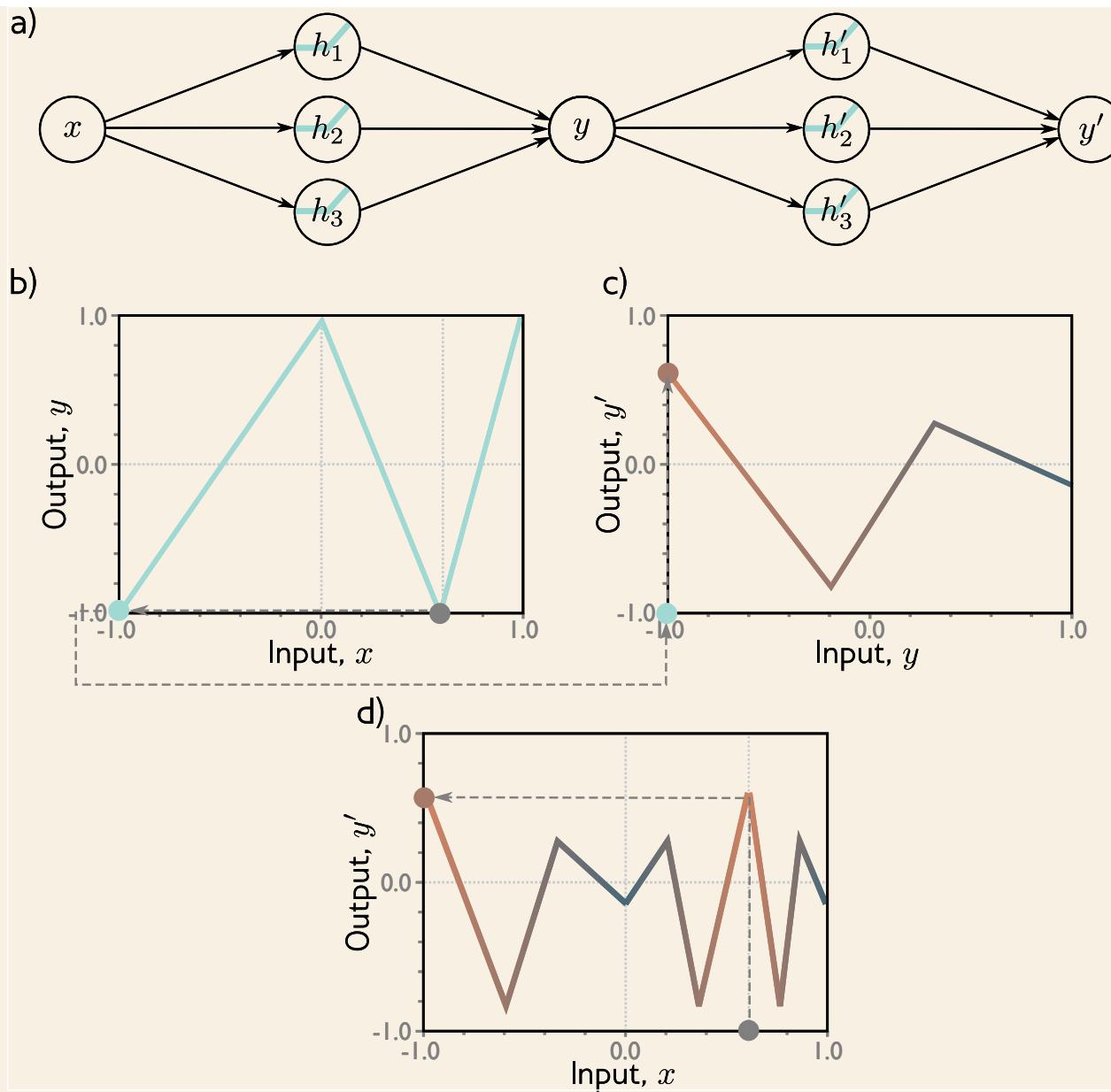








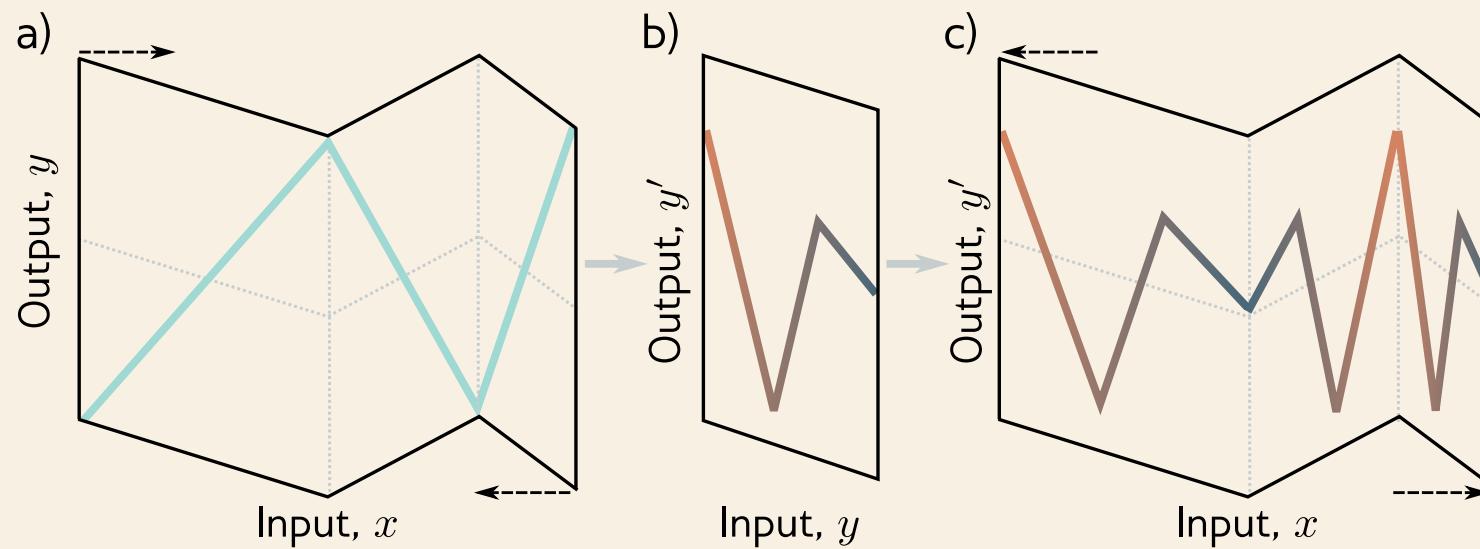




“Folding analogy”

The first network folds the input space x back onto itself so that multiple inputs generate the same output.

The second network applies a function which is replicated at all points that were folded on top of one another.



Combine two networks into one

$$\begin{array}{ll} h_1 = a[\theta_{10} + \theta_{11}x] & \\ \text{Network 1:} & h_2 = a[\theta_{20} + \theta_{21}x] \quad y = \phi_0 + \phi_1 h_1 + \phi_2 h_2 + \phi_3 h_3 \\ & h_3 = a[\theta_{30} + \theta_{31}x] \end{array}$$

$$\begin{array}{ll} h'_1 = a[\theta'_{10} + \theta'_{11}y] & \\ \text{Network 2:} & h'_2 = a[\theta'_{20} + \theta'_{21}y] \quad y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3 \\ & h'_3 = a[\theta'_{30} + \theta'_{31}y] \end{array}$$

Hidden units of second network in terms of first:

$$\begin{aligned} h'_1 &= a[\theta'_{10} + \theta'_{11}y] = a[\theta'_{10} + \theta'_{11}\phi_0 + \theta'_{11}\phi_1 h_1 + \theta'_{11}\phi_2 h_2 + \theta'_{11}\phi_3 h_3] \\ h'_2 &= a[\theta'_{20} + \theta'_{21}y] = a[\theta'_{20} + \theta'_{21}\phi_0 + \theta'_{21}\phi_1 h_1 + \theta'_{21}\phi_2 h_2 + \theta'_{21}\phi_3 h_3] \\ h'_3 &= a[\theta'_{30} + \theta'_{31}y] = a[\theta'_{30} + \theta'_{31}\phi_0 + \theta'_{31}\phi_1 h_1 + \theta'_{31}\phi_2 h_2 + \theta'_{31}\phi_3 h_3] \end{aligned}$$

Create new variables

$$h'_1 = a[\theta'_{10} + \theta'_{11}y] = a[\theta'_{10} + \theta'_{11}\phi_0 + \theta'_{11}\phi_1h_1 + \theta'_{11}\phi_2h_2 + \theta'_{11}\phi_3h_3]$$

$$h'_2 = a[\theta'_{20} + \theta'_{21}y] = a[\theta'_{20} + \theta'_{21}\phi_0 + \theta'_{21}\phi_1h_1 + \theta'_{21}\phi_2h_2 + \theta'_{21}\phi_3h_3]$$

$$h'_3 = a[\theta'_{30} + \theta'_{31}y] = a[\theta'_{30} + \theta'_{31}\phi_0 + \theta'_{31}\phi_1h_1 + \theta'_{31}\phi_2h_2 + \theta'_{31}\phi_3h_3]$$

$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

Two-layer network

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

$$h_2 = a[\theta_{20} + \theta_{21}x]$$

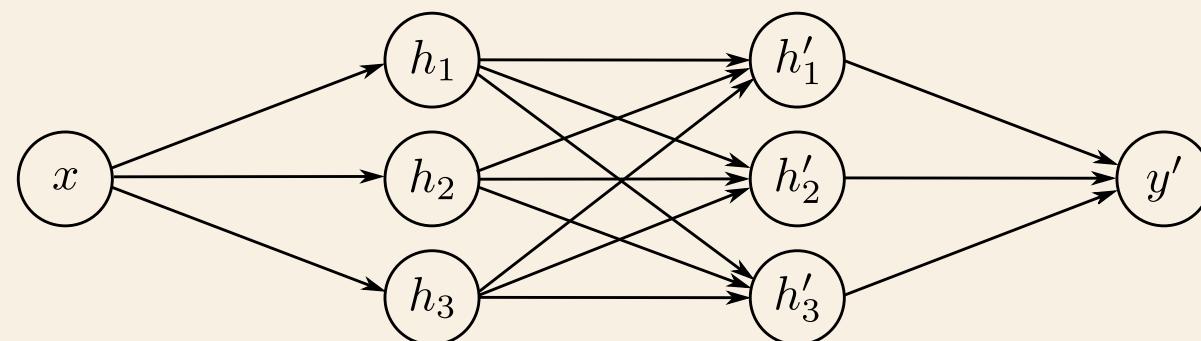
$$h_3 = a[\theta_{30} + \theta_{31}x]$$

$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = a[\psi_{20} + \psi_{21}h_2 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = a[\psi_{30} + \psi_{31}h_2 + \psi_{32}h_2 + \psi_{33}h_3]$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$



Two-layer network as one equation

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

$$h_2 = a[\theta_{20} + \theta_{21}x]$$

$$h_3 = a[\theta_{30} + \theta_{31}x]$$

$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

$$\begin{aligned} y' = & \phi'_0 + \phi'_1 a [\psi_{10} + \psi_{11}a[\theta_{10} + \theta_{11}x] + \psi_{12}a[\theta_{20} + \theta_{21}x] + \psi_{13}a[\theta_{30} + \theta_{31}x]] \\ & + \phi'_2 a [\psi_{20} + \psi_{21}a[\theta_{10} + \theta_{11}x] + \psi_{22}a[\theta_{20} + \theta_{21}x] + \psi_{23}a[\theta_{30} + \theta_{31}x]] \\ & + \phi'_3 a [\psi_{30} + \psi_{31}a[\theta_{10} + \theta_{11}x] + \psi_{32}a[\theta_{20} + \theta_{21}x] + \psi_{33}a[\theta_{30} + \theta_{31}x]] \end{aligned}$$

Networks as composing functions

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

$$h_2 = a[\theta_{20} + \theta_{21}x]$$

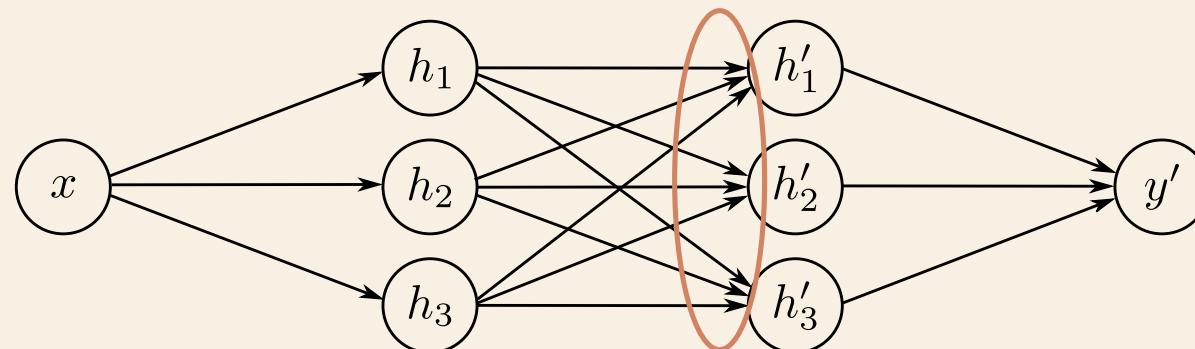
$$h_3 = a[\theta_{30} + \theta_{31}x]$$

$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

Consider the pre-activations at the second hidden units
At this point, it's a one--layer network with three outputs



Networks as composing functions

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

$$h_2 = a[\theta_{20} + \theta_{21}x]$$

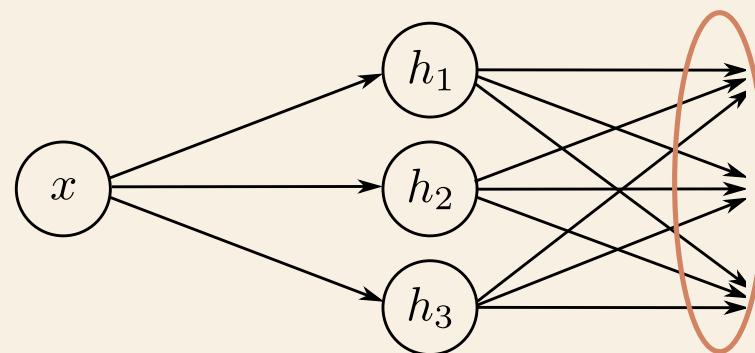
$$h_3 = a[\theta_{30} + \theta_{31}x]$$

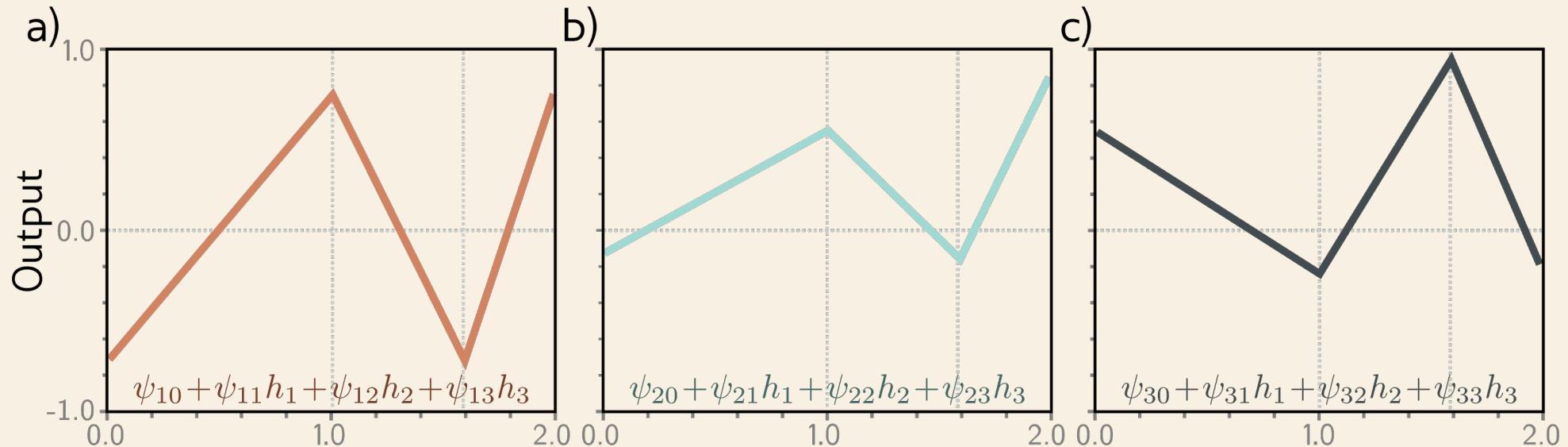
$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

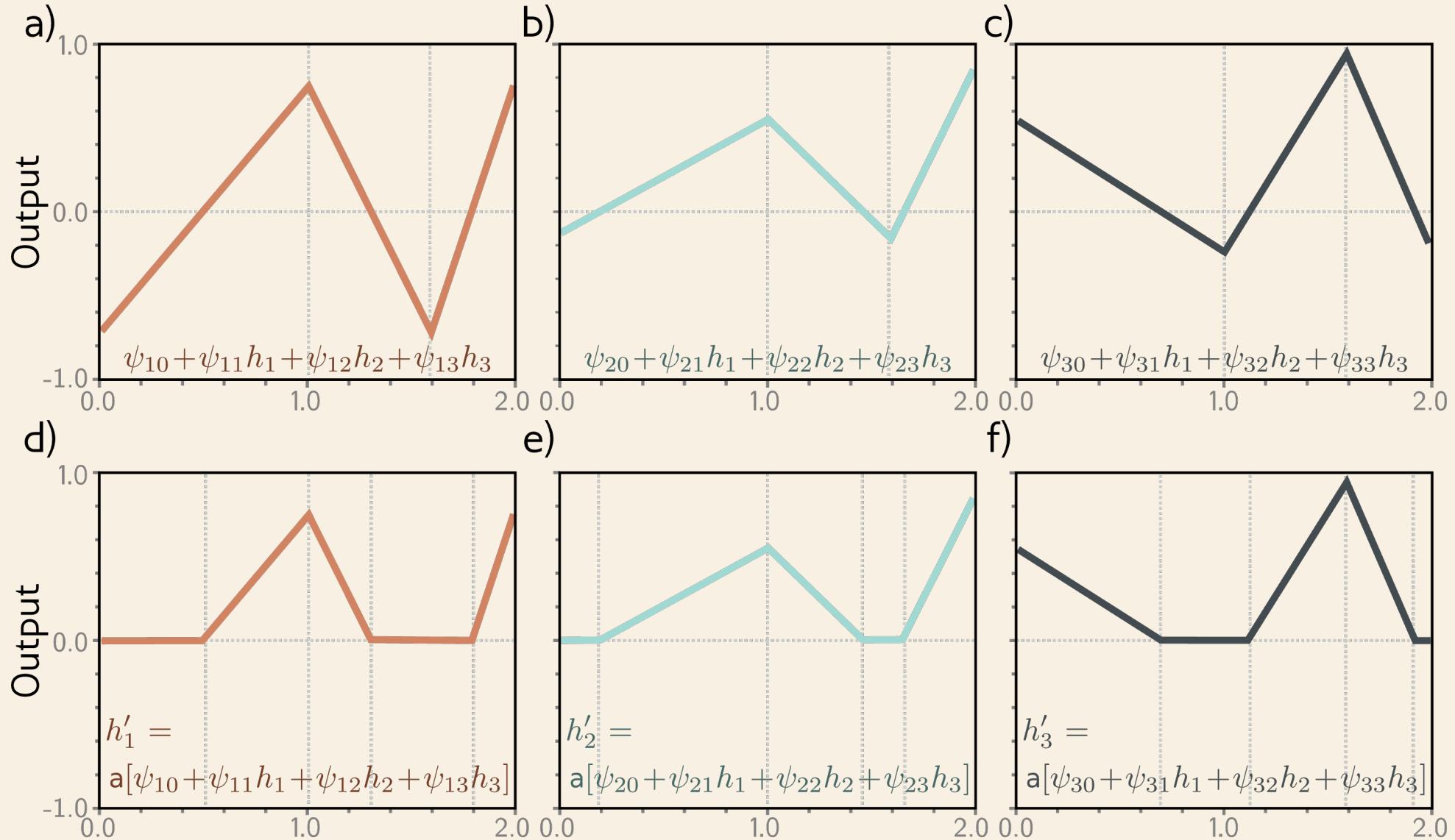
$$h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

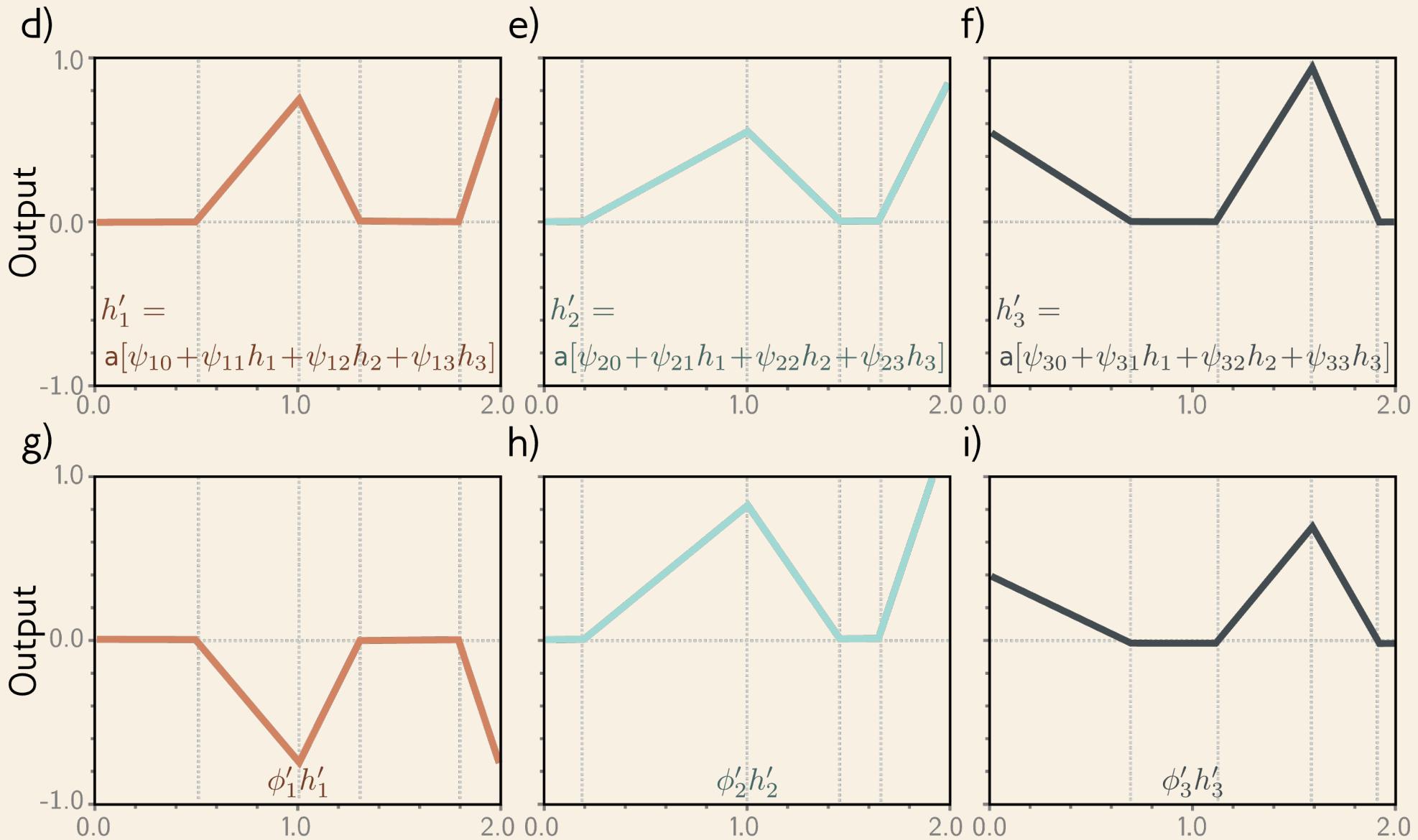
$$h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

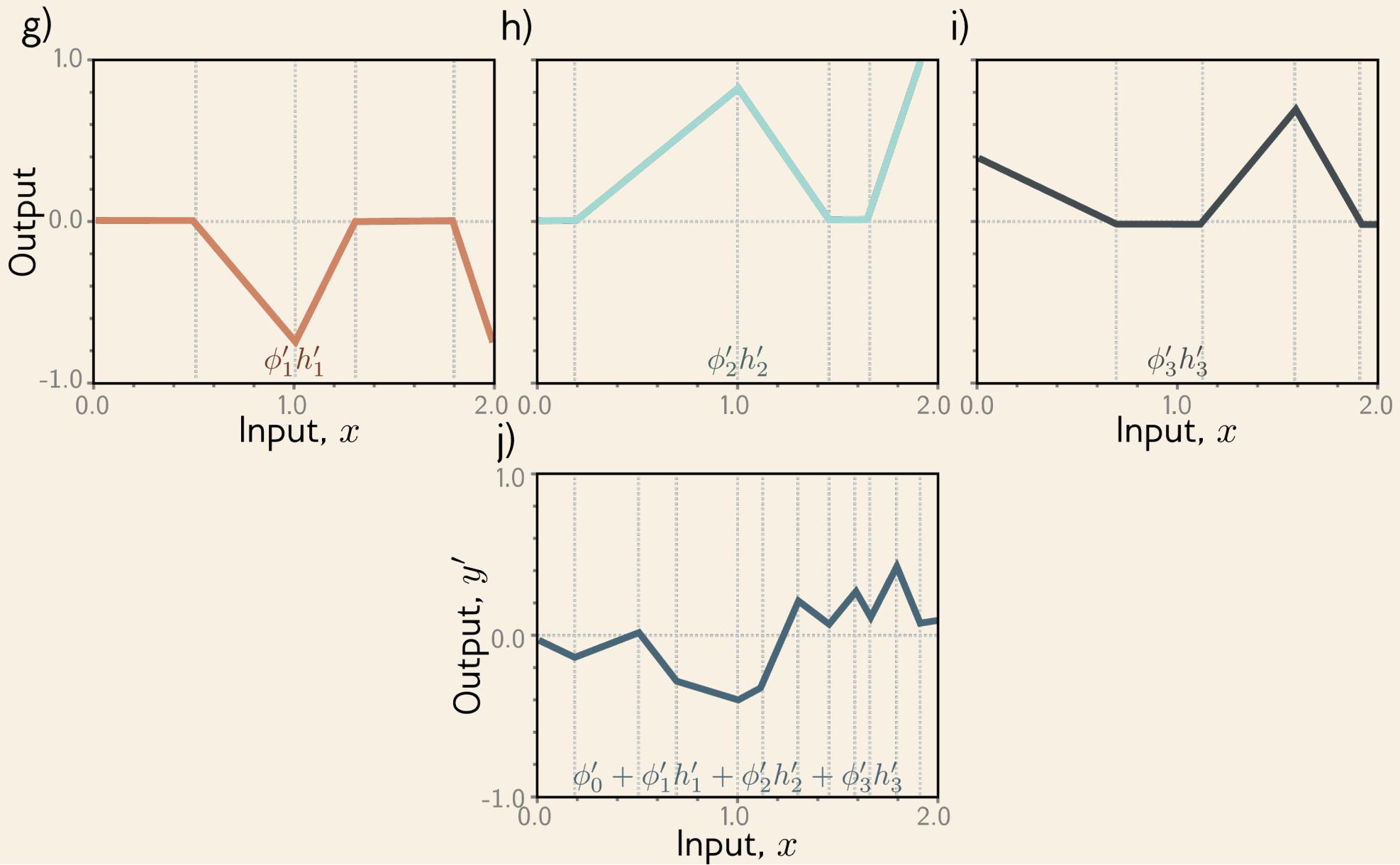
Consider the pre-activations at the second hidden units
At this point, it's a one--layer network with three outputs











Notation change #1

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

$$h_2 = a[\theta_{20} + \theta_{21}x]$$

$$h_3 = a[\theta_{30} + \theta_{31}x]$$

$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

Notation change #1

$$h_1 = \mathbf{a}[\theta_{10} + \theta_{11}x]$$

$$h_2 = \mathbf{a}[\theta_{20} + \theta_{21}x]$$

$$h_3 = \mathbf{a}[\theta_{30} + \theta_{31}x]$$



$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right]$$

$$h'_1 = \mathbf{a}[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = \mathbf{a}[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = \mathbf{a}[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$

Notation change #1

$$h_1 = a[\theta_{10} + \theta_{11}x]$$

$$h_2 = a[\theta_{20} + \theta_{21}x]$$

$$h_3 = a[\theta_{30} + \theta_{31}x]$$



$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right]$$

$$h'_1 = a[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = a[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = a[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$



$$\begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{31} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right]$$

Notation change #1

$$h_1 = \mathbf{a}[\theta_{10} + \theta_{11}x]$$

$$h_2 = \mathbf{a}[\theta_{20} + \theta_{21}x]$$

$$h_3 = \mathbf{a}[\theta_{30} + \theta_{31}x]$$



$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right]$$

$$h'_1 = \mathbf{a}[\psi_{10} + \psi_{11}h_1 + \psi_{12}h_2 + \psi_{13}h_3]$$

$$h'_2 = \mathbf{a}[\psi_{20} + \psi_{21}h_1 + \psi_{22}h_2 + \psi_{23}h_3]$$

$$h'_3 = \mathbf{a}[\psi_{30} + \psi_{31}h_1 + \psi_{32}h_2 + \psi_{33}h_3]$$



$$\begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{32} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right]$$

$$y' = \phi'_0 + \phi'_1 h'_1 + \phi'_2 h'_2 + \phi'_3 h'_3$$



$$y' = \phi'_0 + [\phi'_1 \quad \phi'_2 \quad \phi'_3] \begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix}$$

Notation change #2

$$\begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \theta_{10} \\ \theta_{20} \\ \theta_{30} \end{bmatrix} + \begin{bmatrix} \theta_{11} \\ \theta_{21} \\ \theta_{31} \end{bmatrix} x \right] \longrightarrow \mathbf{h} = \mathbf{a} [\boldsymbol{\theta}_0 + \boldsymbol{\theta}x]$$

$$\begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} = \mathbf{a} \left[\begin{bmatrix} \psi_{10} \\ \psi_{20} \\ \psi_{30} \end{bmatrix} + \begin{bmatrix} \psi_{11} & \psi_{12} & \psi_{13} \\ \psi_{21} & \psi_{22} & \psi_{23} \\ \psi_{32} & \psi_{32} & \psi_{33} \end{bmatrix} \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix} \right] \longrightarrow \mathbf{h}' = \mathbf{a} [\boldsymbol{\psi}_0 + \boldsymbol{\Psi}\mathbf{h}]$$

$$y' = \phi'_0 + [\phi'_1 \quad \phi'_2 \quad \phi'_3] \begin{bmatrix} h'_1 \\ h'_2 \\ h'_3 \end{bmatrix} \longrightarrow y = \boldsymbol{\phi}'_0 + \boldsymbol{\phi}'\mathbf{h}'$$

Notation change #3

$$\mathbf{h} = \mathbf{a} [\theta_0 + \theta x] \longrightarrow \mathbf{h}_1 = \mathbf{a} [\beta_0 + \Omega_0 \mathbf{x}]$$

$$\mathbf{h}' = \mathbf{a} [\psi_0 + \Psi \mathbf{h}] \longrightarrow \mathbf{h}_2 = \mathbf{a} [\beta_1 + \Omega_1 \mathbf{h}_1]$$

$$y = \phi'_0 + \phi' \mathbf{h}' \longrightarrow \mathbf{y} = \beta_2 + \Omega_2 \mathbf{h}_2$$

Notation change #3

$$\mathbf{h} = \mathbf{a} [\theta_0 + \theta x]$$



Bias vector
Weight matrix

$$\mathbf{h}_1 = \mathbf{a} [\beta_0 + \Omega_0 \mathbf{x}]$$

$$\mathbf{h}' = \mathbf{a} [\psi_0 + \Psi \mathbf{h}]$$



$$\mathbf{h}_2 = \mathbf{a} [\beta_1 + \Omega_1 \mathbf{h}_1]$$

$$y = \phi'_0 + \phi' \mathbf{h}'$$



$$\mathbf{y} = \beta_2 + \Omega_2 \mathbf{h}_2$$

General equations for deep network

$$\mathbf{h}_1 = \mathbf{a}[\boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0 \mathbf{x}]$$

$$\mathbf{h}_2 = \mathbf{a}[\boldsymbol{\beta}_1 + \boldsymbol{\Omega}_1 \mathbf{h}_1]$$

$$\mathbf{h}_3 = \mathbf{a}[\boldsymbol{\beta}_2 + \boldsymbol{\Omega}_2 \mathbf{h}_2]$$

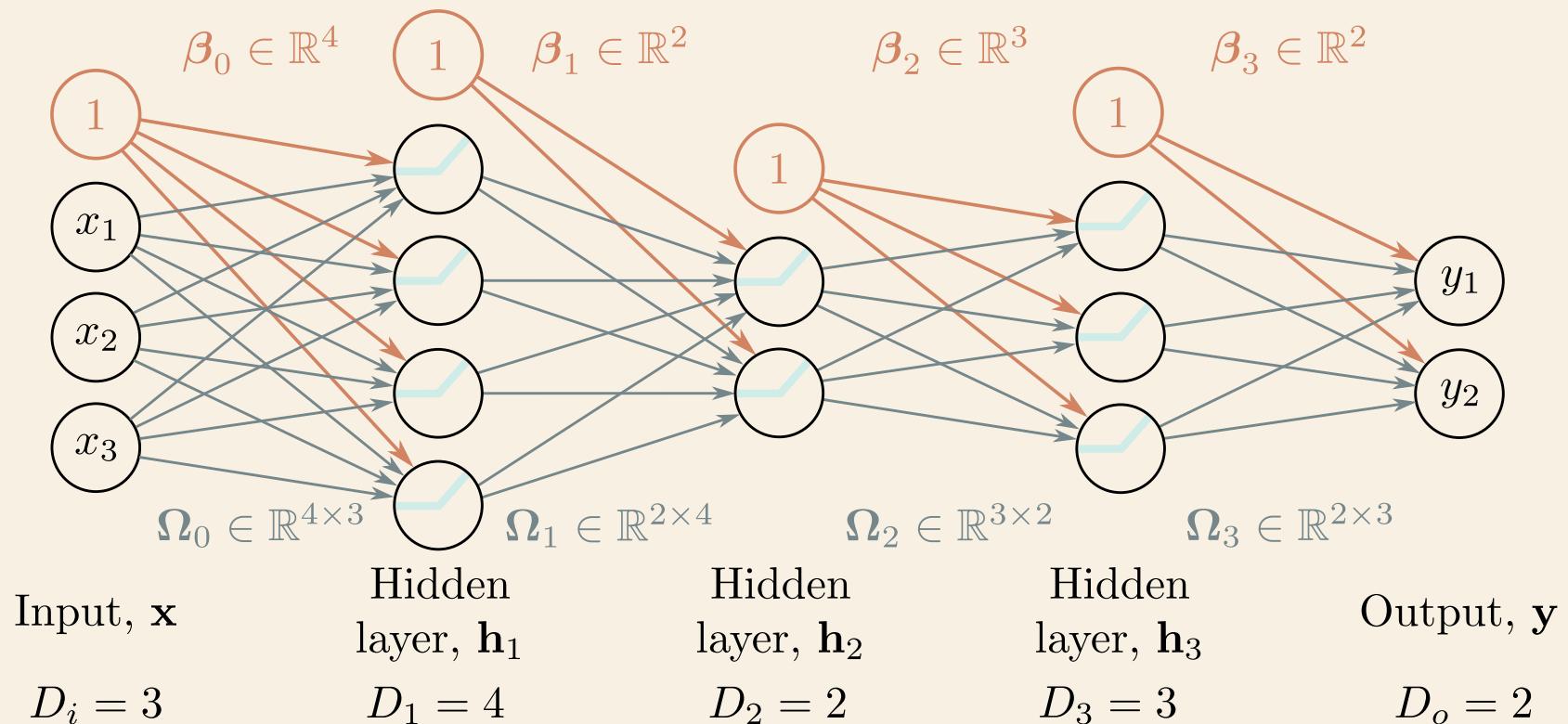
⋮

$$\mathbf{h}_K = \mathbf{a}[\boldsymbol{\beta}_{K-1} + \boldsymbol{\Omega}_{K-1} \mathbf{h}_{K-1}]$$

$$\mathbf{y} = \boldsymbol{\beta}_K + \boldsymbol{\Omega}_K \mathbf{h}_K,$$

$$\mathbf{y} = \boldsymbol{\beta}_K + \boldsymbol{\Omega}_K \mathbf{a} [\boldsymbol{\beta}_{K-1} + \boldsymbol{\Omega}_{K-1} \mathbf{a} [\dots \boldsymbol{\beta}_2 + \boldsymbol{\Omega}_2 \mathbf{a} [\boldsymbol{\beta}_1 + \boldsymbol{\Omega}_1 \mathbf{a} [\boldsymbol{\beta}_0 + \boldsymbol{\Omega}_0 \mathbf{x}]] \dots]]$$

Example

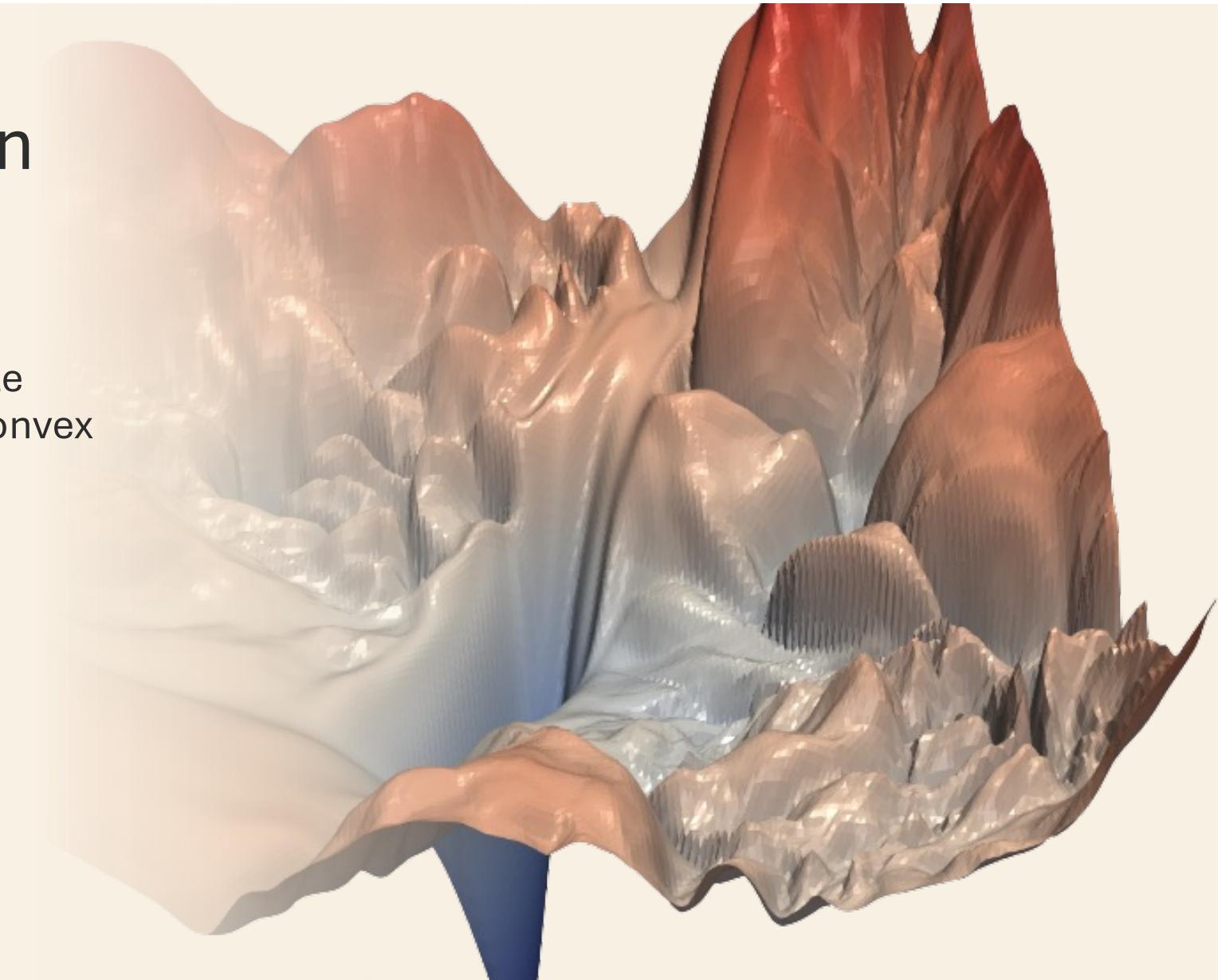


Optimization Challenges

Local minima, saddle points due to non-convex loss

Narrow valleys

Exploding/vanishing gradient



Minibatch Gradient descent

- While not converged, on dev, sample data in pieces (batches) updating model
- Training Set
- EPOCH

$$\theta = \theta - \mu \nabla Loss(\text{batch 1})$$

$$\theta = \theta - \mu \nabla Loss(\text{batch 2})$$

$$\theta = \theta - \mu \nabla Loss(\text{batch 3})$$

$$\theta = \theta - \mu \nabla Loss(\text{batch 4})$$

Momentum-Based Optimization

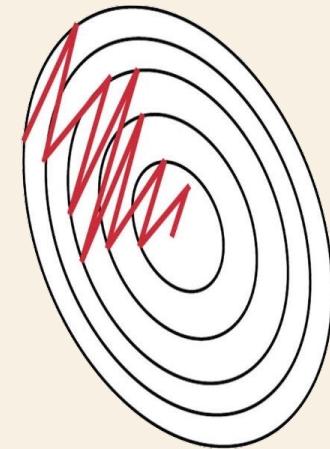
Gradient Descent:

$$\theta \leftarrow \theta - \alpha \nabla_{\theta} L(f_{\theta}(x), y)$$

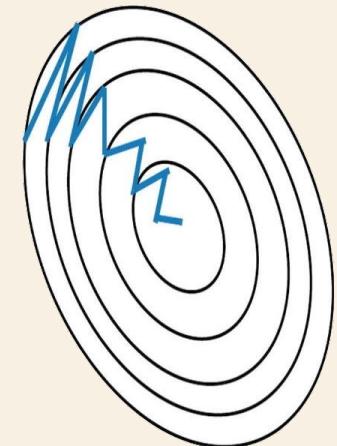
Momentum Update:

$$\begin{aligned}\rho &\leftarrow \mu \rho - \alpha \nabla_{\theta} L \\ \theta &\leftarrow \theta + \rho\end{aligned}$$

- Momentum accumulates direction of past gradients
- Helps escape small local minima and narrow valleys
- Typical $\mu \approx 0.9 - 0.99$



Stochastic Gradient
Descent **without**
Momentum



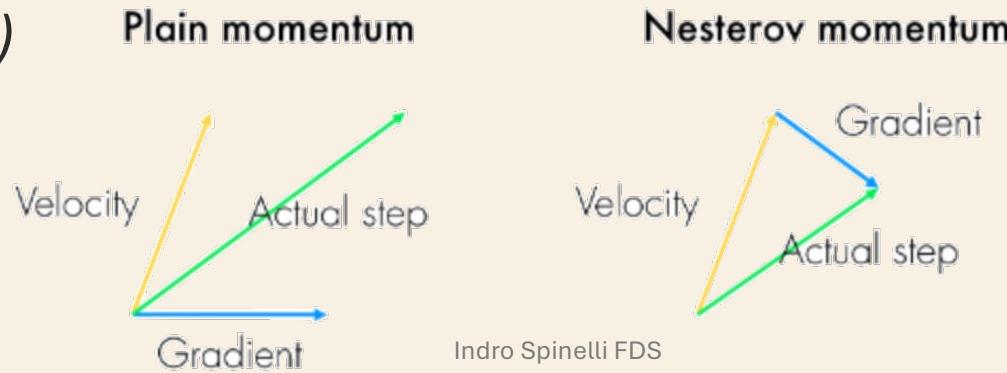
Stochastic Gradient
Descent **with**
Momentum

Nesterov Momentum

Lookahead gradient: compute gradient at future position

$$\rho \leftarrow \mu\rho - \alpha \nabla_{\theta} L(f_{\theta+\mu\rho}(x), y)$$

- Prevents overshooting when near optimum
- Often converges faster and smoother than vanilla momentum
- (*Visual: arrows showing “look-ahead” step vs standard momentum*)



Adaptive Learning Rates

Key idea: adjust the step size per parameter

AdaGrad: divides update by accumulated squared gradients → good for sparse features

$$\begin{aligned} g_t &= \nabla_{\theta} L(f_{\theta}(x_t), y_t) \\ G_t &= G_{t-1} + g_t \odot g_t \\ \theta_t &= \theta_{t-1} - \frac{\alpha}{\sqrt{G_t + \epsilon}} \odot g_t \end{aligned}$$

Adaptive Learning Rates

Key idea: adjust the step size per parameter

AdaGrad: divides update by accumulated squared gradients → good for sparse features

$$\begin{aligned} g_t &= \nabla_{\theta} L(f_{\theta}(x_t), y_t) \\ G_t &= G_{t-1} + g_t^2 \\ \theta_t &= \theta_{t-1} - \frac{\alpha}{\sqrt{G_t + \epsilon}} g_t \end{aligned}$$

g_t : gradient at step t

G_t : element-wise sum of squared gradients

ϵ : small constant for numerical stability

Adaptive Learning Rates

RMSProp: Adapts learning rate per parameter, avoiding AdaGrad's rapid decay. uses exponential moving average instead of full history

$$g_t = \nabla_{\theta} L(f_{\theta}(x_t), y_t)$$

$$E[g^2]_t = \gamma E[g^2]_{t-1} + (1 - \gamma) g_t^2$$

$$\theta_t = \theta_{t-1} - \frac{\alpha}{\sqrt{E[g^2]_t + \epsilon}} g_t$$

γ : decay rate (typically 0.9)

$E[g^2]_t$: exponentially decaying average of past squared gradients

Adaptive Learning Rates

Adam: Combines momentum (*first moment*) and adaptive learning rate (*second moment*)

$$\begin{aligned} g_t &= \nabla_{\theta} L(f_{\theta}(x_t), y_t) \\ m_t &= \beta_1 m_{t-1} + (1 - \beta_1) g_t \\ v_t &= \beta_2 v_{t-1} + (1 - \beta_2) g_t^2 \\ \hat{m}_t &= \frac{m_t}{1 - \beta_1^t}, \hat{v}_t = \frac{v_t}{1 - \beta_2^t} \\ \theta_t &= \theta_{t-1} - \frac{\alpha}{\sqrt{\hat{v}_t} + \epsilon} \hat{m}_t \end{aligned}$$

Typical values: $\beta_1 = 0.9$, $\beta_2 = 0.999$, $\epsilon = 10^{-8}$

Default optimizer in many deep learning frameworks

Learning Rate Schedules

The *most critical hyperparameter!!!*

Use validation loss to guide tuning

- *Decay schedules*: reduce LR when progress stalls
 - Popular strategies:
 - Step decay
 - Exponential decay
 - Cosine annealing

Weight Initialization

Why initialization matters:

- Prevents gradients from exploding or vanishing
- Ensures all neurons start from diverse states

Bad practice: zero initialization → identical neurons

Long history of trick based on “fan in” aka input dimension d_{in}

Xavier (Glorot) Init: for tanh activations $W_j \sim \mathcal{N}\left(0, \frac{1}{d_{in} + d_{out}}\right)$

Kaiming (He) Init: for ReLU activations $W_j \sim \mathcal{N}\left(0, \frac{2}{d_{in}}\right)$

Activation Functions: Why They Matter

Purpose: introduce non-linearity to neural networks

- Enable modeling of complex input–output mappings
- Affect both *gradient flow* and *training stability*
- Choice of activation impacts network depth and performance

Issues with sigmoid

Sigmoid function:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

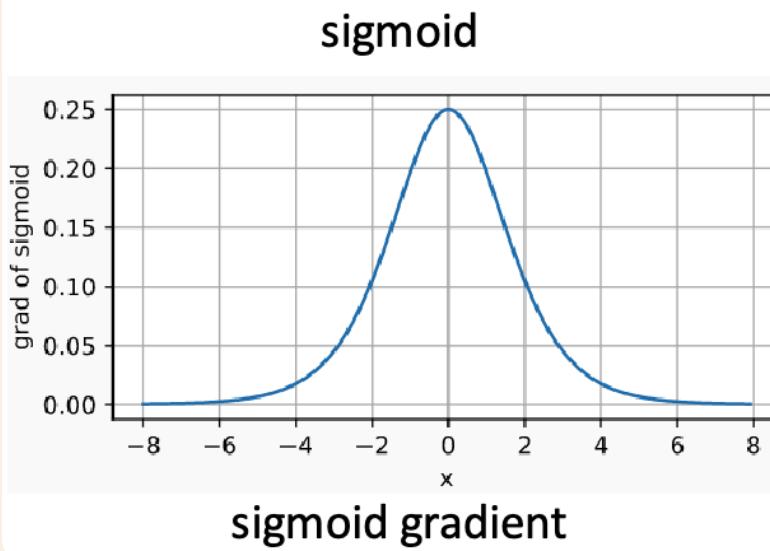
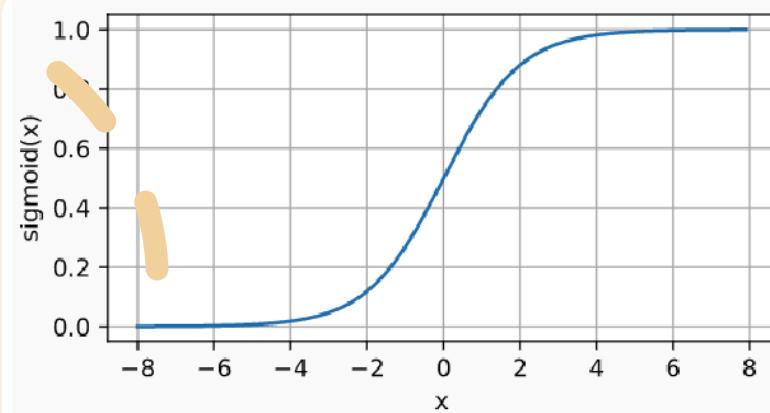
Derivative:

$$\sigma'(x) = \sigma(x)(1 - \sigma(x))$$

Problem: $\sigma'(x) \leq 0.25$, so gradient multiplies by ≤ 0.25 at each layer.

Consequence: In deep networks, gradients shrink exponentially \rightarrow early layers learn very slowly.

Example: For 10 layers, max gradient $\approx 0.25^{10} \approx 10^{-6}$



Relu as solution

ReLU activation: $\text{ReLU}(x) = \max(0, x)$

Derivative: $f'(x) = \begin{cases} 1 & x > 0 \\ 0 & x \leq 0 \end{cases}$

Why it solves vanishing gradient:

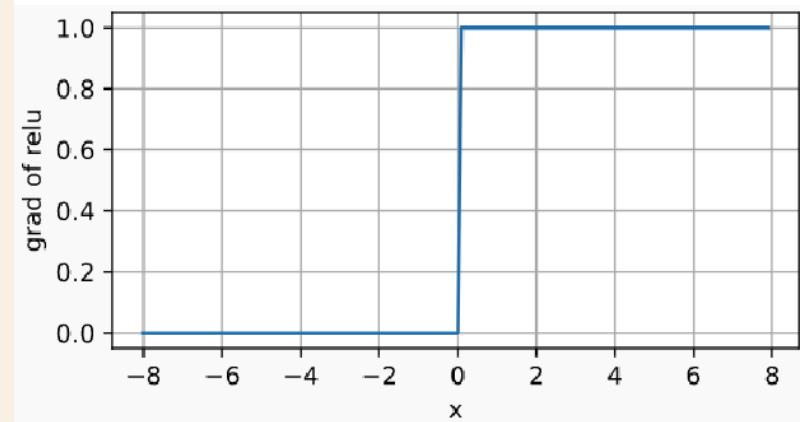
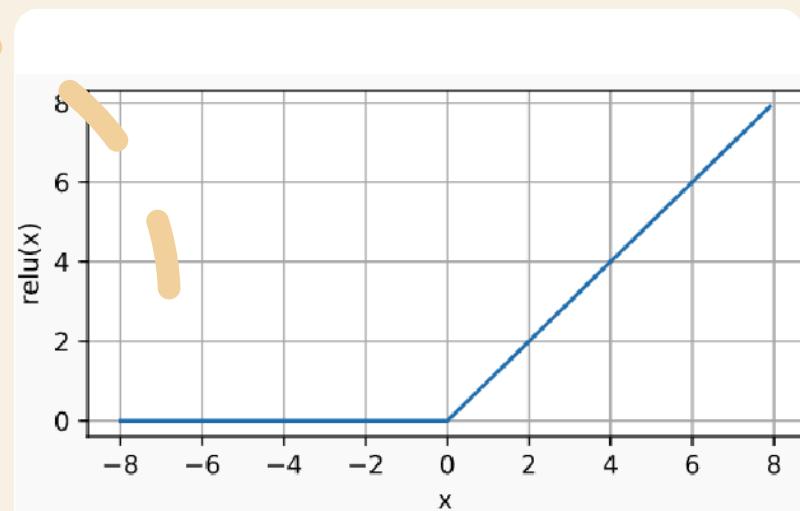
Derivative = 1 for positive activations \rightarrow gradient does **not shrink** in backprop.

Enables deep networks to train efficiently.

Sparse activations \rightarrow more efficient computation.

Simpler gradient propagation \rightarrow faster convergence.

Caveat: “Dying ReLU” problem (neurons stuck at 0).
Can be mitigated with **Leaky**



Residual Connections

- **Introduction:** Introduced in Residual Networks (ResNets) to solve the vanishing gradient problem.
- **Core Idea:** Instead of forcing a layer to learn a complete transformation $H(x)$, we have it learn a *residual* (or *change*) $F(x)$ to *add* to the input x .

$$h = F(x) + x$$

Intuition: It's easier for a layer to learn to "do nothing" (just output $F(x) = 0$) than to learn the identity function ($H(x) = x$). This "skip connection" passes the original input x straight through.

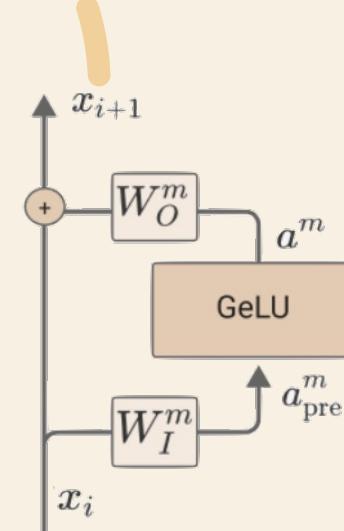
Residual Connections

During backpropagation:

$$\nabla_x(F(x) + x) = \nabla_x F(x) + 1$$

The **+1** creates a direct, uninterrupted "highway" for gradients to flow back through the network.

- Gradients can propagate through hundreds of layers without vanishing.
- The network can **easily learn the identity function** ($F(x) = 0$), so adding more layers doesn't hurt performance.



$$x_{i+1} = x_i + W_O^m a^m$$

$$a^m = \text{GeLU}(a_{\text{pre}}^m)$$

$$a_{\text{pre}}^m = W_I^m x_i$$

Deep Residual Networks (ResNet)

What they are: ResNets are built by stacking these residual units. It is typically a block of 2-3 convolutional layers (next week).

The Results (ImageNet):

- A 152-layer ResNet achieved 3.57% top-5 error, winning ImageNet 2015.
- This was the first time a model surpassed human-level performance (which was ~5.1%).
- Researchers were able to train ResNets over 1000 layers deep.

Key Takeaway: Residual connections finally allowed for the successful training of *extremely* deep neural networks.

Refresher: Standard Scaling

What is it? A common pre-processing step for an *entire dataset X*.

Goal: Transform each feature (column) to have a zero mean and unit variance.

$$X' = (X - \mu)/\sigma$$

- μ is the mean of each feature.
- σ is the standard deviation of each feature.

Helps optimization algorithms converge faster.

Batch Normalization (BN): The Concept

The Idea: What if we applied this same "Standard Scaling" principle *inside* the network, for every layer's output?

The Problem:

- We don't have access to the *entire* dataset's statistics for each layer.
- These statistics would "shift" during training anyway as the network's parameters change (this is called "Internal Covariate Shift").

The Solution: Normalize the layer's output using the statistics (mean and variance) of the **current mini-batch**.

Batch Normalization (BN): During Training

Reduce **internal covariate shift** → stabilize and accelerate training.
Normalizes layer inputs across the **batch**:

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

- Then scales and shifts:

$$y_i = \alpha \hat{x}_i + \beta$$

- Keeps mean ≈ 0 and variance ≈ 1 for each layer input
- Allows **higher learning rates**
- They allow the network to **adapt the mean and variance** for optimal performance.

Batch Normalization (BN): During Inference

Problem: We can't use mini-batch stats at test time. The output for one sample shouldn't depend on what other samples are in its "batch".

Standard Solution: Moving Averages

- During training, the network maintains a **moving average** of the mini-batch means (μ) and variances (σ^2).
- At inference time, these final, fixed averages are used to normalize the test input.

$$\hat{x}_i = \frac{x_i - \mu_{\text{avg}}}{\sqrt{\sigma_{\text{avg}}^2 + \epsilon}} \quad y_i = \alpha \hat{x}_i + \beta$$

Batch Normalization: Challenges

Mini-Batch Dependency:

- BN creates a dependency between samples in a mini-batch.
- This can be problematic for certain tasks (e.g., contrastive learning) or distributed training.

Fails on Small Batch Sizes:

- If the batch size is very small (e.g., 1 or 2), the mean and variance estimates are extremely noisy and unstable.
- This is a major problem for large models (like in NLP) that require small batches to fit in memory.

Layer Normalization

Normalizes **all features of a single sample** rather than across the batch:

$$\hat{x}_i = \frac{x_i - \mu_L}{\sqrt{\sigma_L^2 + \epsilon}}$$

- Then scales and shifts:

$$y_i = \gamma \hat{x}_i + \beta$$

- Works well for **RNNs or small batch sizes**
- Gradient flow is stabilized **per sample**, independent of batch size
- Improves training stability in deep architectures

Dropout Regularization: The Core Idea

Goal: A regularization technique to prevent overfitting.

How it Works: During training, randomly "drop" (set to 0) a fraction of neurons at a given layer.

Hyperparameter p : The "dropout rate" (e.g., $p = 0.5$), or the probability a neuron is dropped. The "keep probability" is $(1 - p)$.

- **Intuition:**

- Forces the network to learn redundant representations.
- Prevents neurons from "co-adapting" or relying too much on specific other neurons.
- It's like training many different "thinned" networks in parallel.

The Standard Implementation (Inverted Dropout)

The Problem: If we *only* drop neurons at train time, the layer's output will be much larger at test time (since all neurons are active).

Naïve Solution: At test time, scale all activations down by the keep probability $(1 - p)$.

$$y = (1 - p)x$$

This is inefficient! It requires different code for training and testing.

Standard Solution: Inverted Dropout

We scale the activations *during training* instead.

We divide the activations of the *kept* neurons by the keep probability.

$$\text{Train: } h = \frac{(\text{mask} \odot x)}{(1-p)}$$

Test: $h = x$

The test-time (inference) code is identical to a network *without* dropout. No scaling needed.

Slides from Understanding deep learning S. Prince Chapter 4,
Prof. O. Bastani and Prof. M. Zhao