

Lecture 2: Neural Networks



Today:



Today:

- Revisit feature transform (5%)



Today:

- Revisit feature transform (5%)
- What is a neural net? (10%)



Today:

- Revisit feature transform (5%)
- What is a neural net? (10%)
- Derivatives and chain-rule reminder (10%)



Today:

- Revisit feature transform (5%)
- What is a neural net? (10%)
- Derivatives and chain-rule reminder (10%)
- Training a vanilla network (back-prop) (40%)



Today:

- Revisit feature transform (5%)
- What is a neural net? (10%)
- Derivatives and chain-rule reminder (10%)
- Training a vanilla network (back-prop) (40%)
- Differential computational graph (25%)

Today:

- Revisit feature transform (5%)
- What is a neural net? (10%)
- Derivatives and chain-rule reminder (10%)
- Training a vanilla network (back-prop) (40%)
- Differential computational graph (25%)
- Demo (10%)

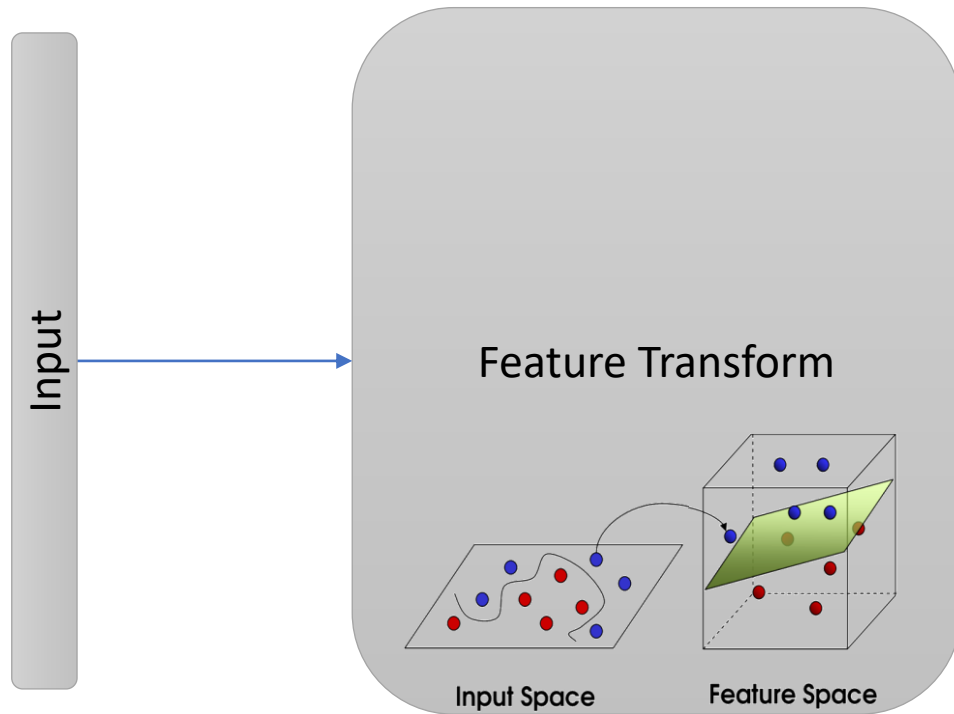


Feature transform

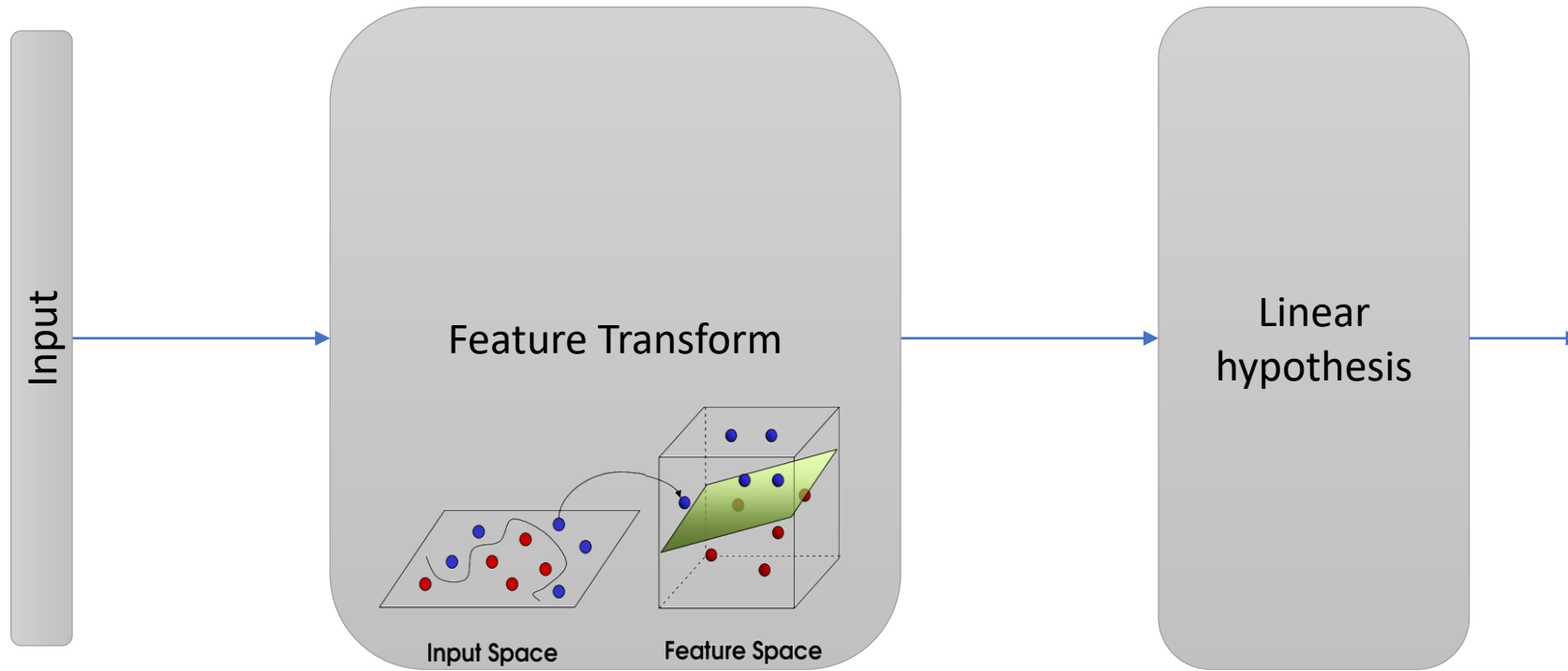
Input



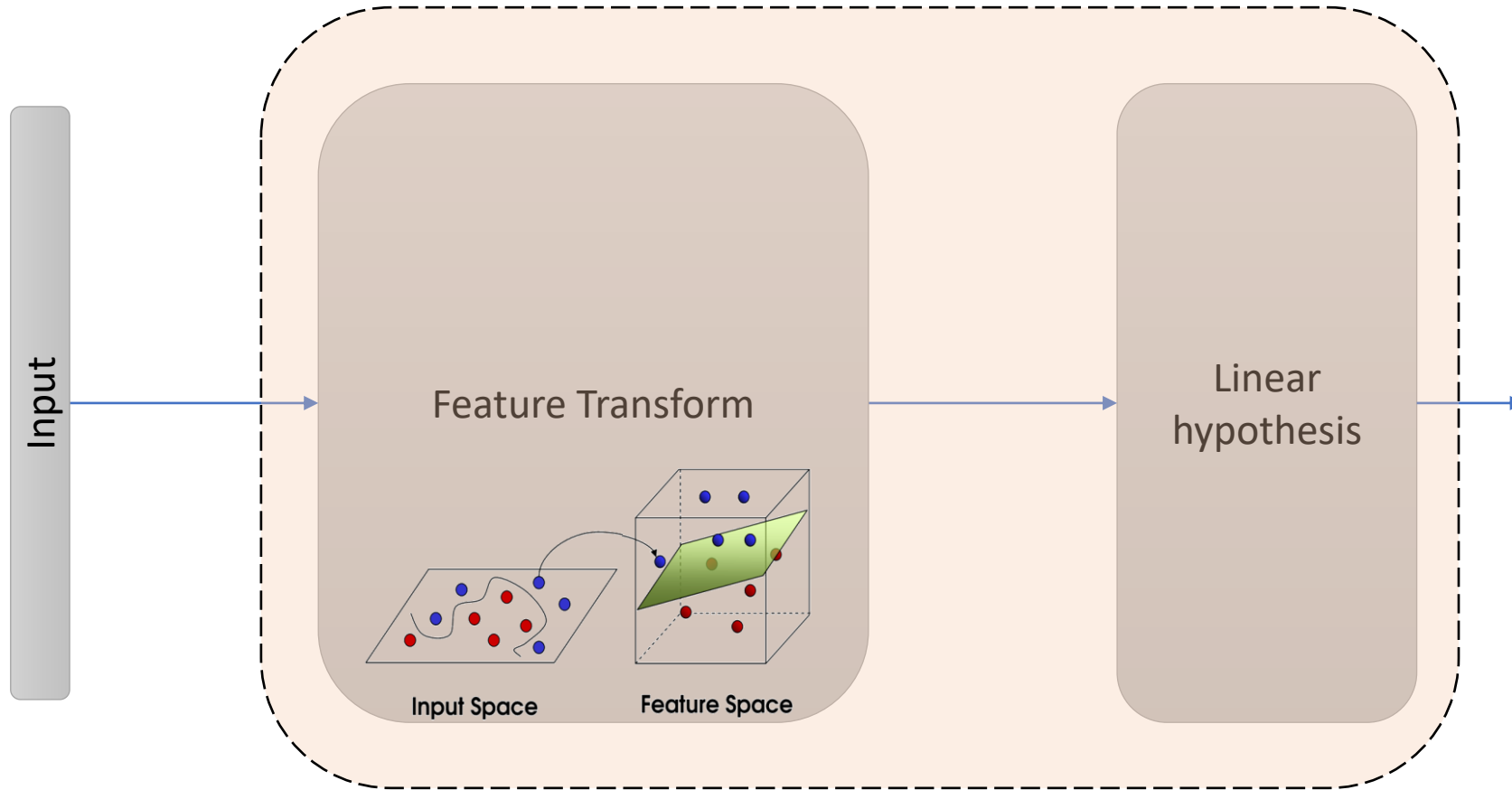
Feature transform



Feature transform

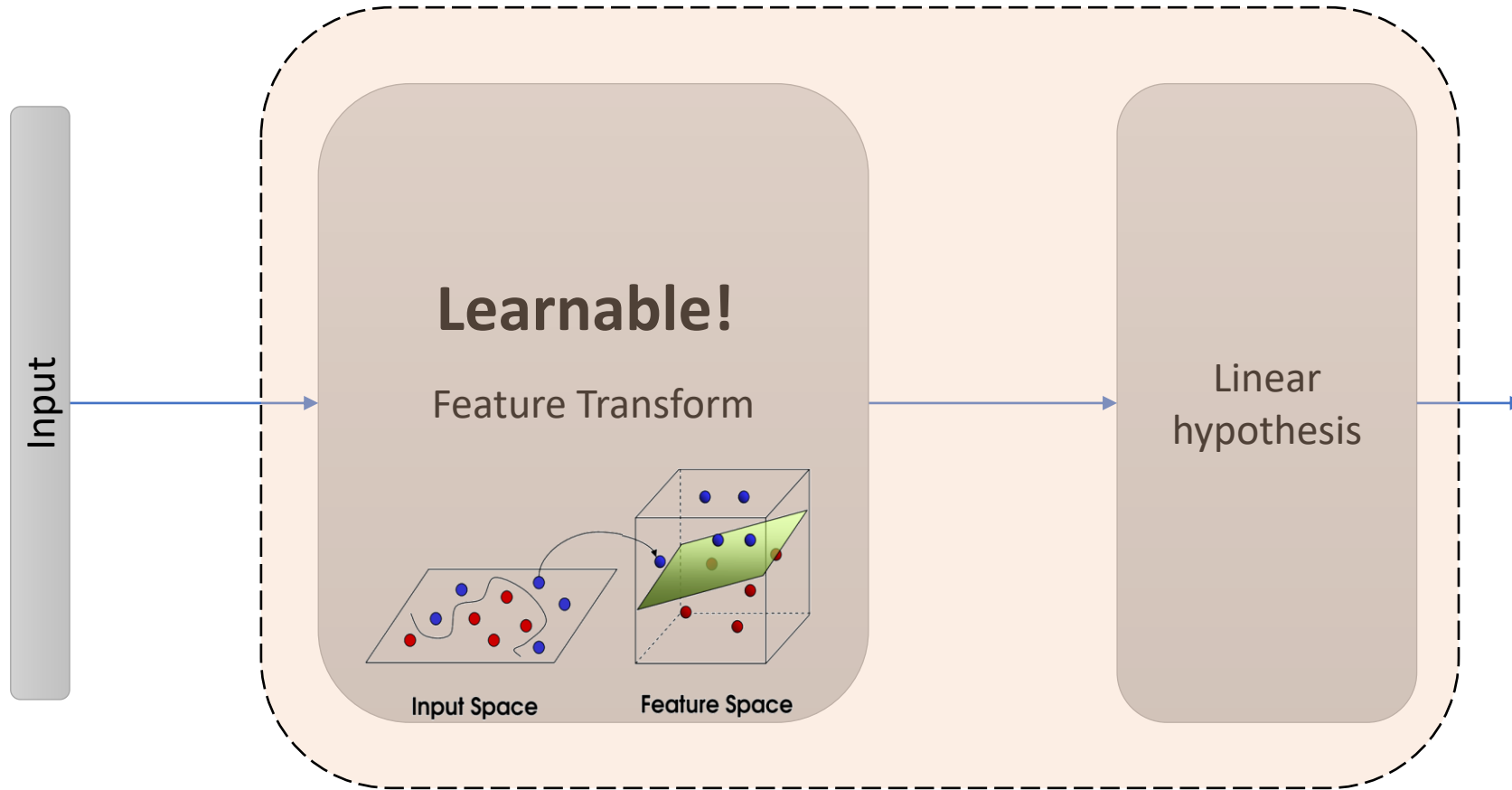


Feature transform



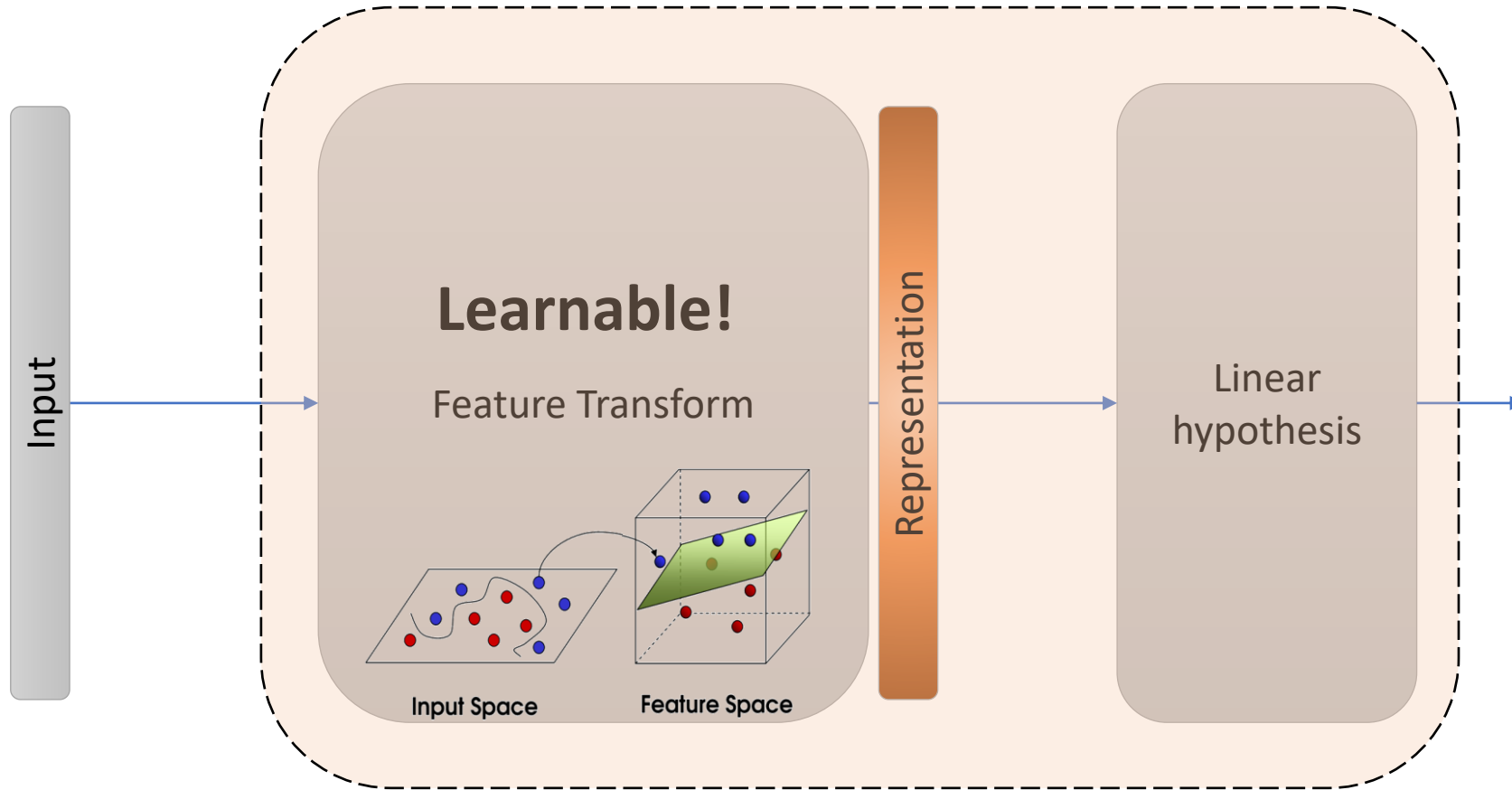
Non-linear hypothesis!

Feature transform



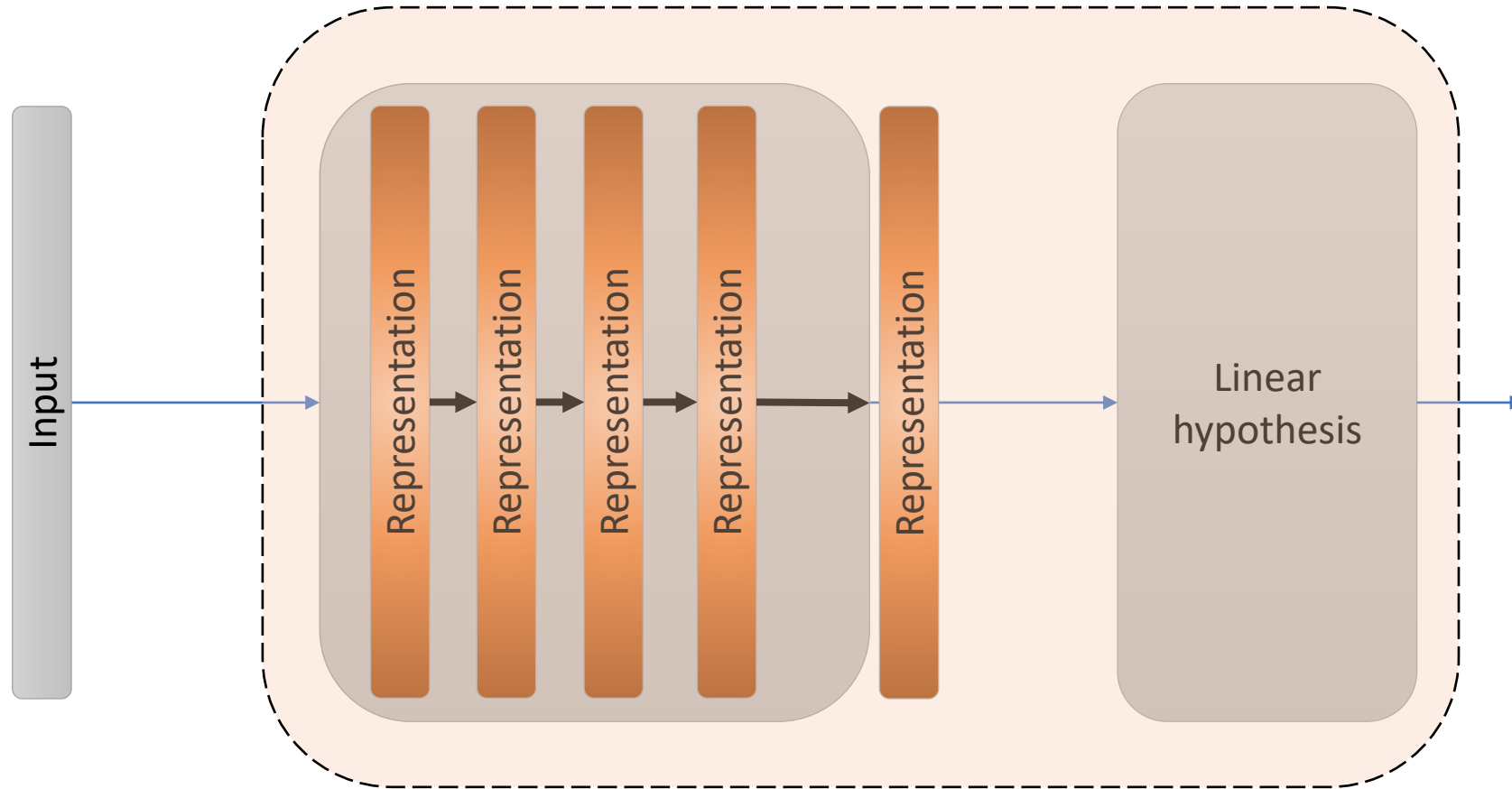
Non-linear hypothesis!

Feature transform



Non-linear hypothesis!

Feature transform



Non-linear hypothesis!

Artificial Neural Networks



Artificial Neural Networks

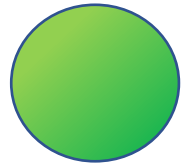
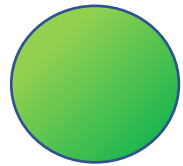
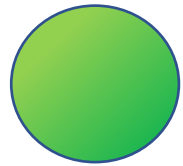
Vaguely inspired by biological neural networks



Artificial Neural Networks

Vaguely inspired by biological neural networks

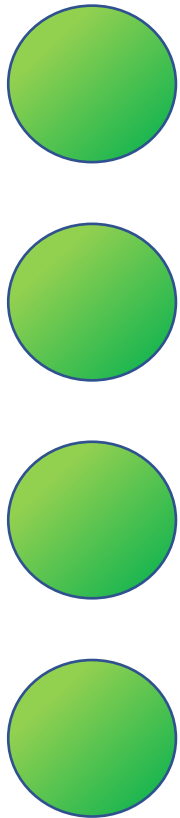
Input vector



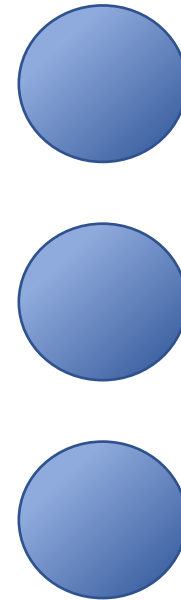
Artificial Neural Networks

Vaguely inspired by biological neural networks

Input vector



Output vector



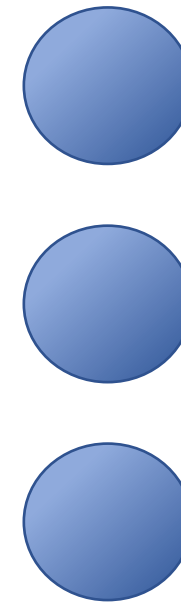
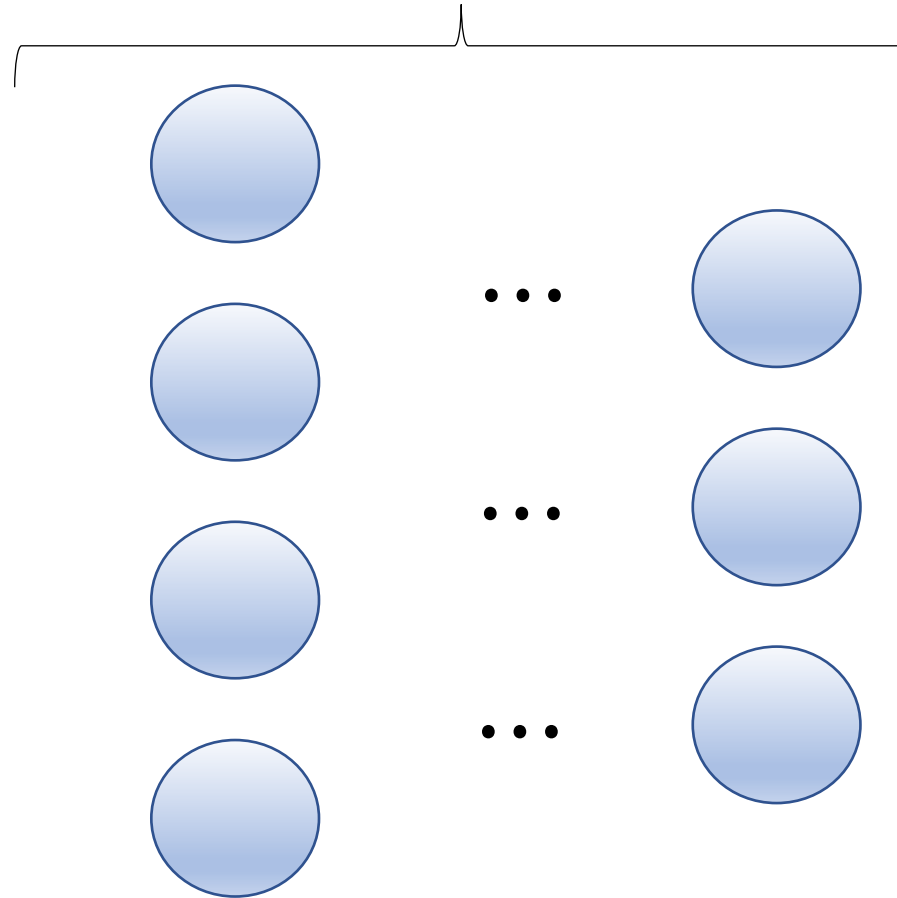
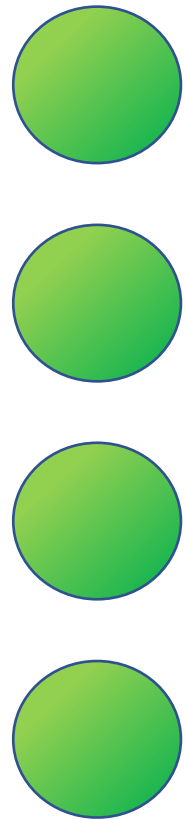
Artificial Neural Networks

Vaguely inspired by biological neural networks

Input vector

Hidden layers

Output vector



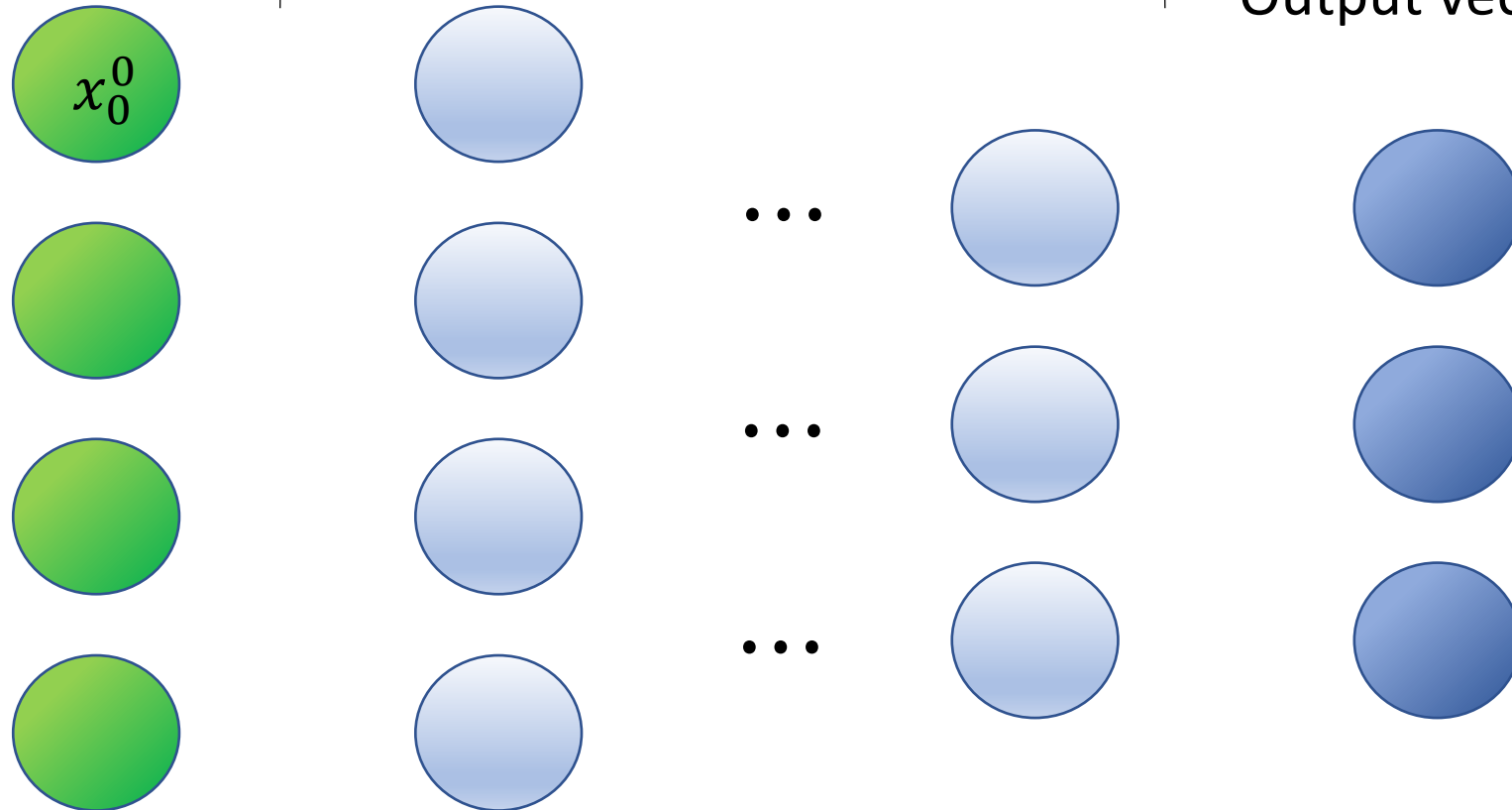
Artificial Neural Networks

Vaguely inspired by biological neural networks

Input vector

Hidden layers

Output vector



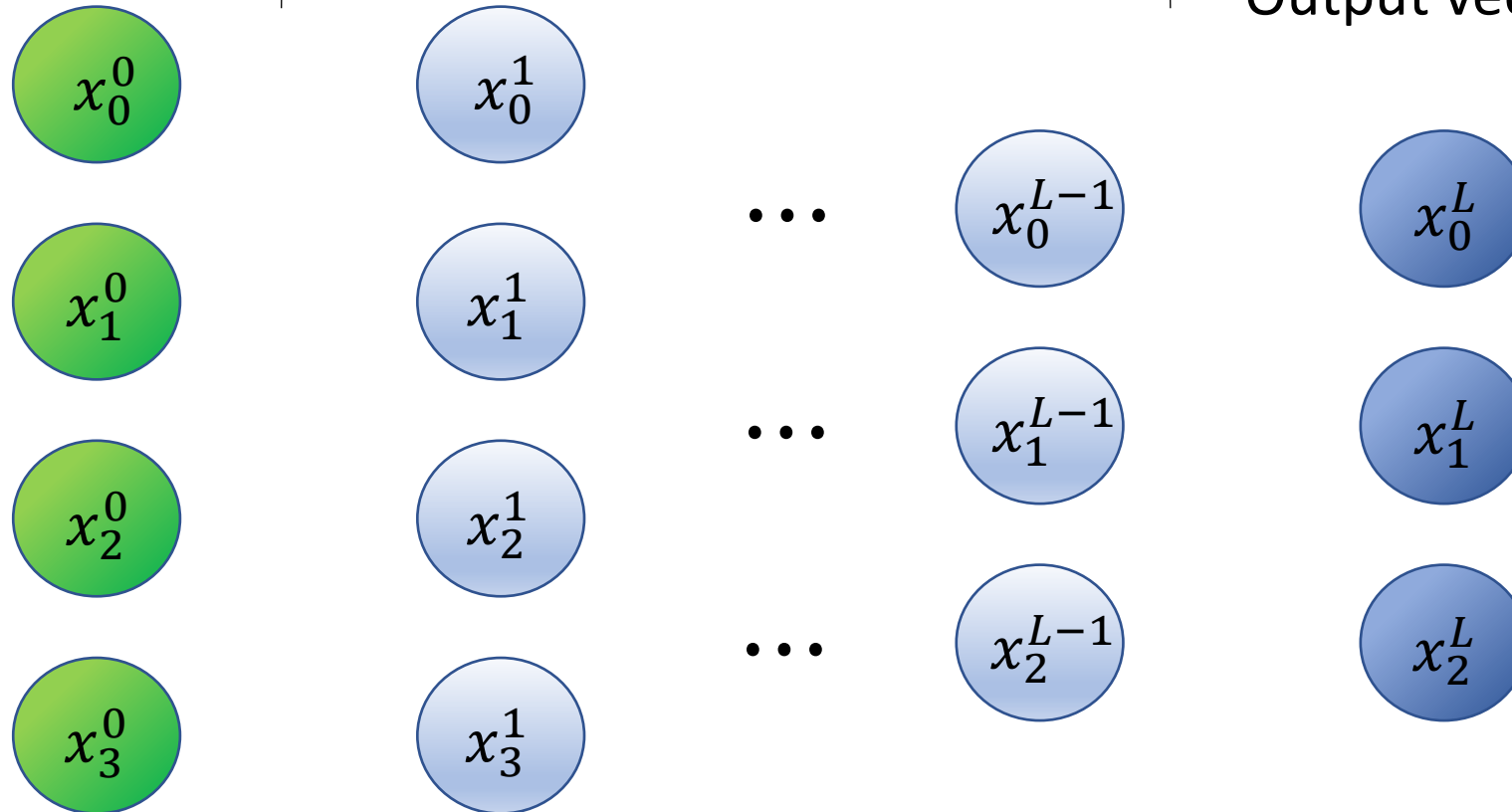
Artificial Neural Networks

Vaguely inspired by biological neural networks

Input vector

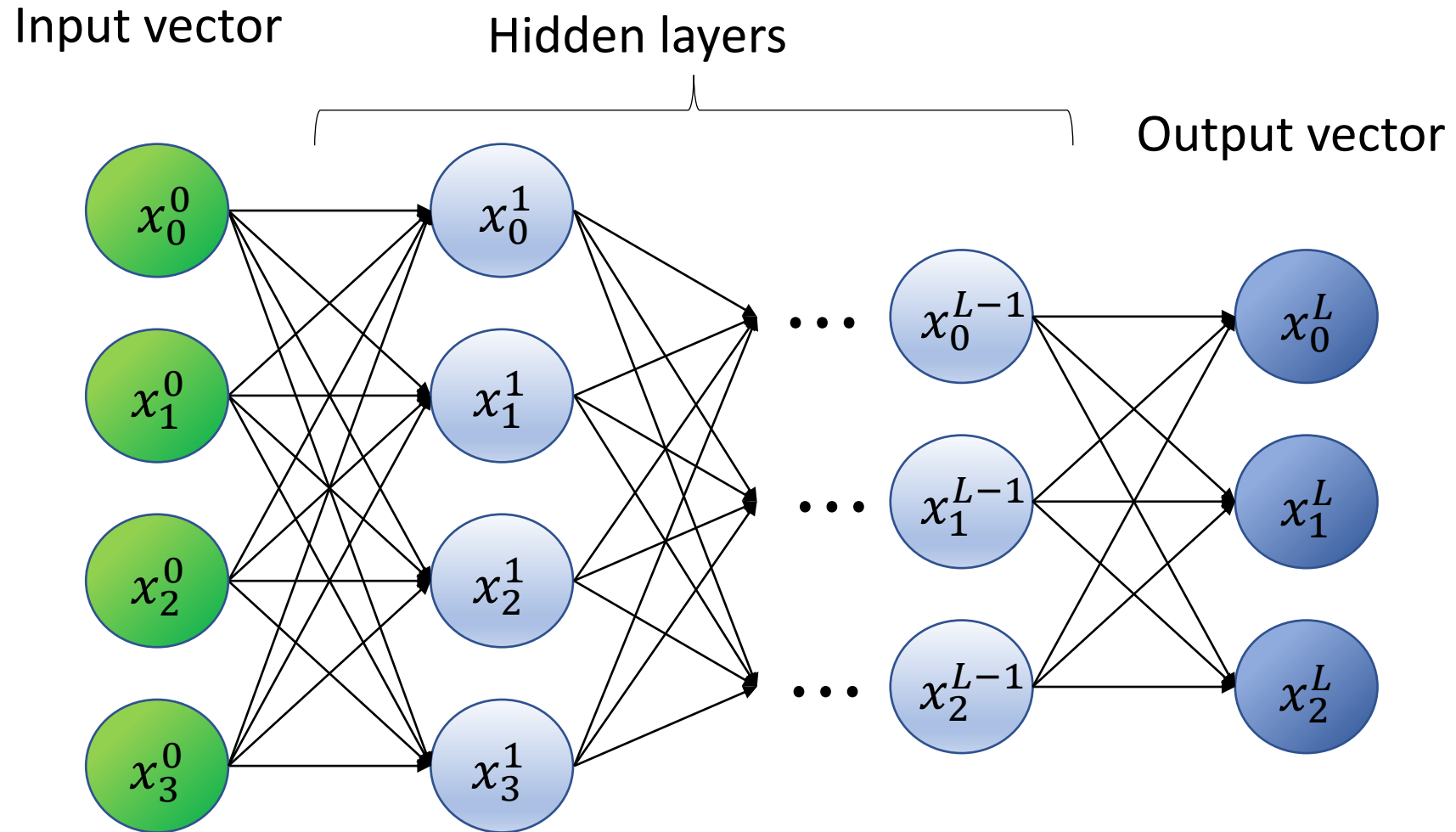
Hidden layers

Output vector



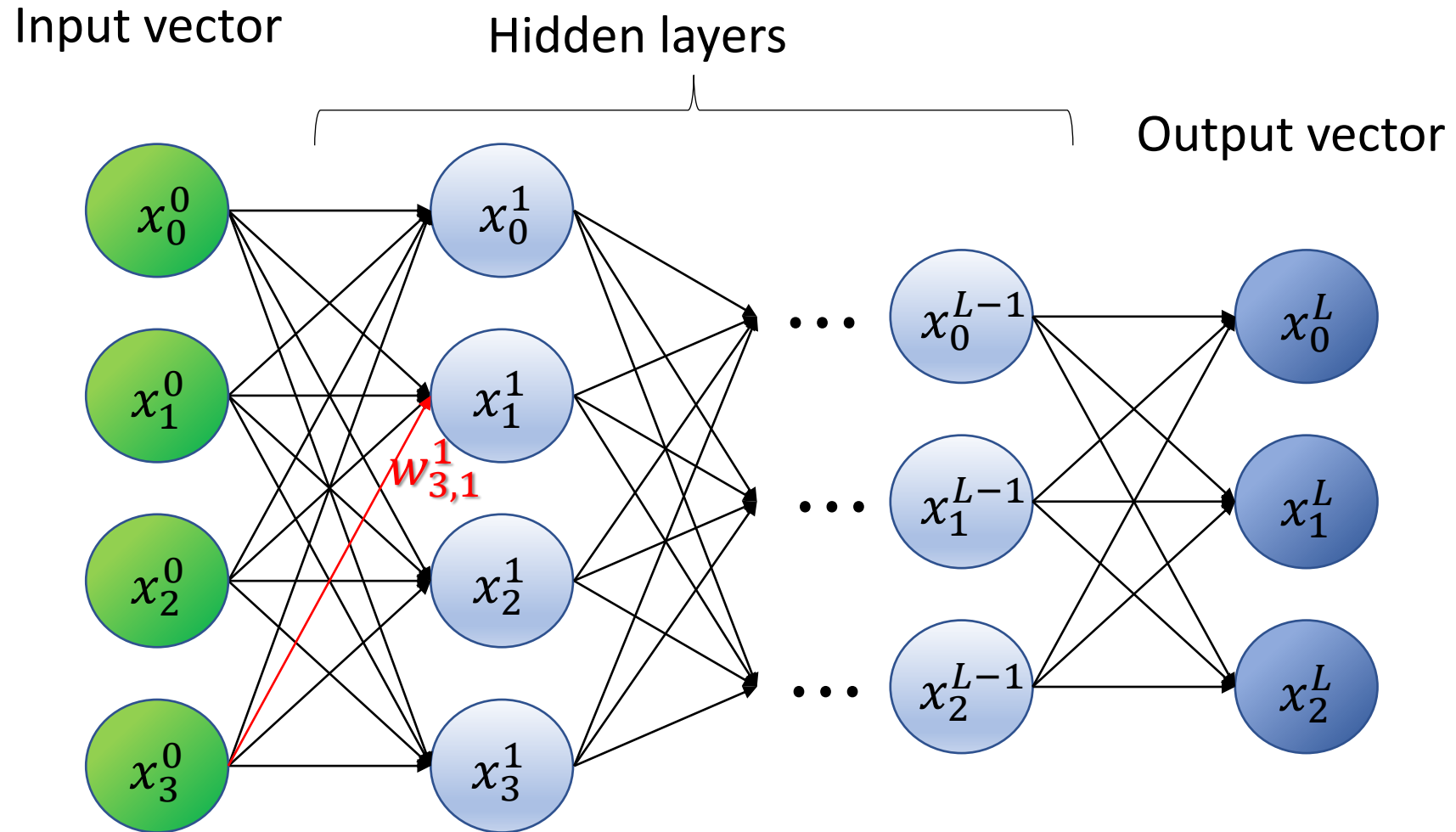
Artificial Neural Networks

Vaguely inspired by biological neural networks



Artificial Neural Networks

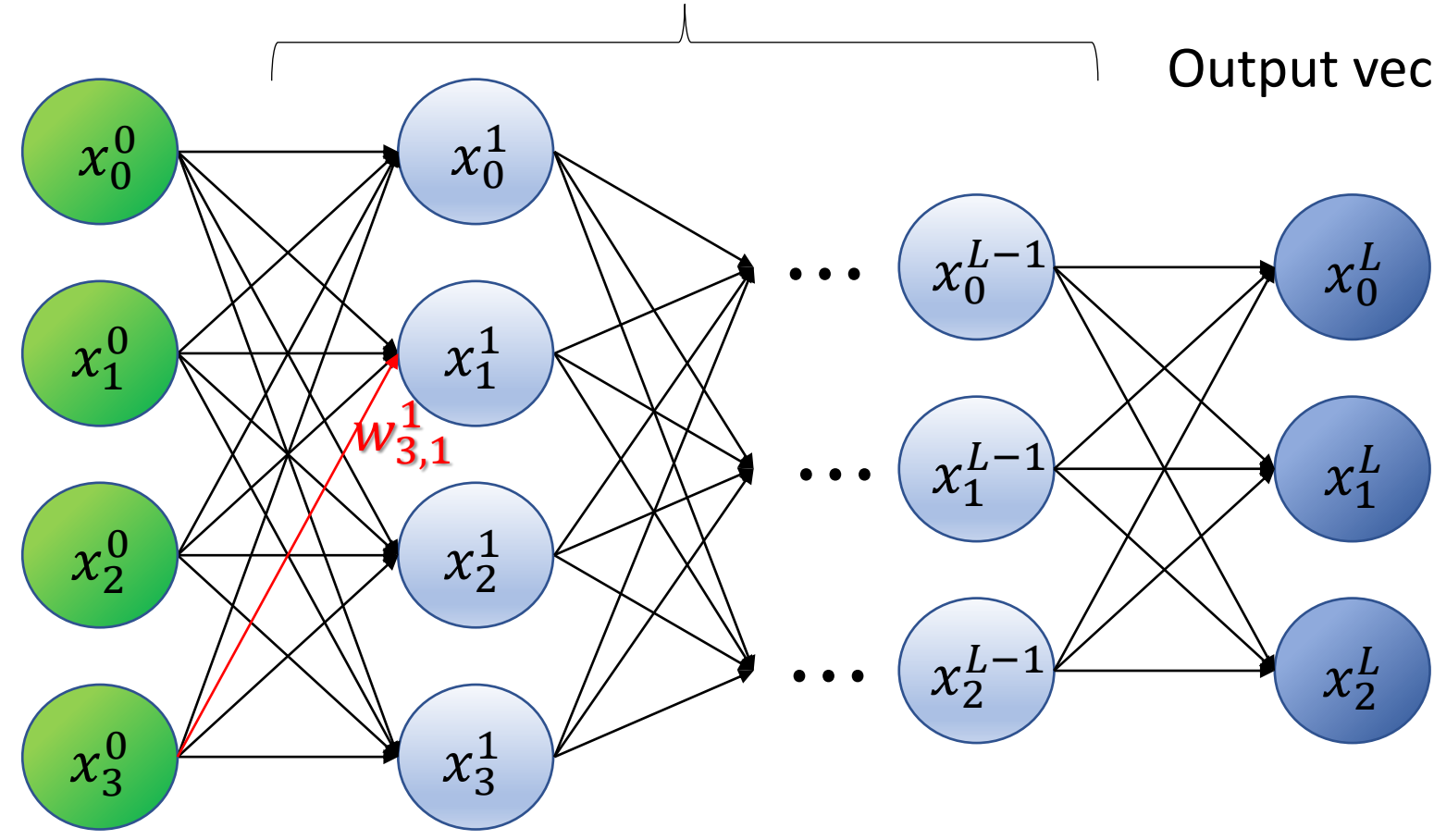
Vaguely inspired by biological neural networks

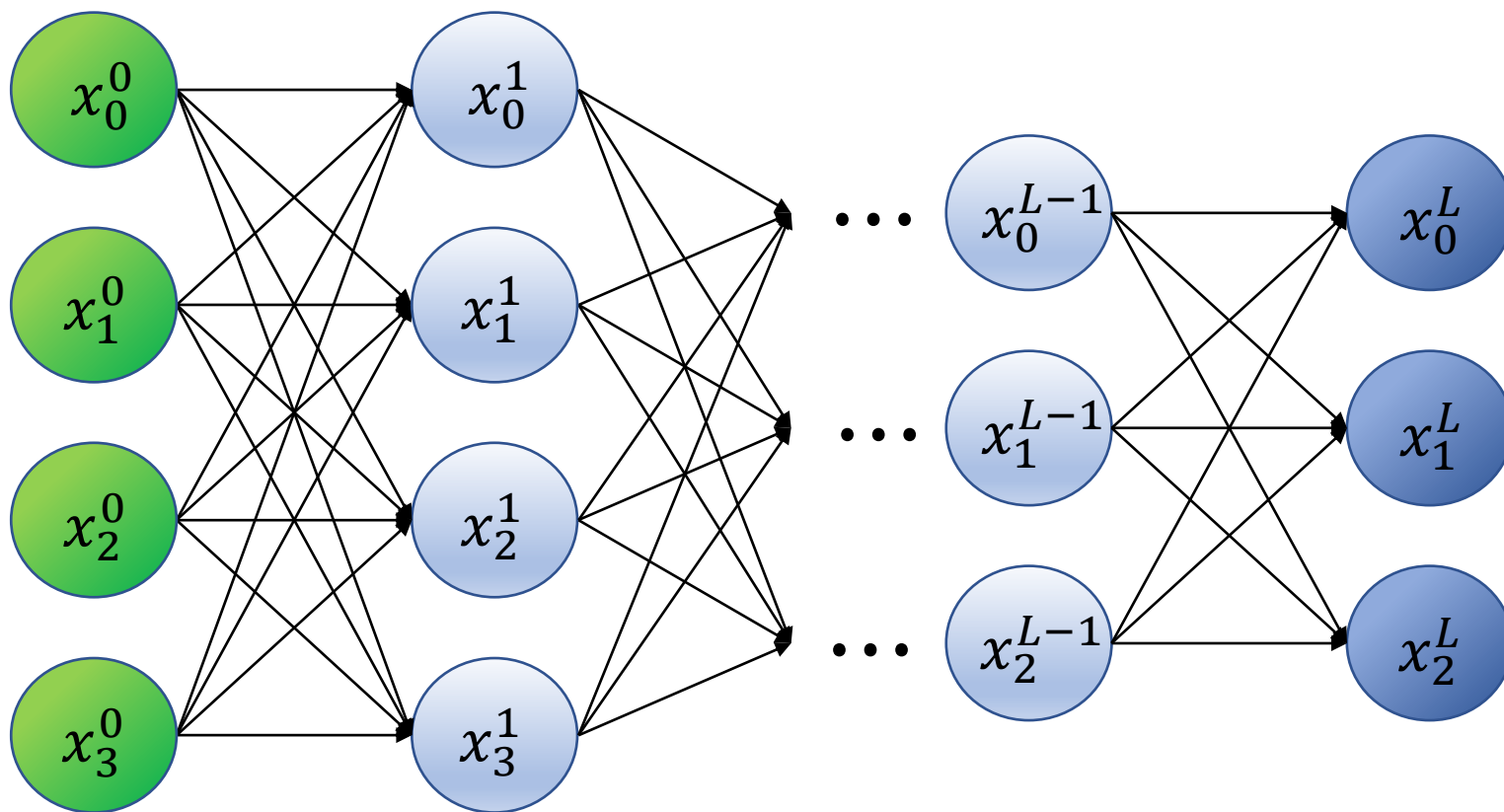


Input vector

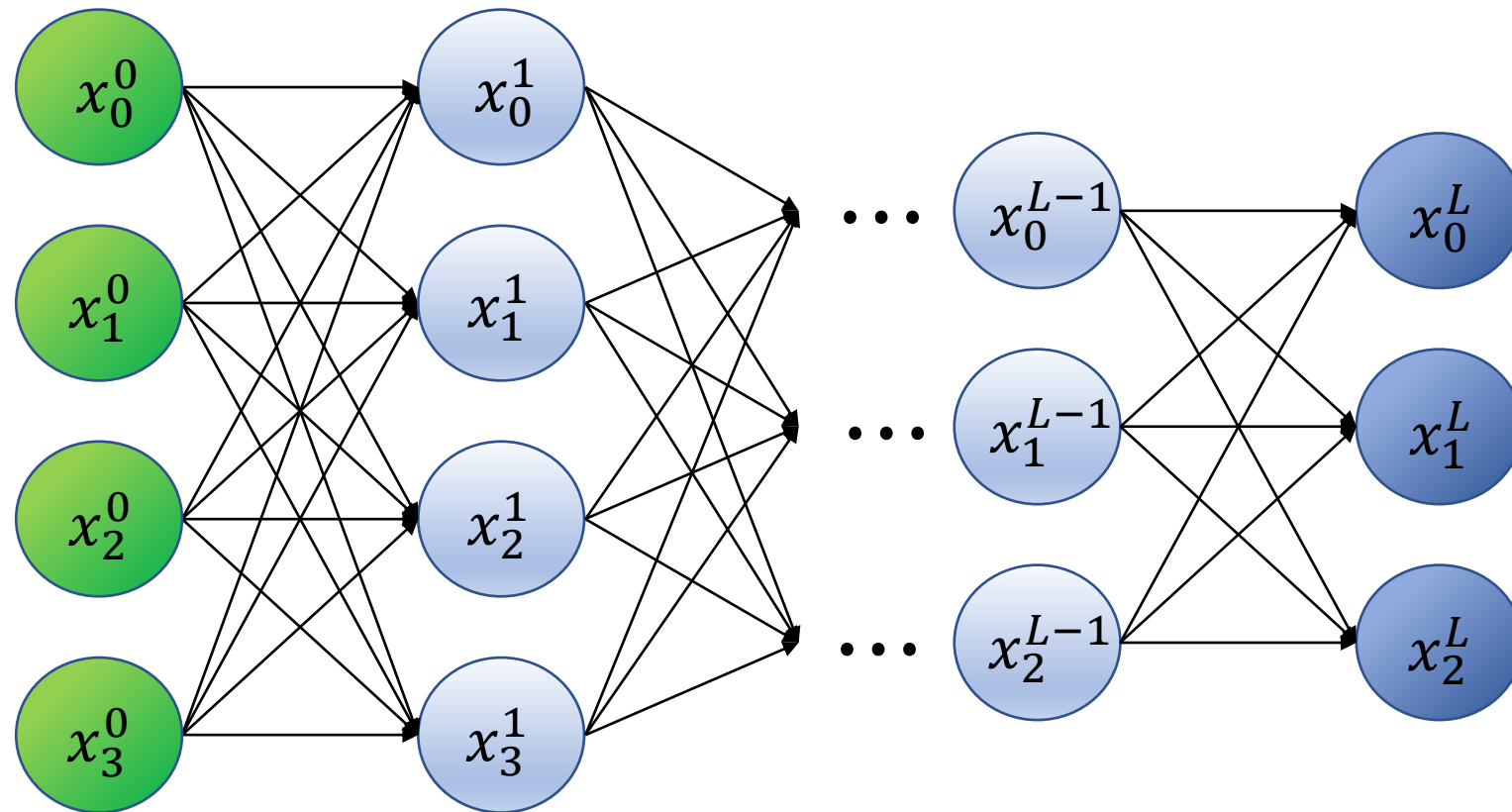
Hidden layers

Output vector

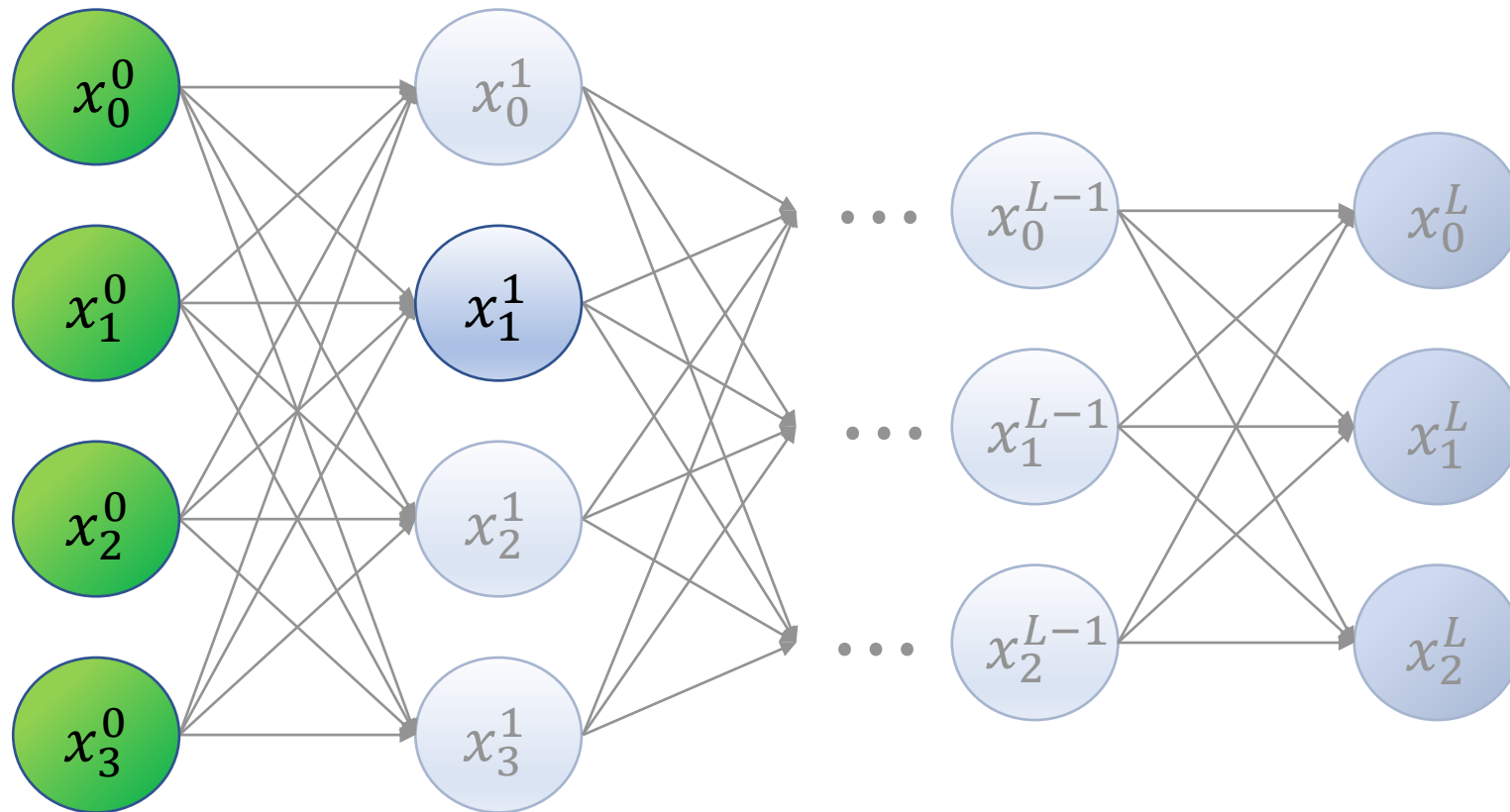




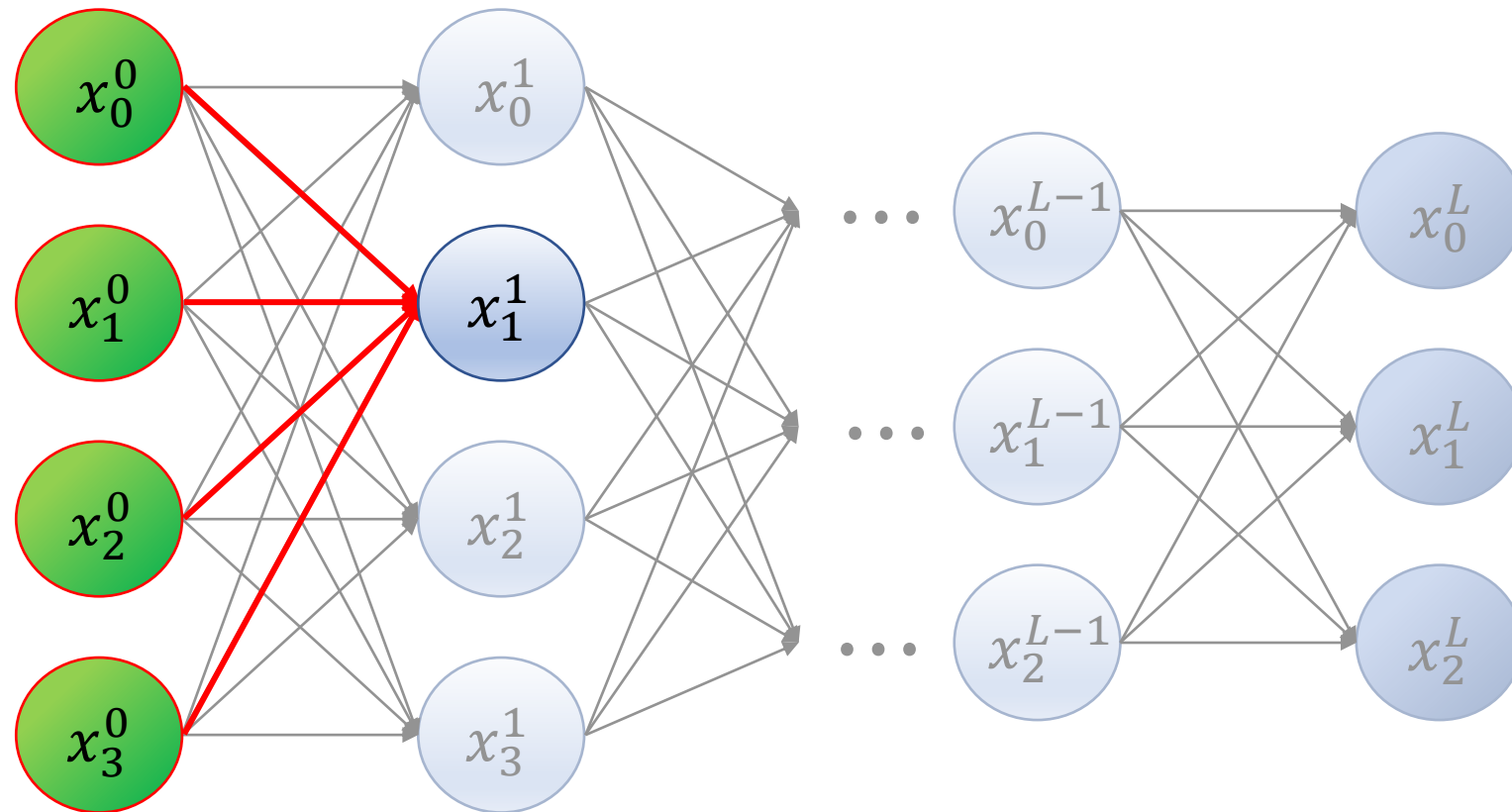
$$x_j^l =$$



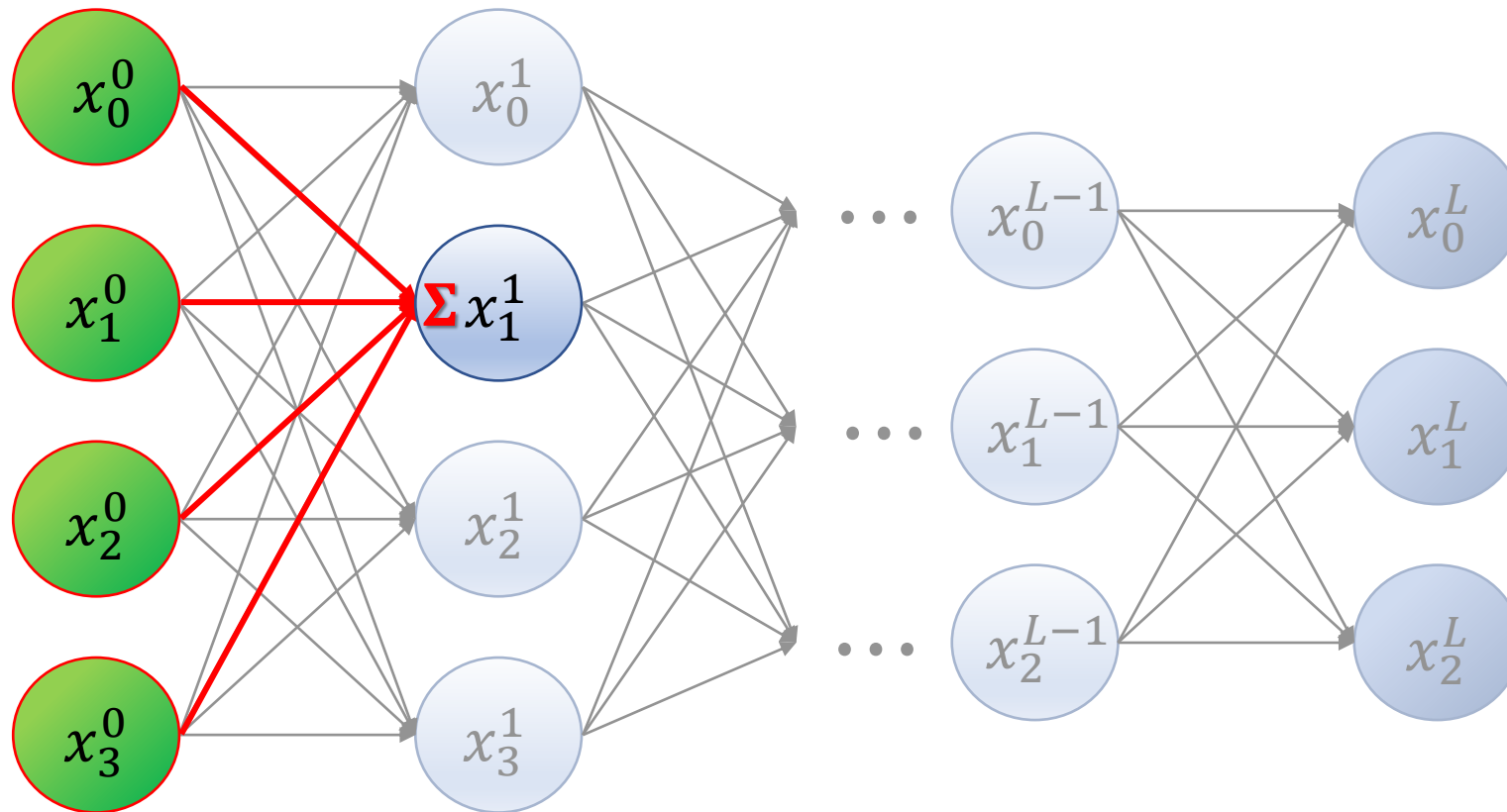
$$x_j^l =$$



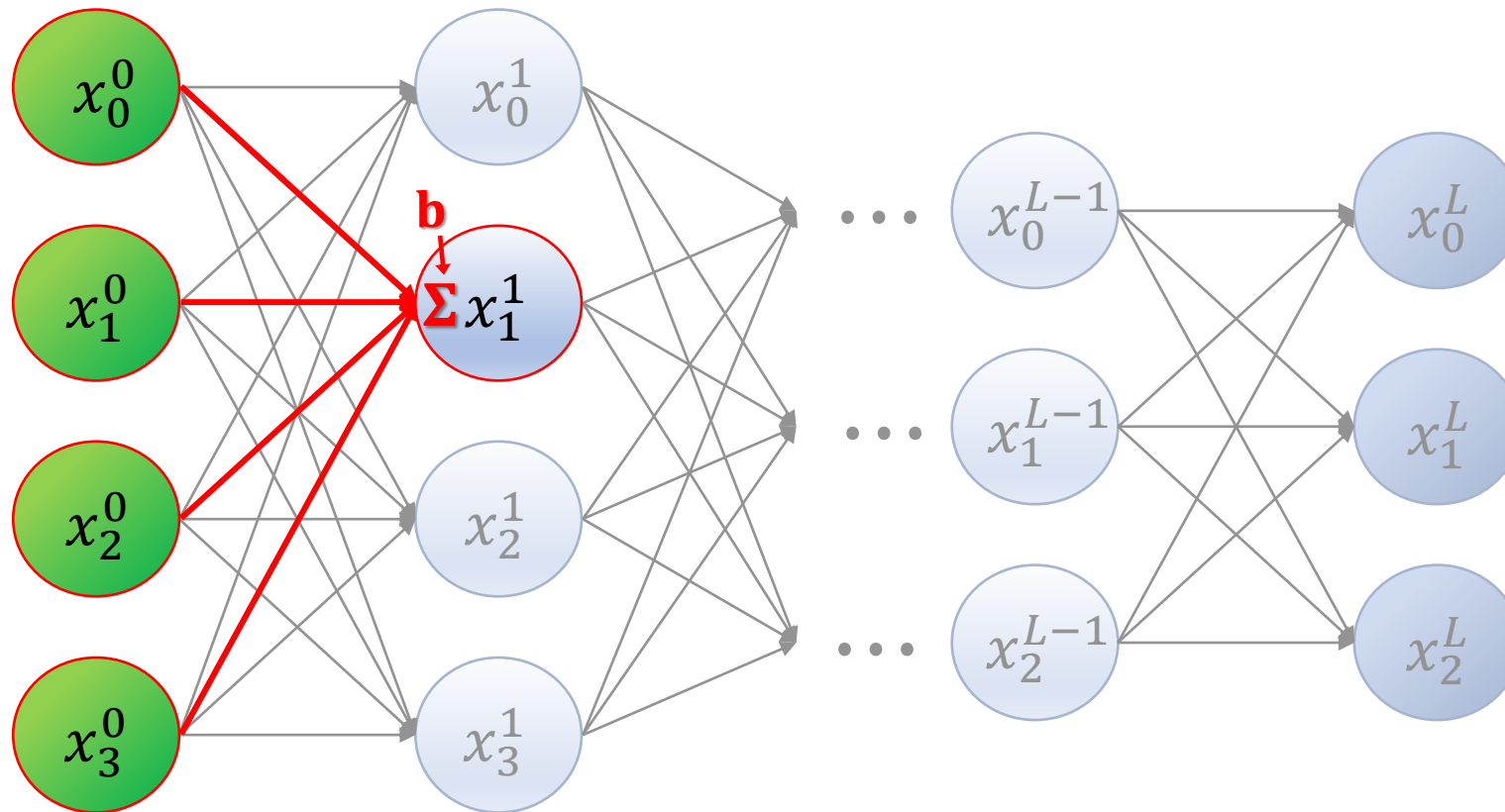
$$x_j^l =$$



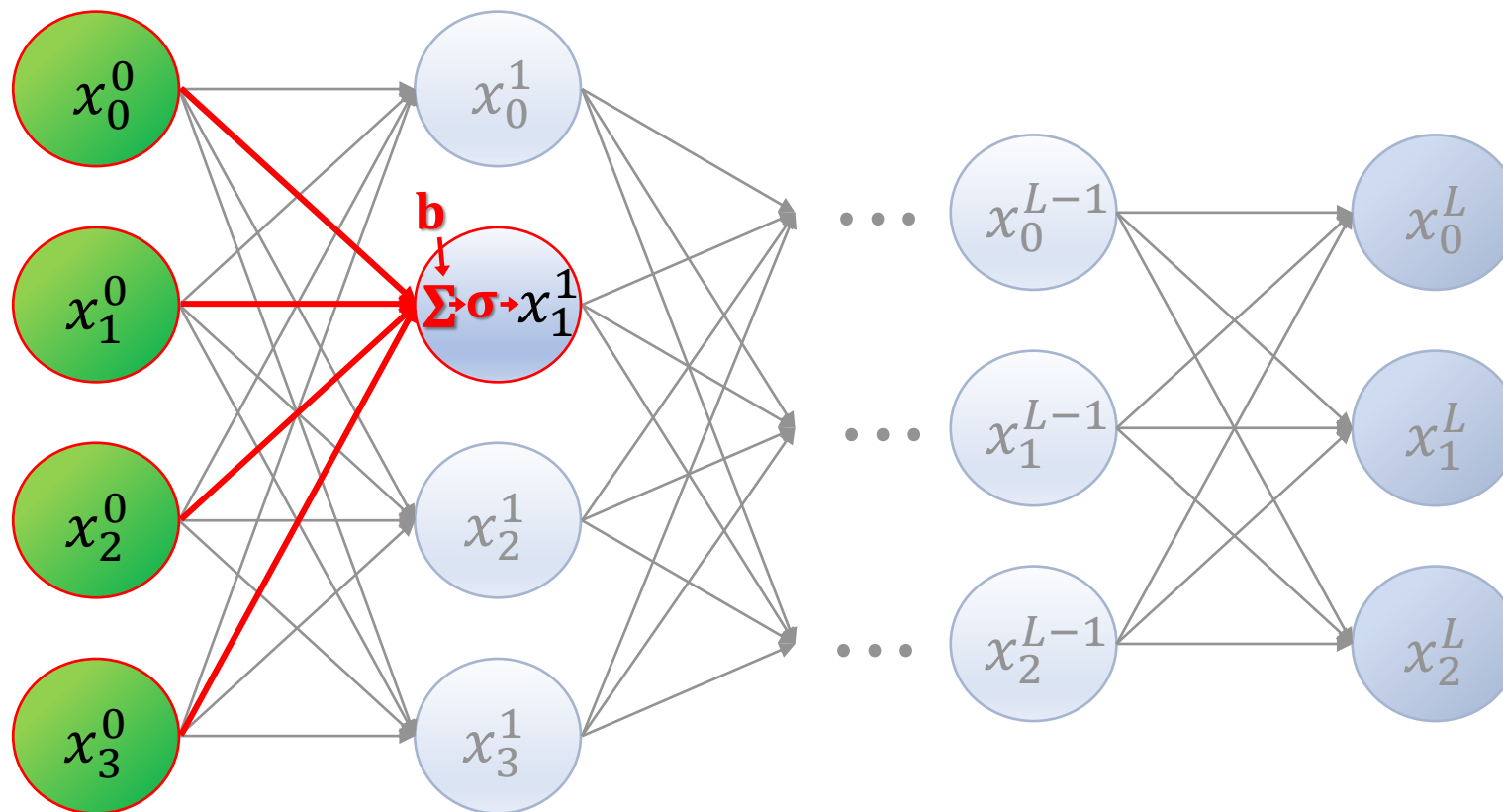
$$x_j^l = \sum_i w_{ij}^l \cdot x_i^{l-1}$$



$$x_j^l = \sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l$$

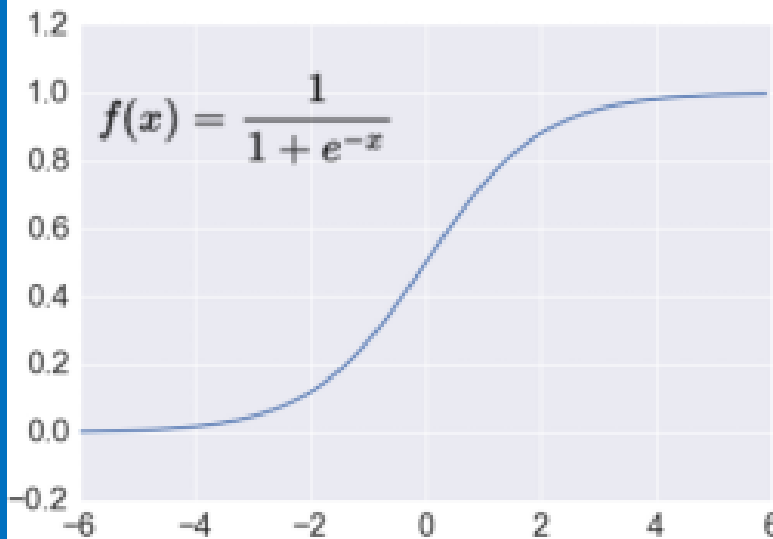


$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$

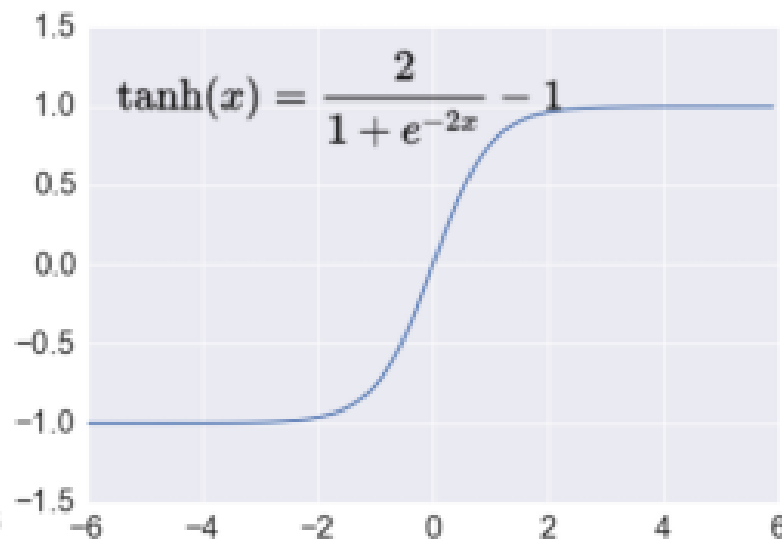


$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$

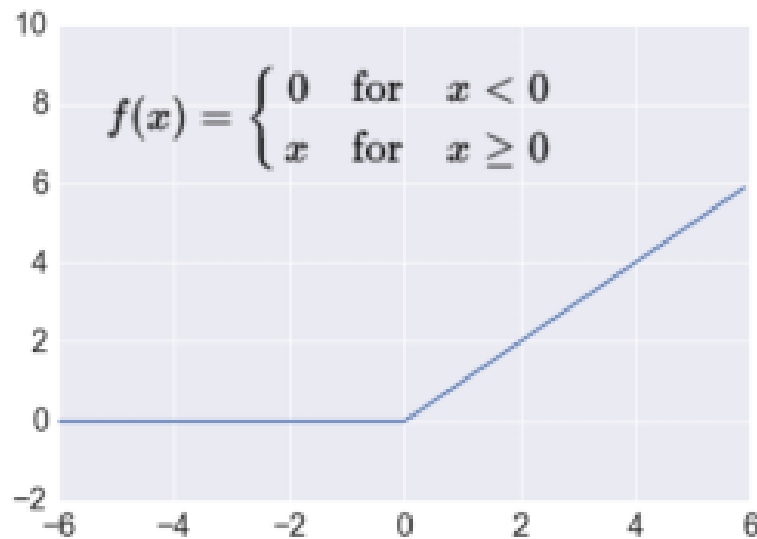
Sigmoid



TanH



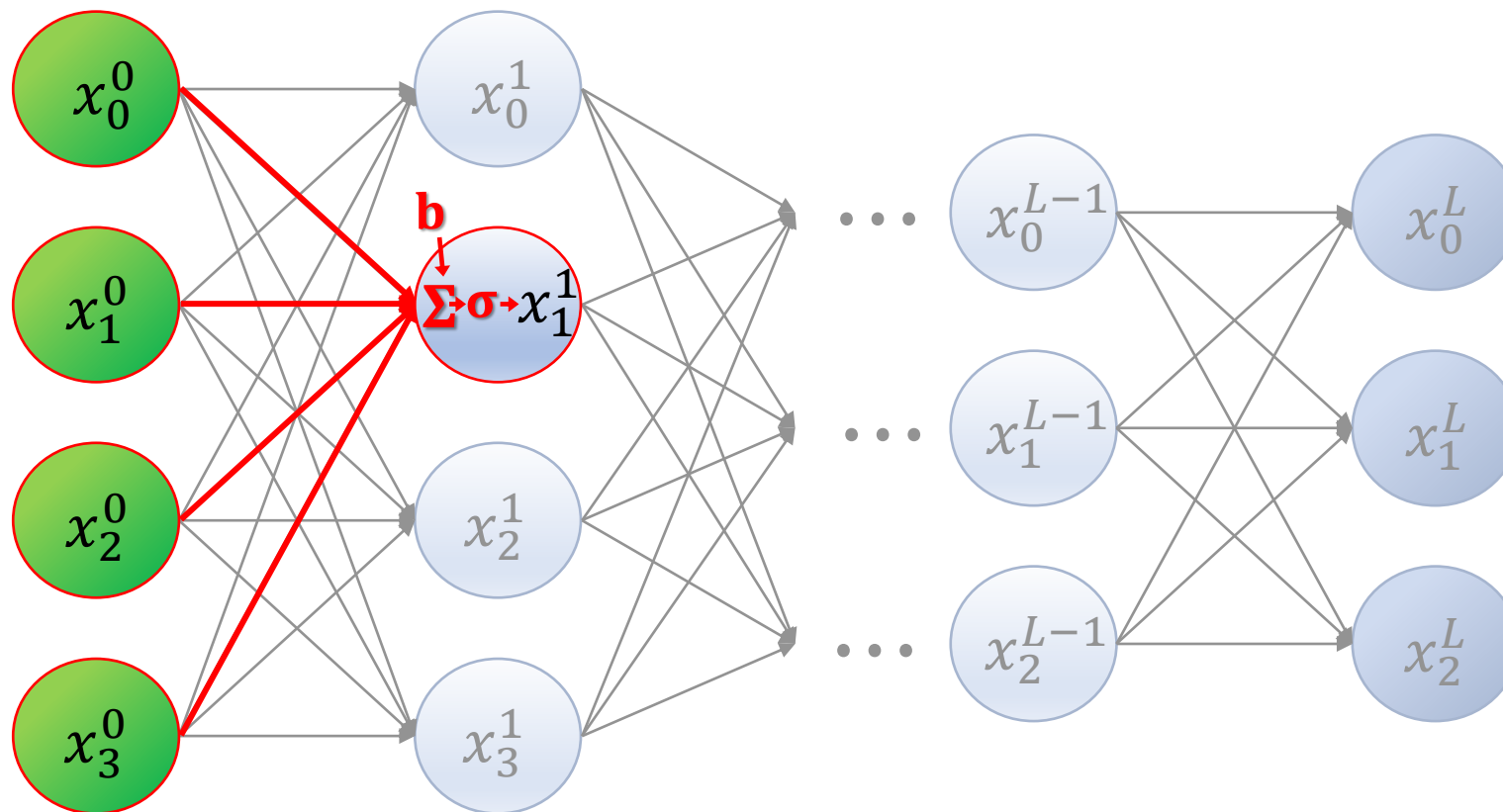
ReLU



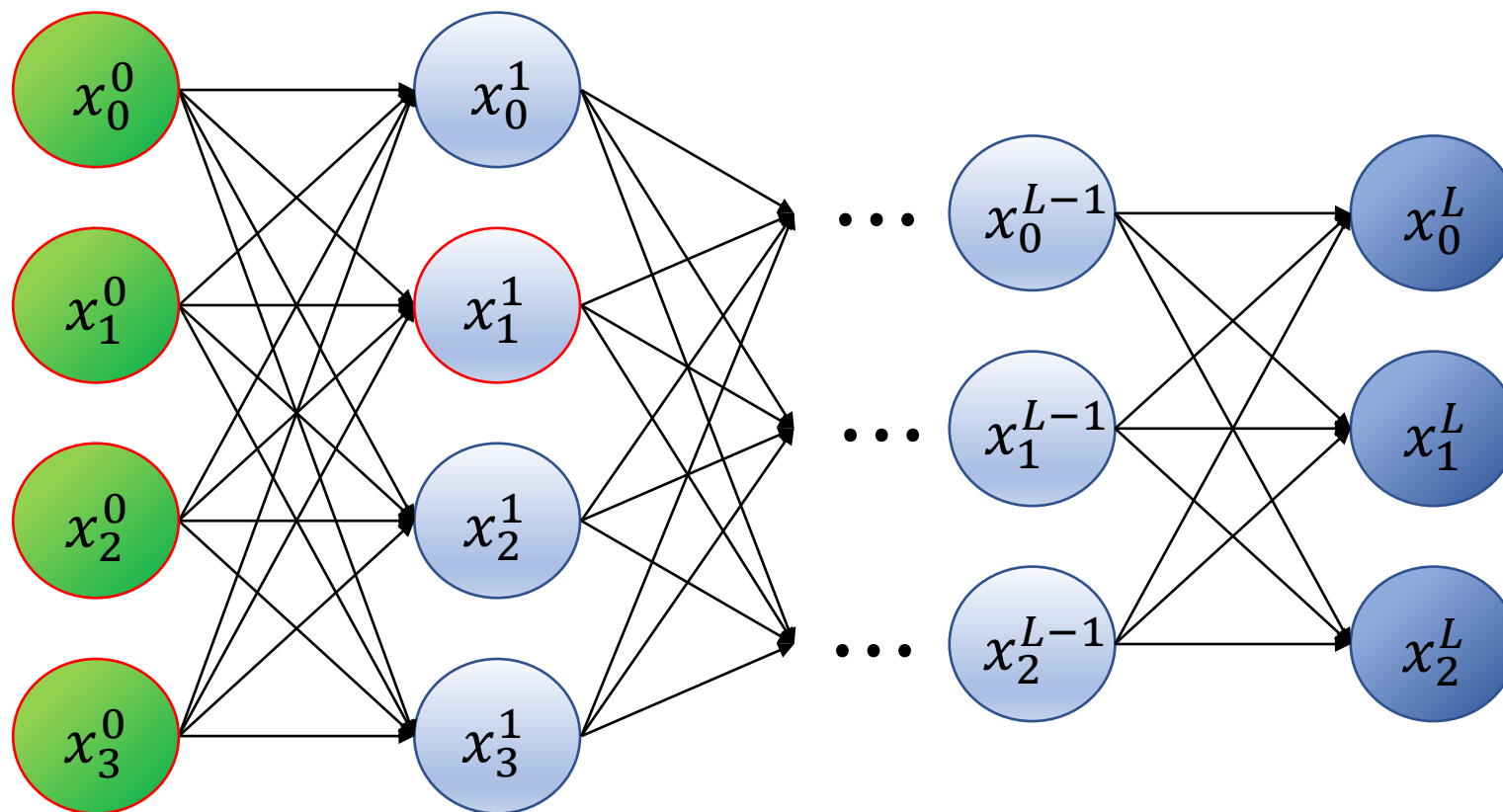
λ_3

λ_3

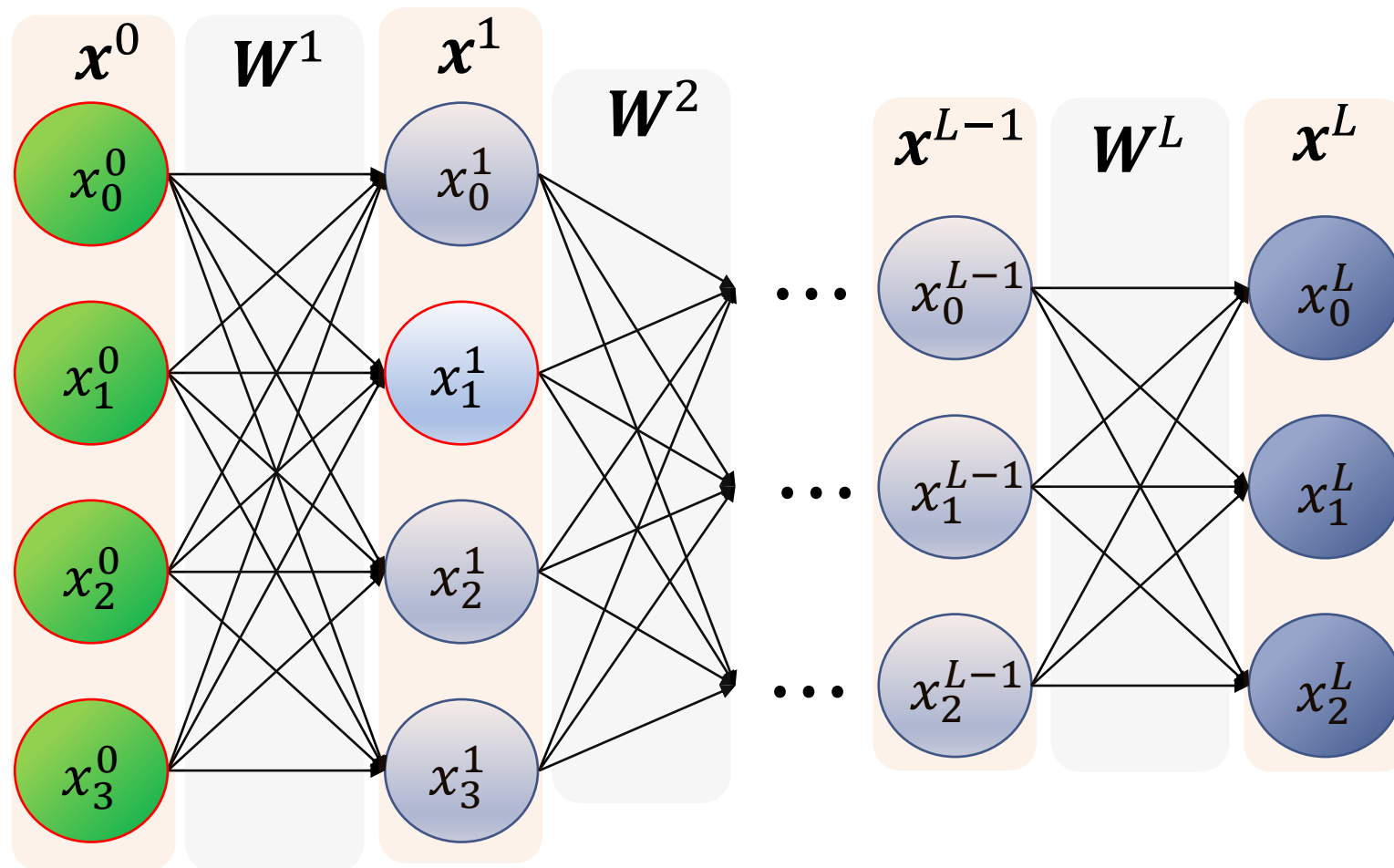
$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$



$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$

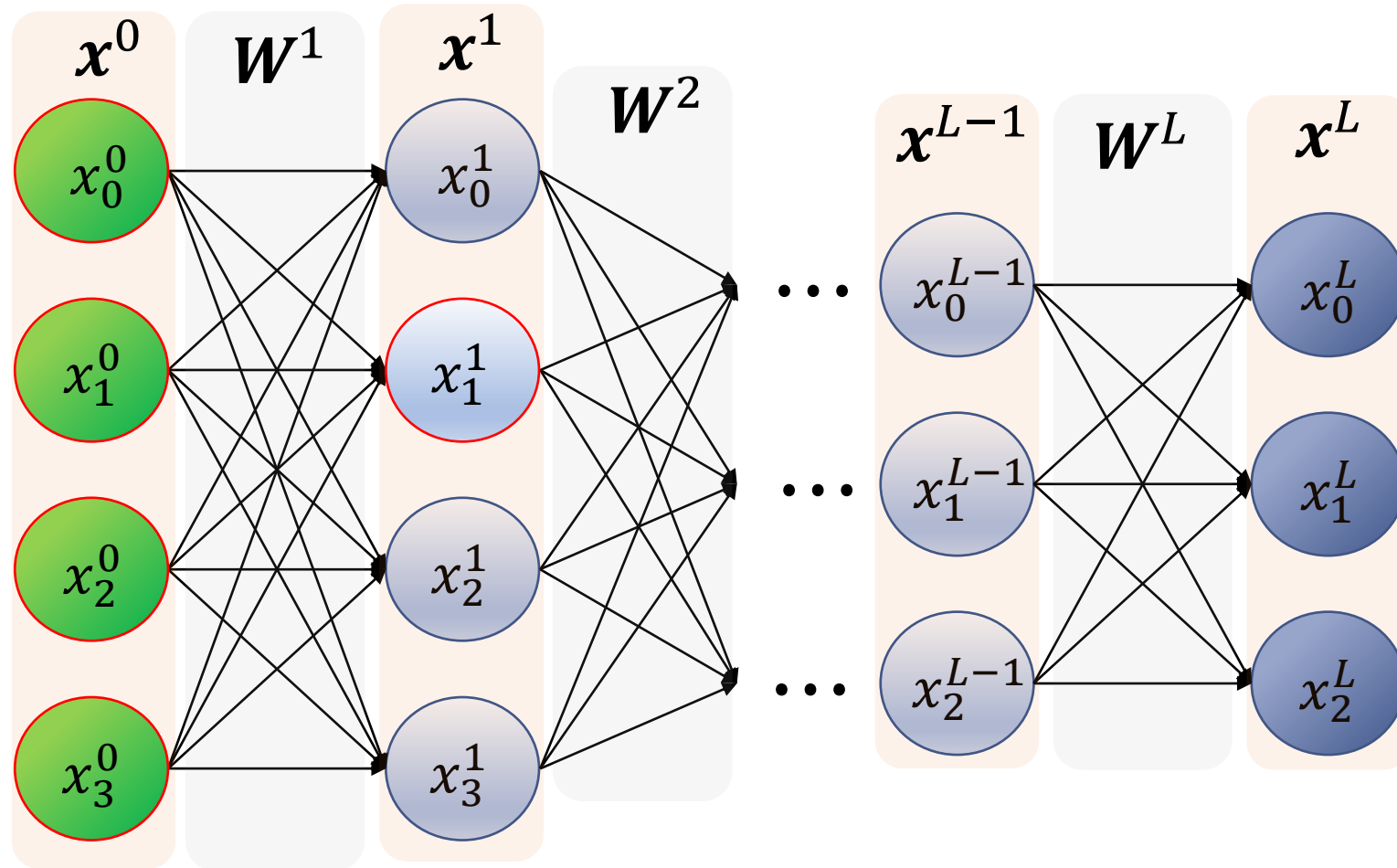


$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$



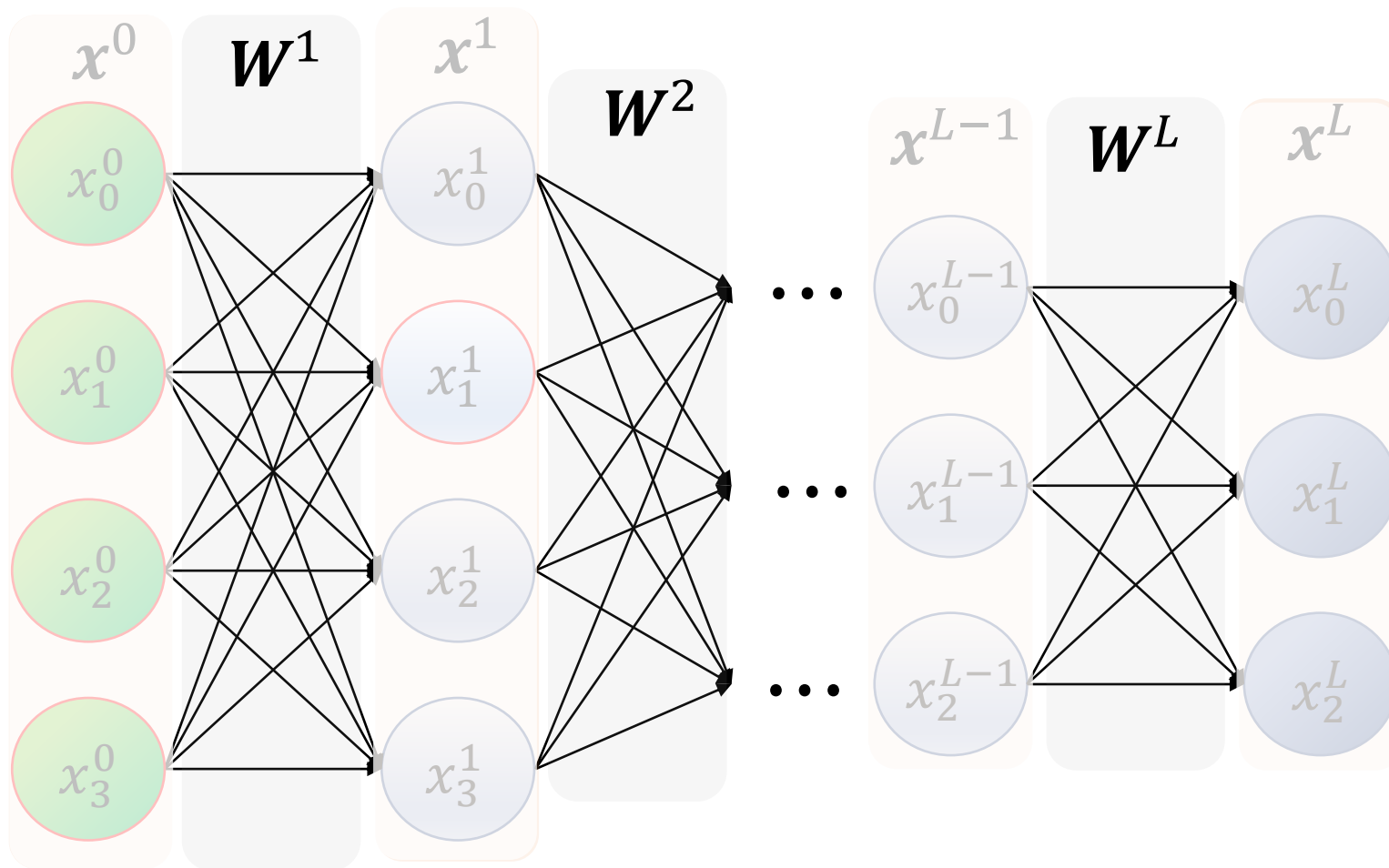
$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$

$$\mathbf{x}^l = \sigma \left(\mathbf{W}^{lT} \mathbf{x}^{l-1} + \mathbf{b}^l \right)$$



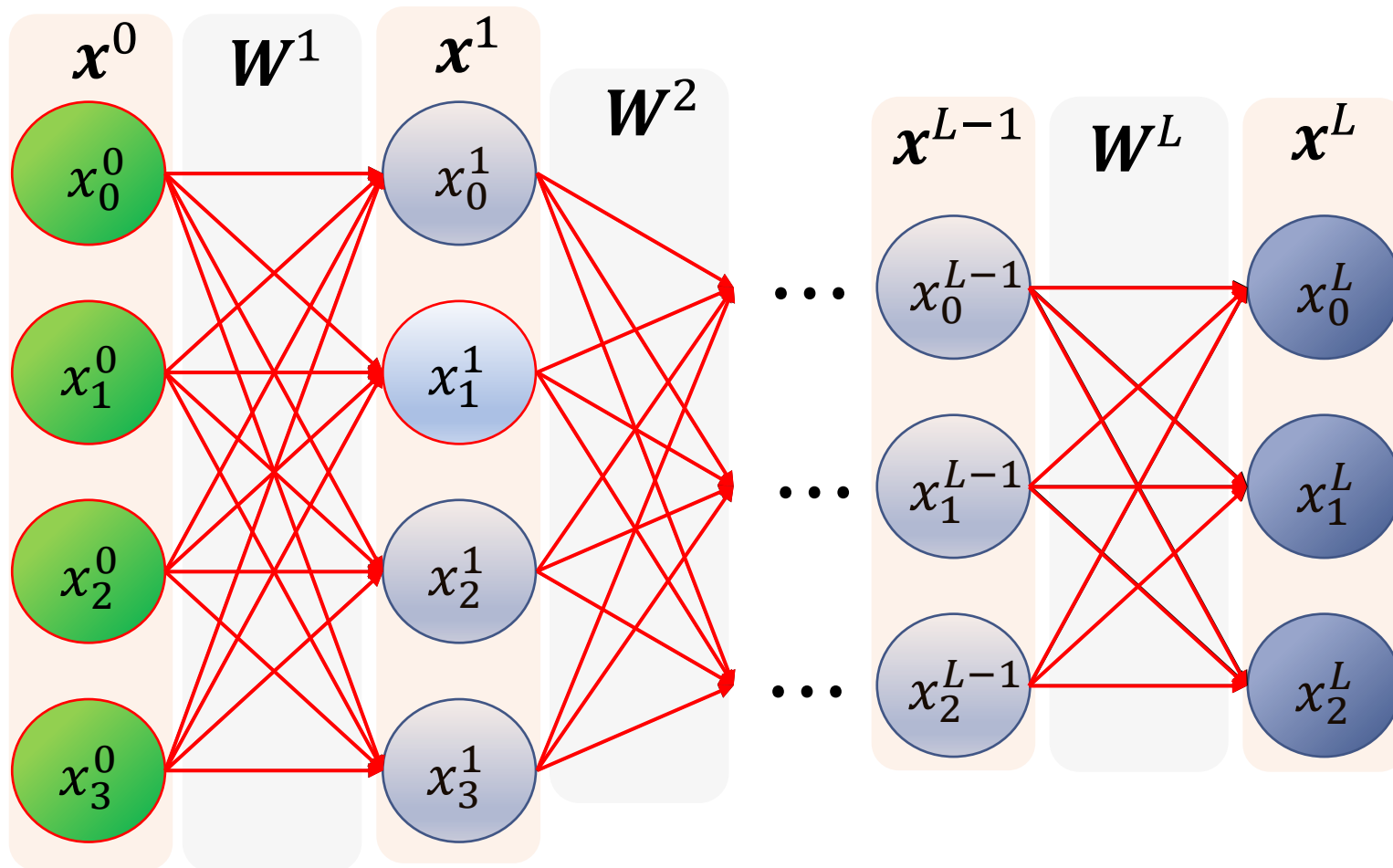
$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$

$$\mathbf{x}^l = \sigma \left(\mathbf{W}^{lT} \mathbf{x}^{l-1} + \mathbf{b}^l \right)$$



$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$

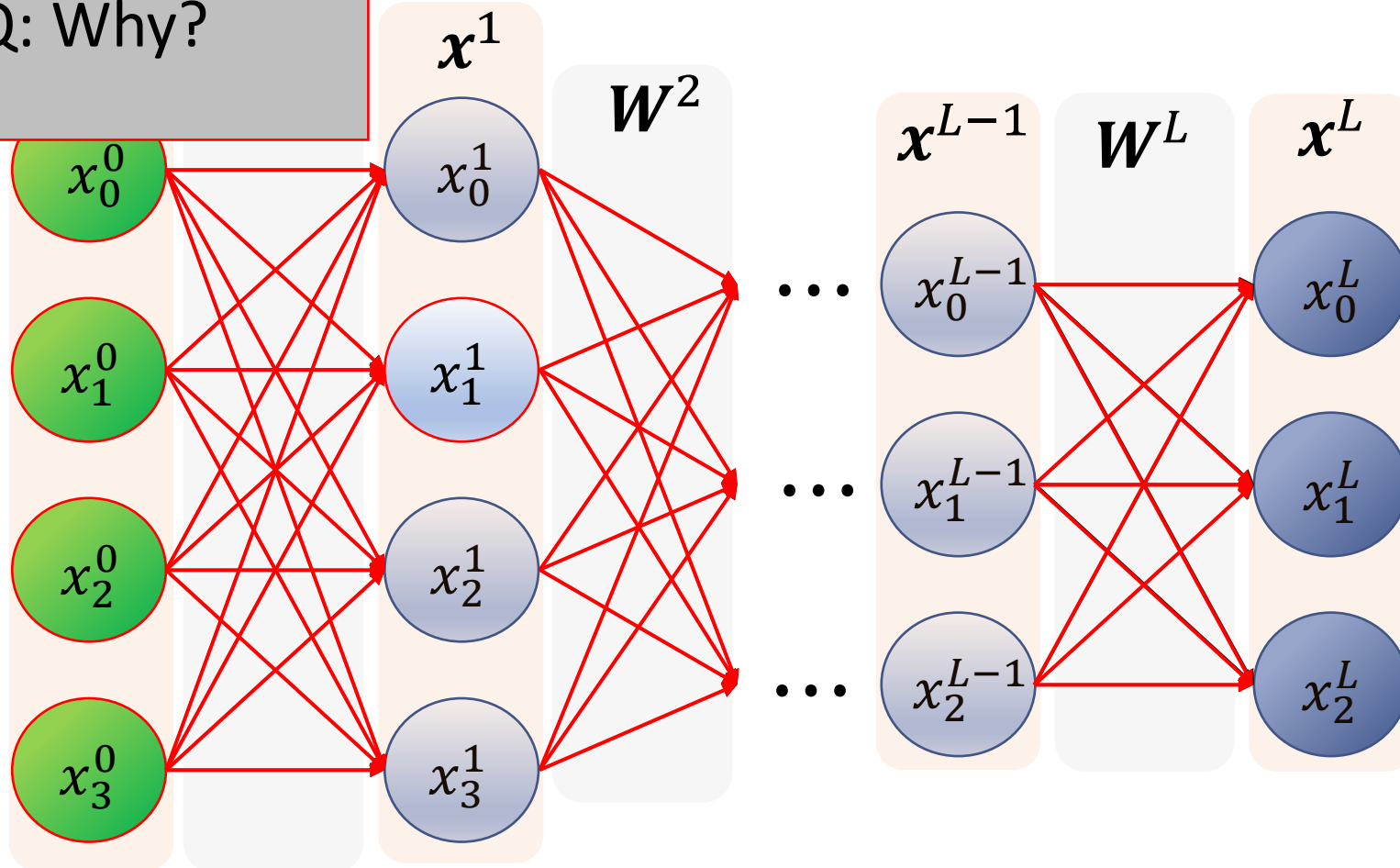
$$\mathbf{x}^l = \sigma \left(\mathbf{W}^{lT} \mathbf{x}^{l-1} + \mathbf{b}^l \right)$$



$$x_j^l = \sigma \left(\sum w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$

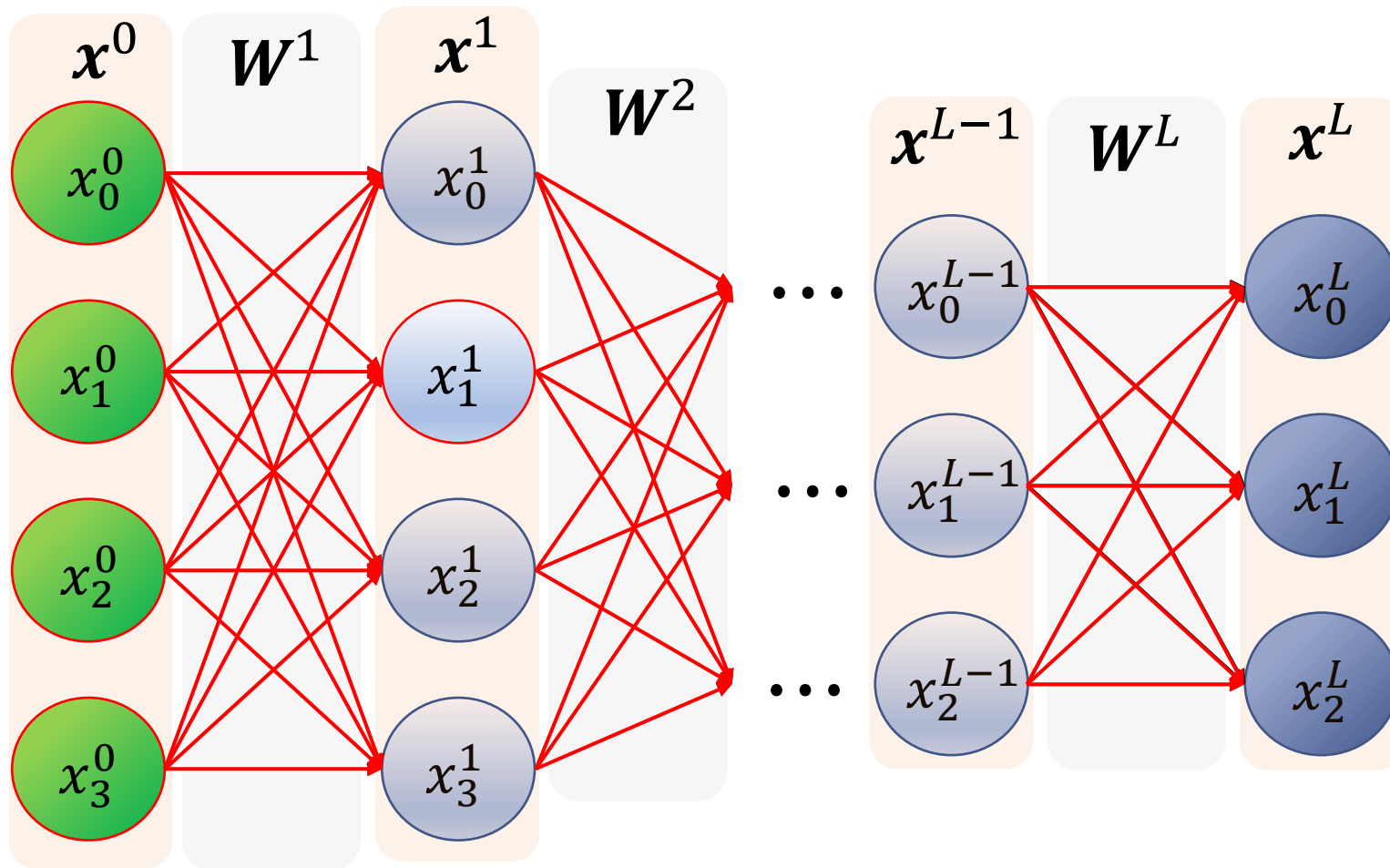
$$\mathbf{x}^l = \sigma \left(\mathbf{W}^{lT} \mathbf{x}^{l-1} + \mathbf{b}^l \right)$$

Q: Why?

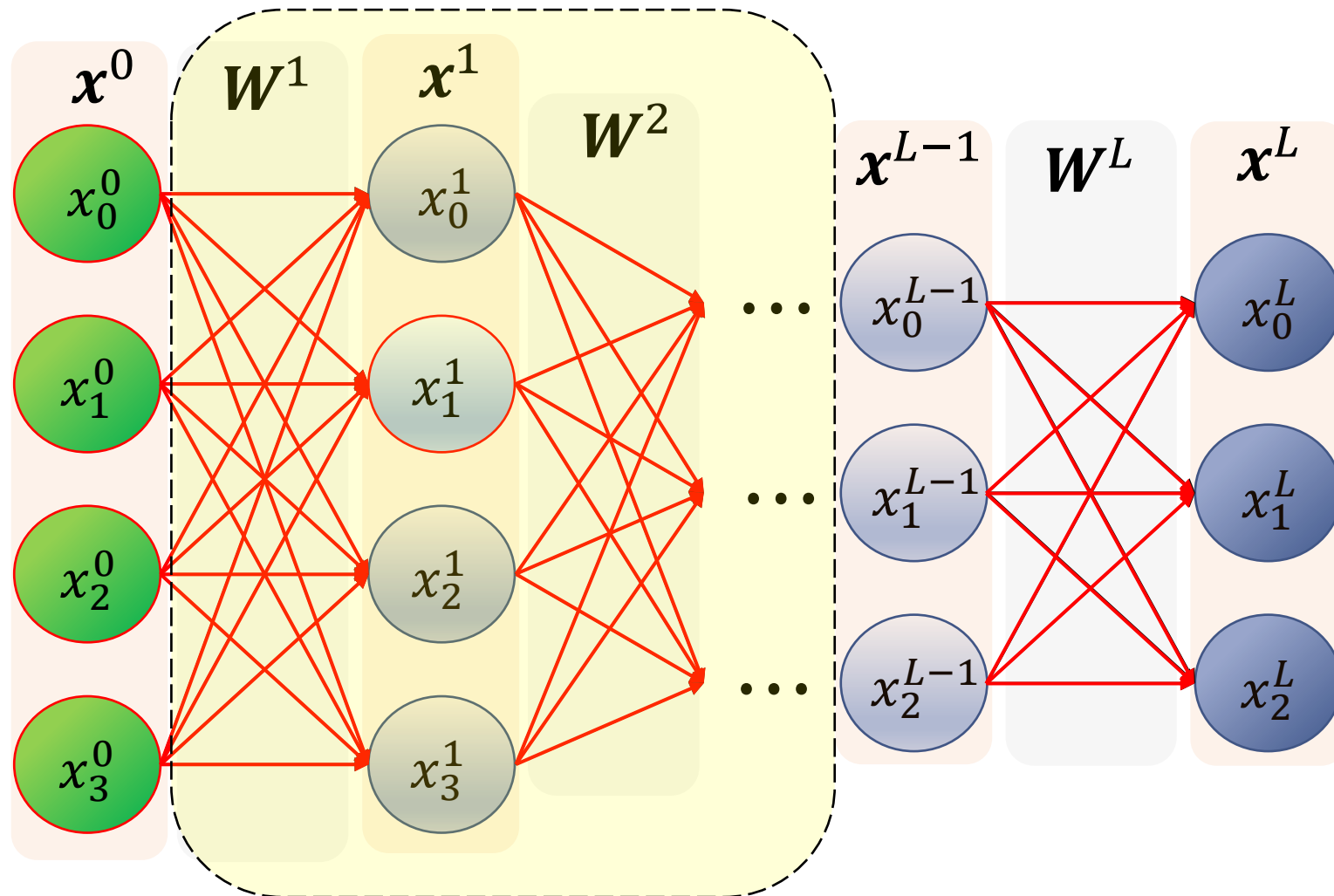


$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right)$$

$$\mathbf{x}^l = \sigma \left(\mathbf{W}^{lT} \mathbf{x}^{l-1} + \mathbf{b}^l \right)$$

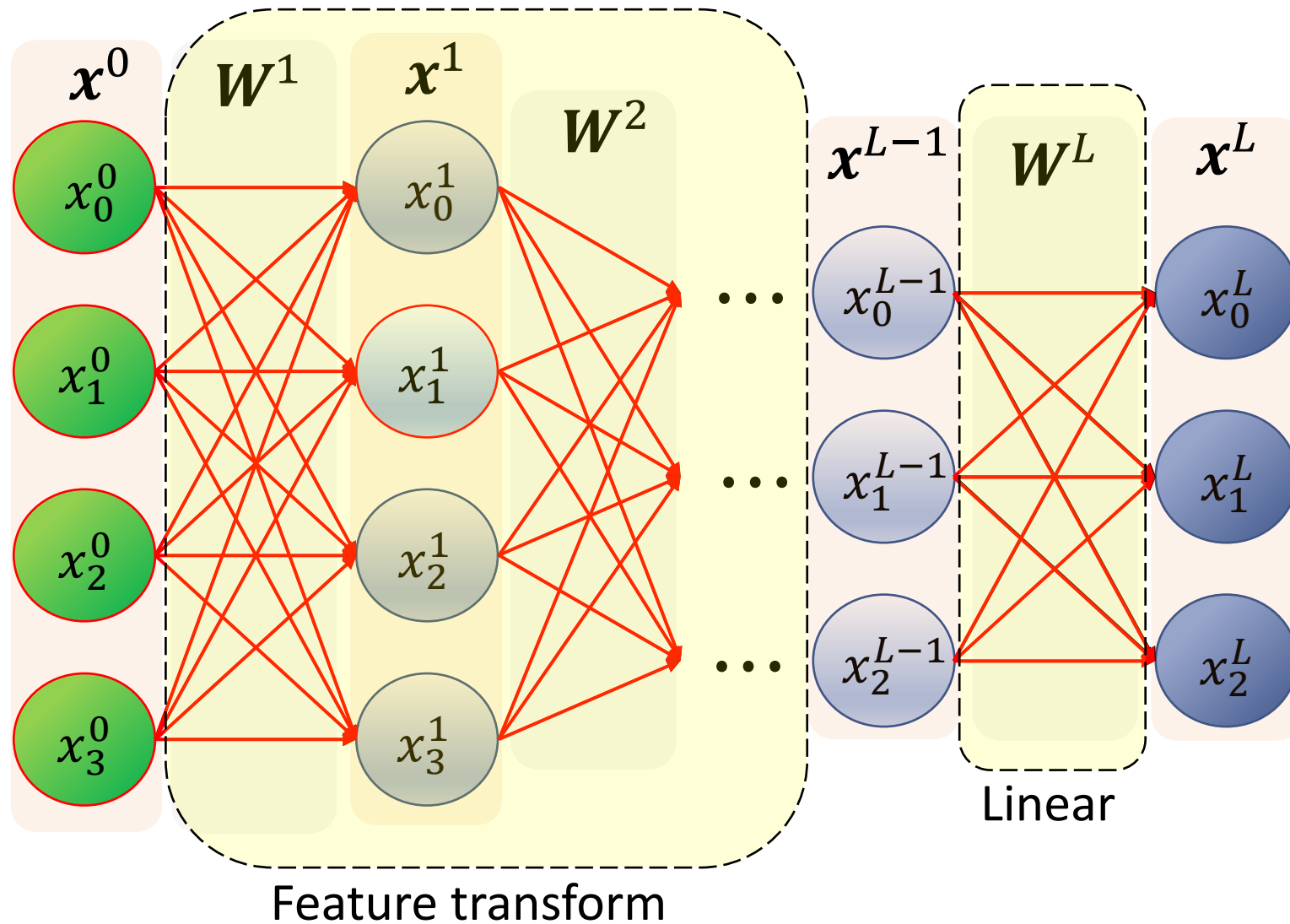


$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right) \quad \mathbf{x}^l = \boldsymbol{\sigma} \left(\mathbf{W}^{lT} \mathbf{x}^{l-1} + \mathbf{b}^l \right)$$

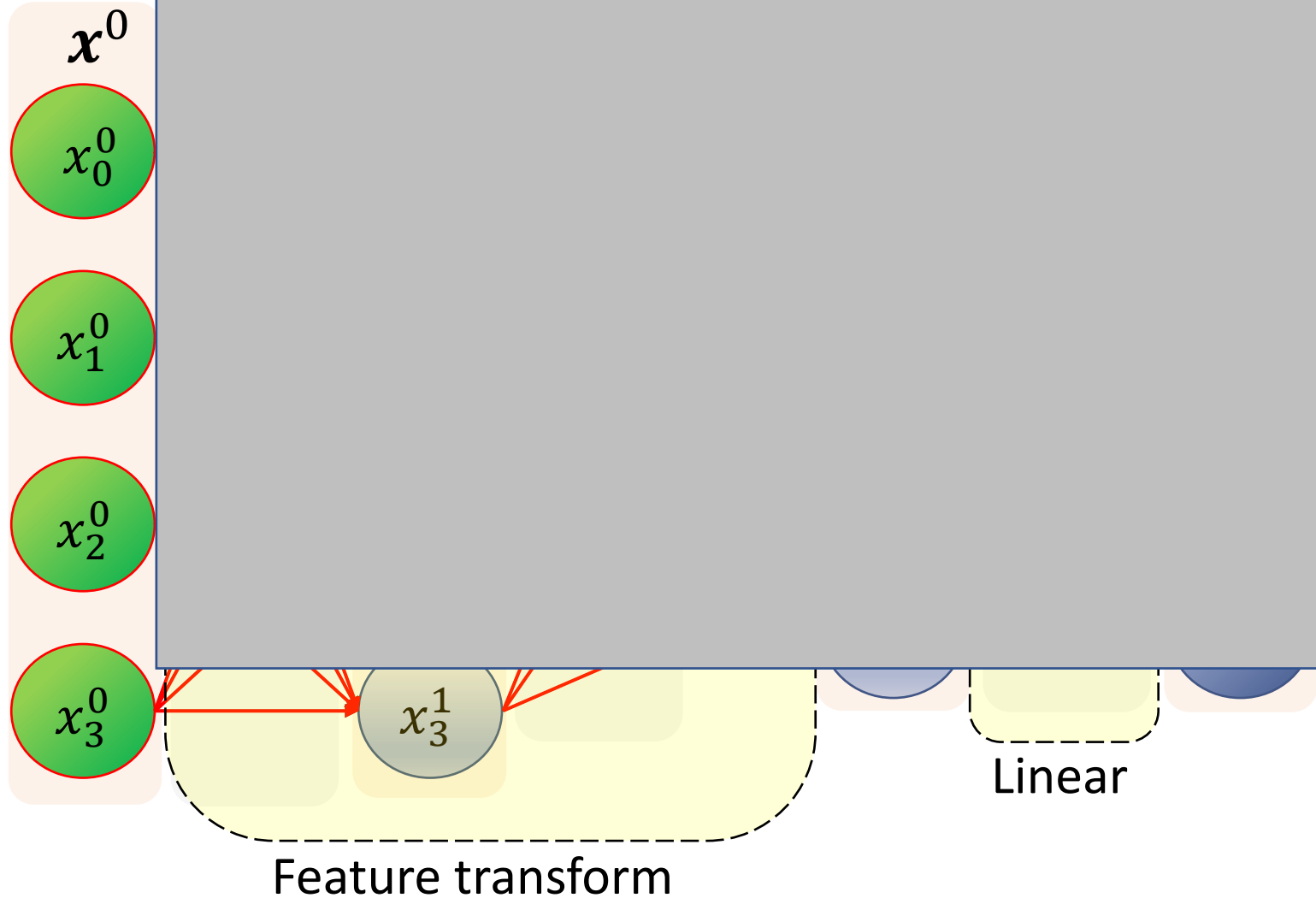


Feature transform

$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j^l \right) \quad \mathbf{x}^l = \boldsymbol{\sigma} \left(\mathbf{W}^{lT} \mathbf{x}^{l-1} + \mathbf{b}^l \right)$$

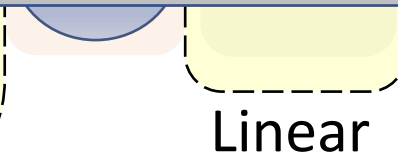
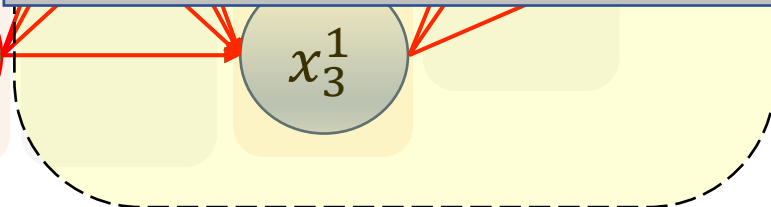
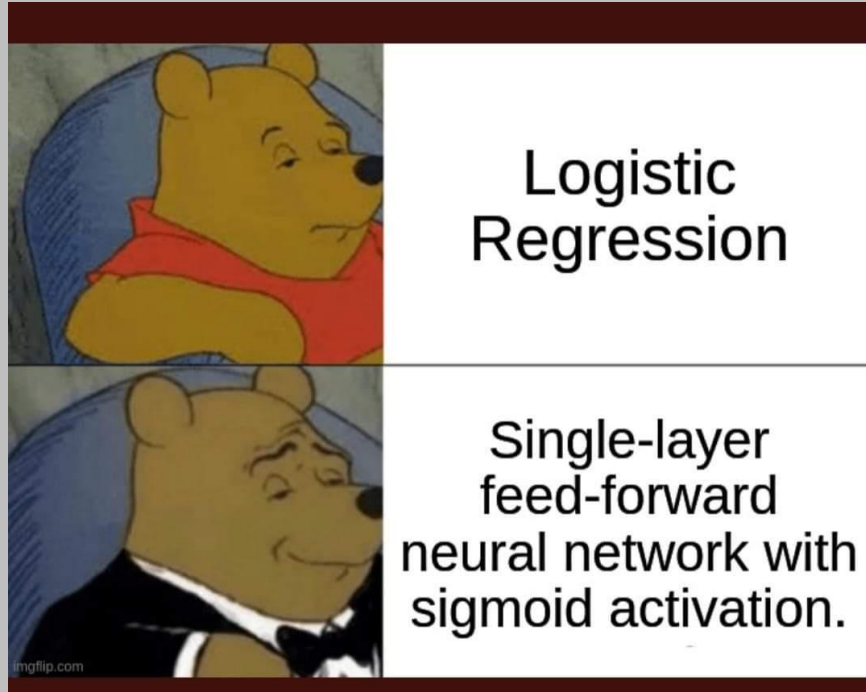
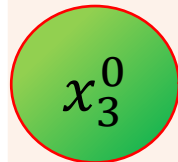
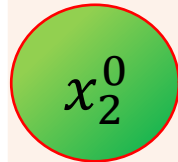
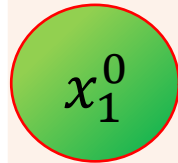
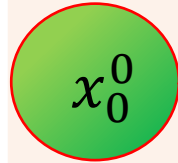


$$x_j^l = \sigma \left(\text{Q: What do you call a single layered net? } b^l \right)$$

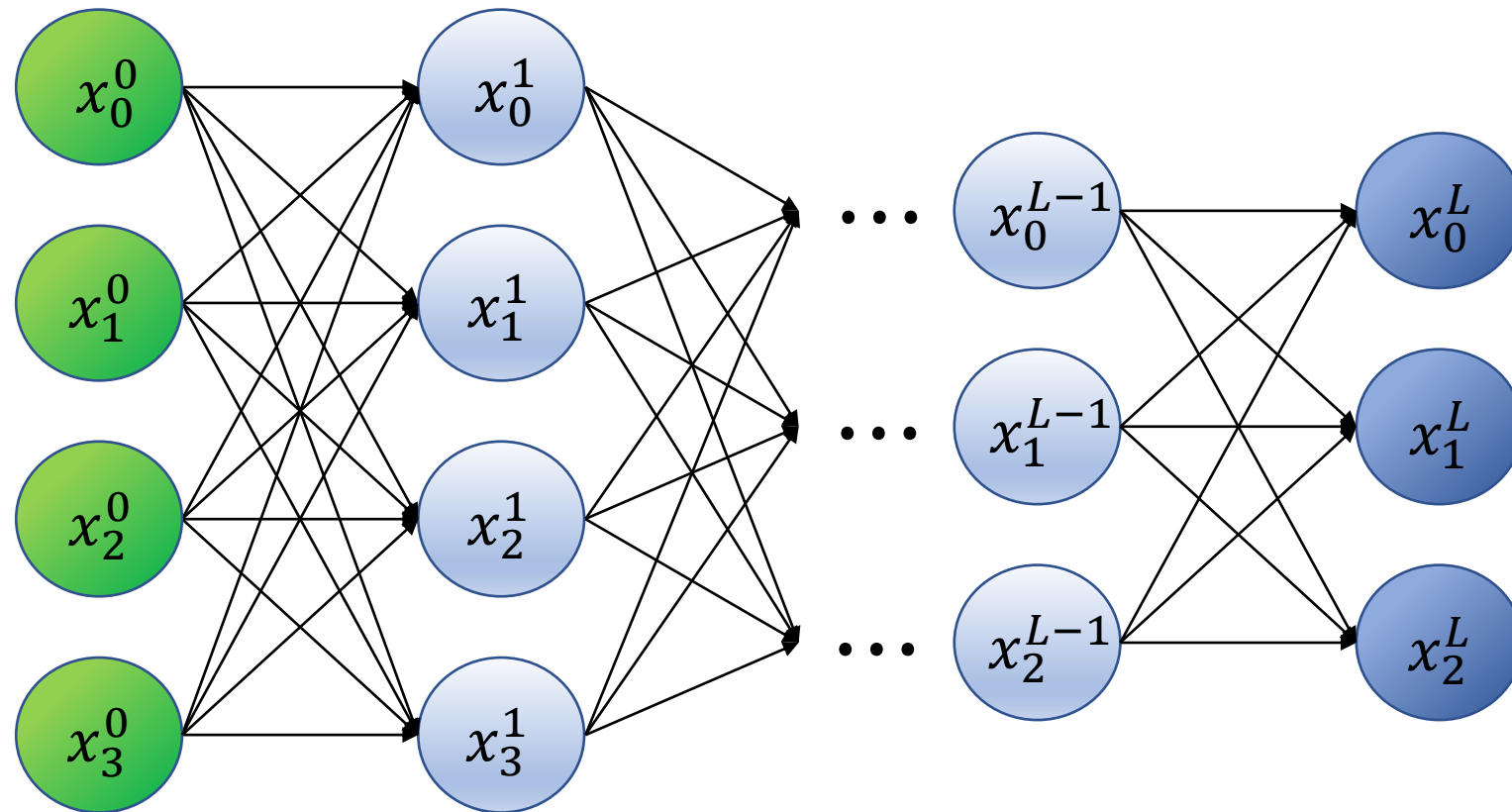


$$x_j^l = \sigma \left(\text{Q: What do you call a single layered net?} \right)$$

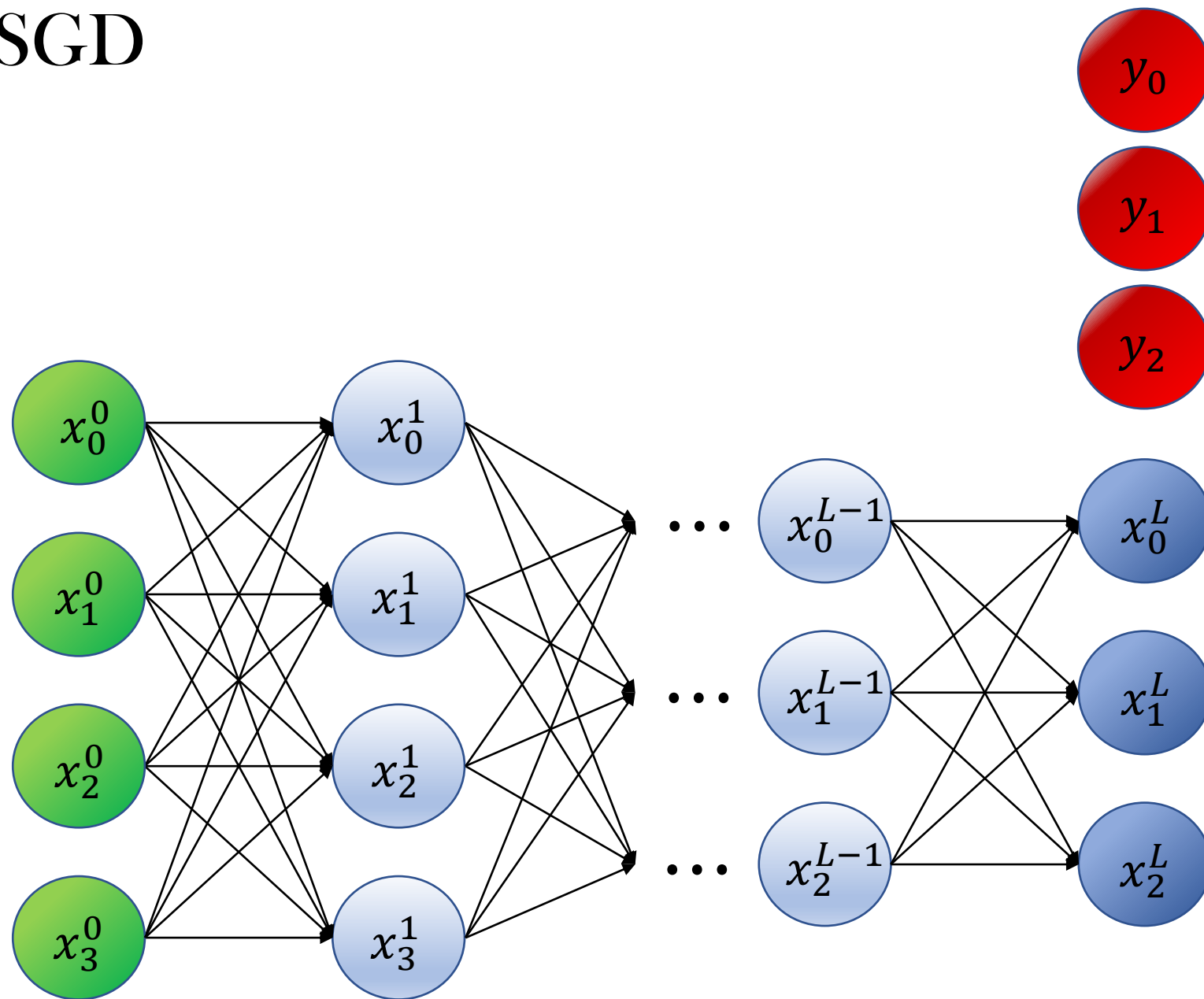
x^0



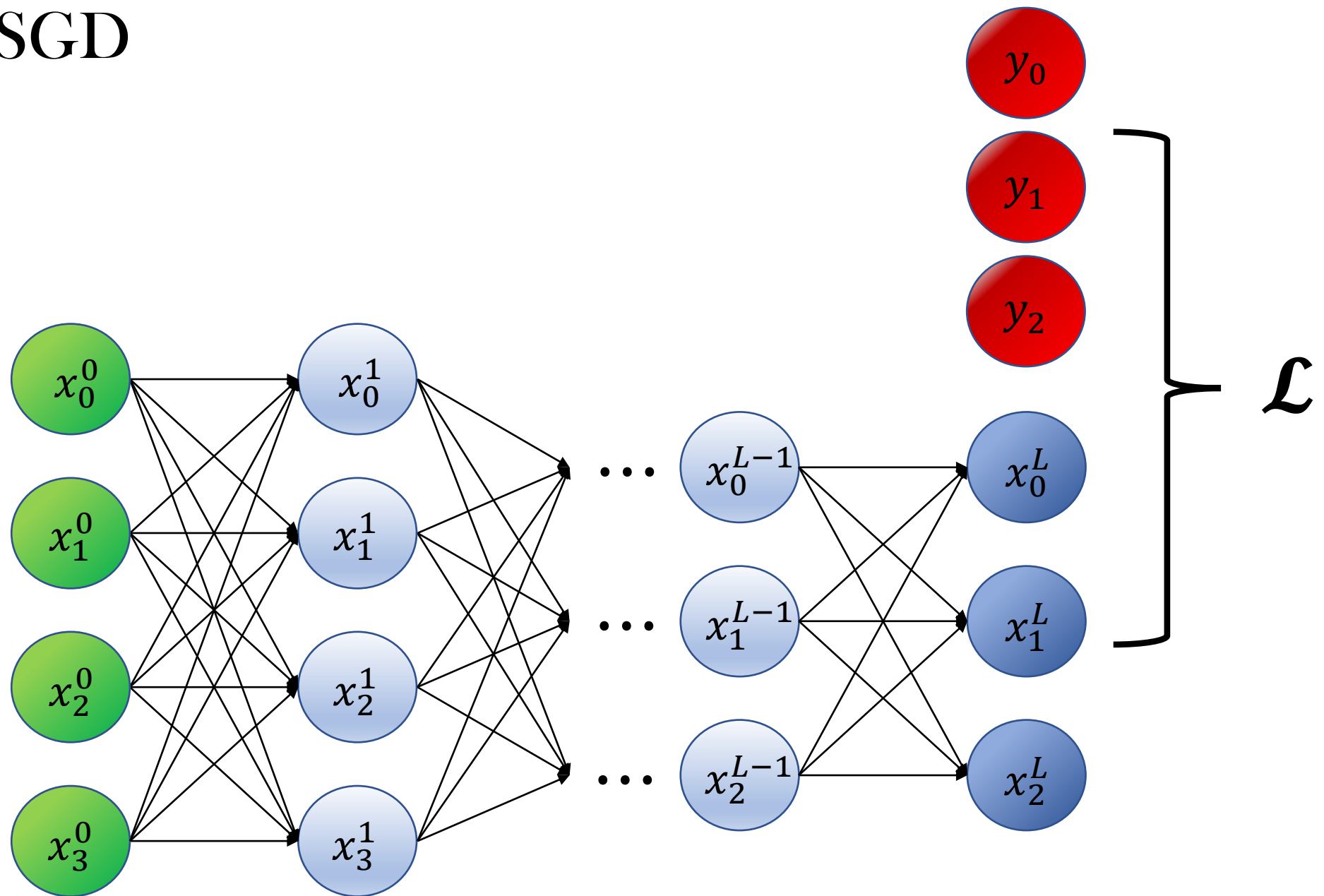
Learning by SGD



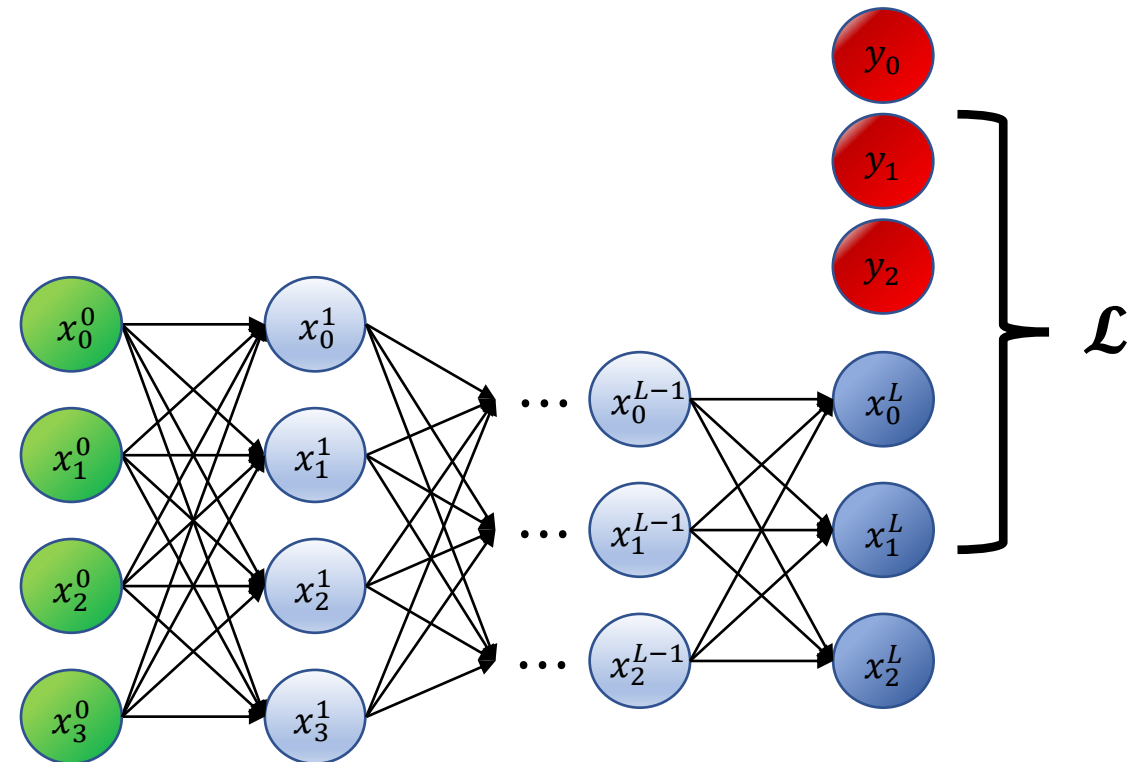
Learning by SGD



Learning by SGD

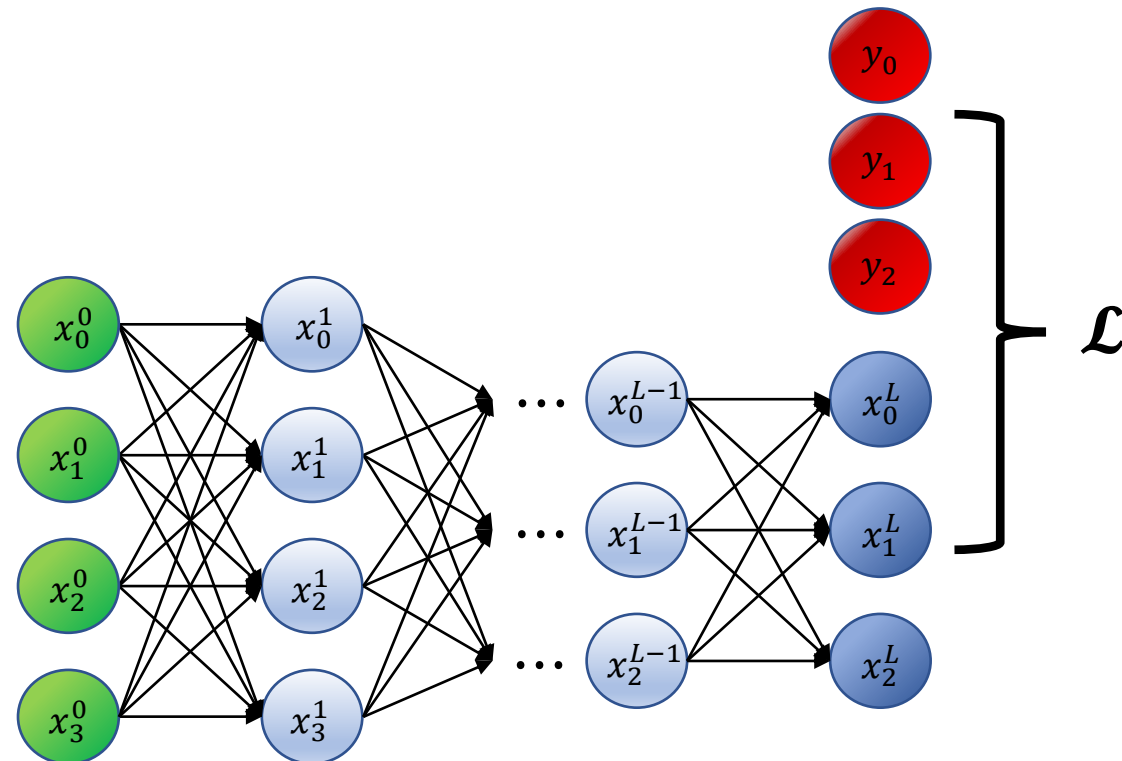


Learning by SGD



Learning by SGD

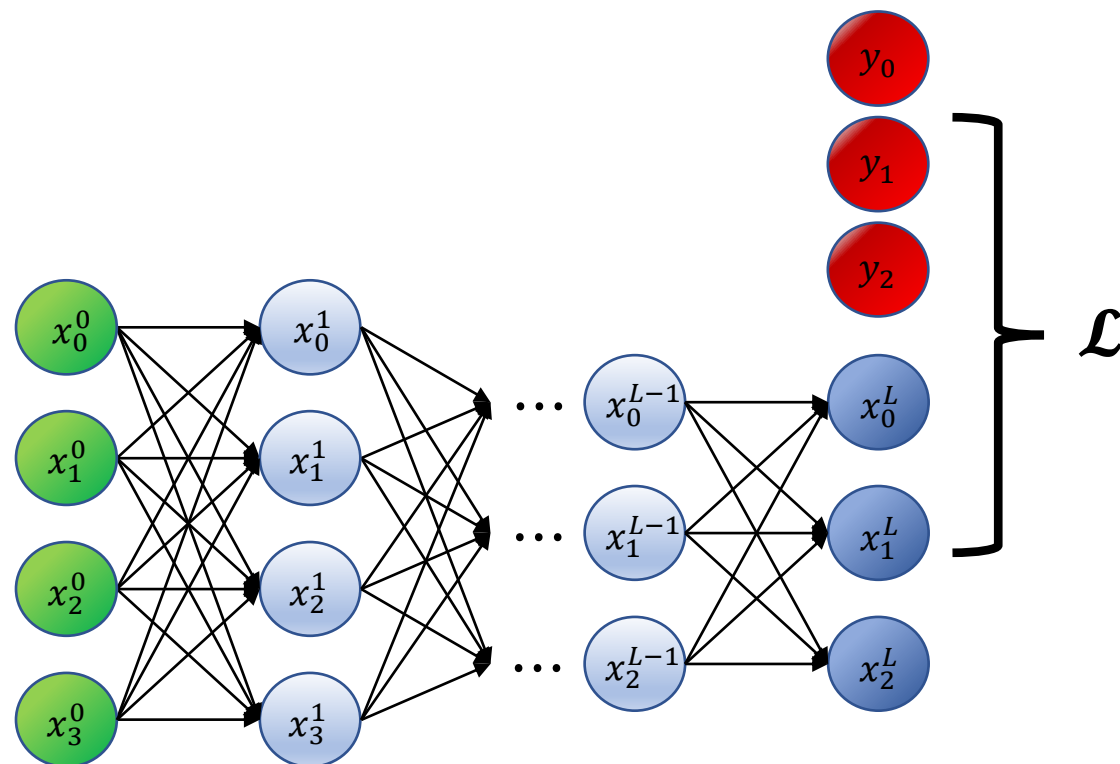
We need $\frac{\partial \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})}{\partial w_{ij}^l}$, $\frac{\partial \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})}{\partial b_j^l}$ To all l, i, j



Learning by SGD

We need $\frac{\partial \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})}{\partial w_{ij}^l}$, $\frac{\partial \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})}{\partial b_j^l}$ To all l, i, j

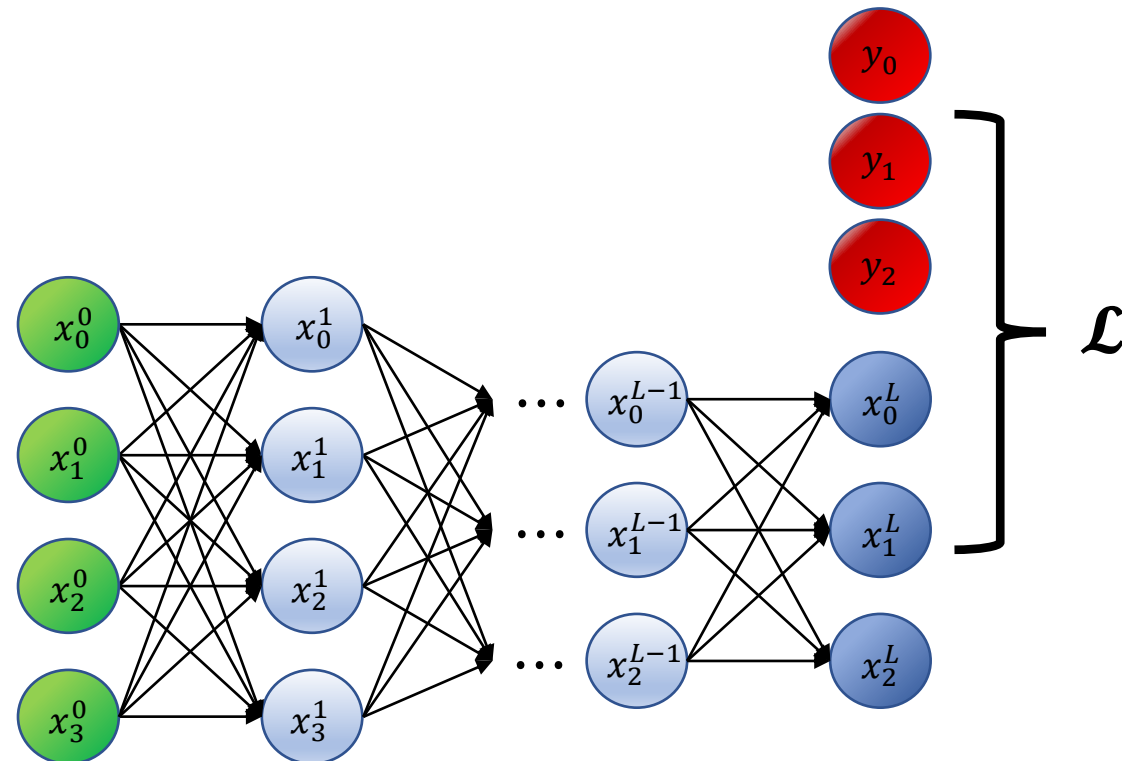
$$x_j^L = \sum_i W_{ij}^L x_i^{L-1} + b_j^L$$



Learning by SGD

We need $\frac{\partial \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})}{\partial w_{ij}^l}$, $\frac{\partial \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})}{\partial b_j^l}$ To all l, i, j

$$x_j^L = \sum_i w_{ij}^L x_i^{L-1} + b_j^L$$

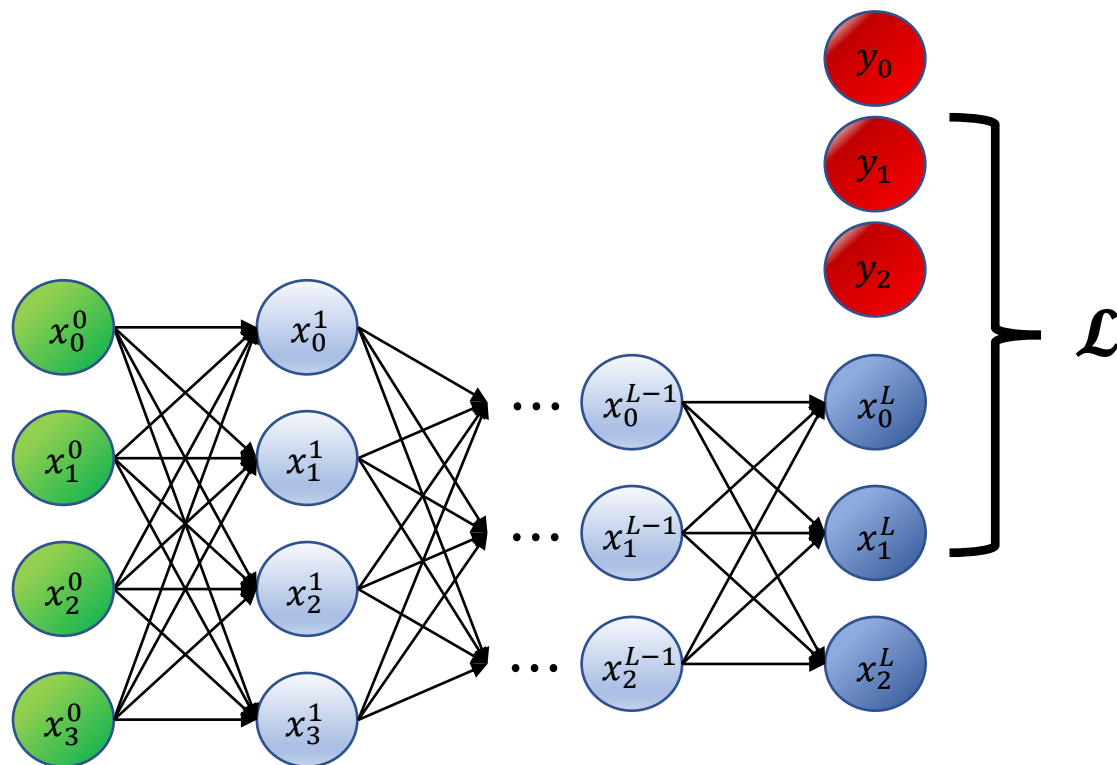


Learning by SGD

We need $\frac{\partial \mathcal{L}(\theta; \mathbf{x}, \mathbf{y})}{\partial w_{ij}^l}$, $\frac{\partial \mathcal{L}(\theta; \mathbf{x}, \mathbf{y})}{\partial b_j^l}$

To all l, i, j

$$\begin{aligned} x_j^L &= \sum_i W_{ij}^L x_i^{L-1} + b_j^L \\ &= \sum_i W_{ij}^L \sigma \left(\sum_k W_{ki}^{L-1} x_k^{L-2} + b_i^{L-1} \right) + b_j^L \end{aligned}$$

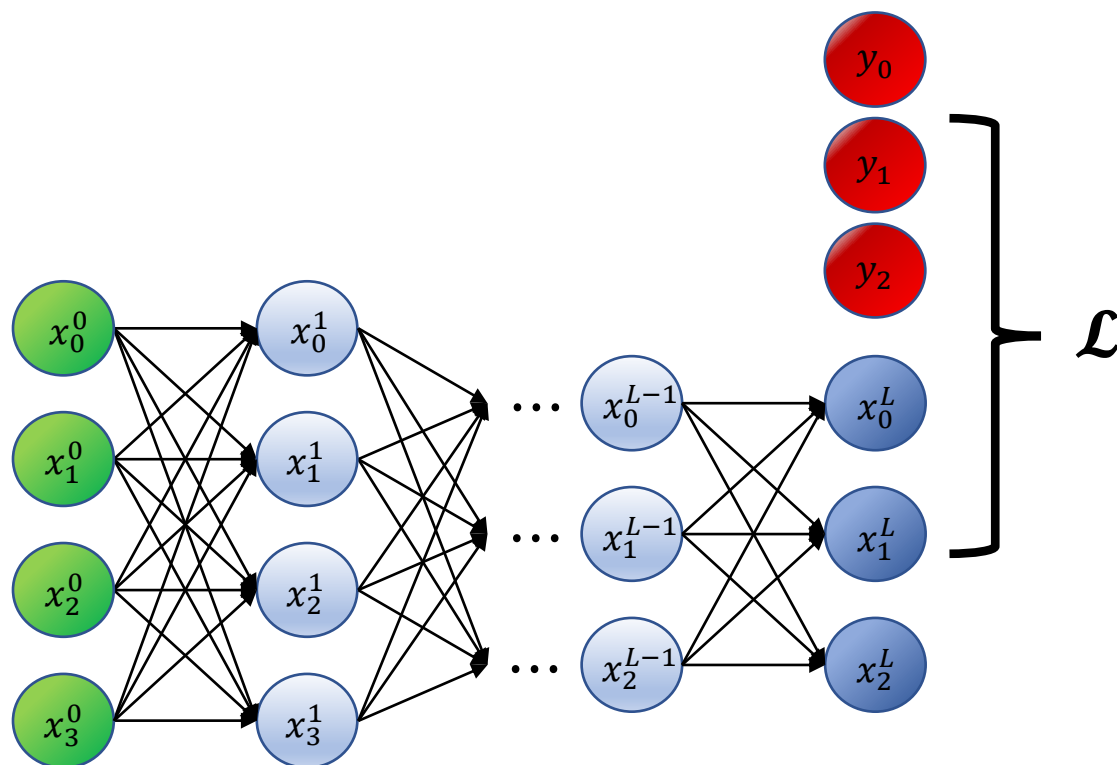


Learning by SGD

We need $\frac{\partial \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})}{\partial w_{ij}^l}$, $\frac{\partial \mathcal{L}(\boldsymbol{\theta}; \mathbf{x}, \mathbf{y})}{\partial b_j^l}$

To all l, i, j

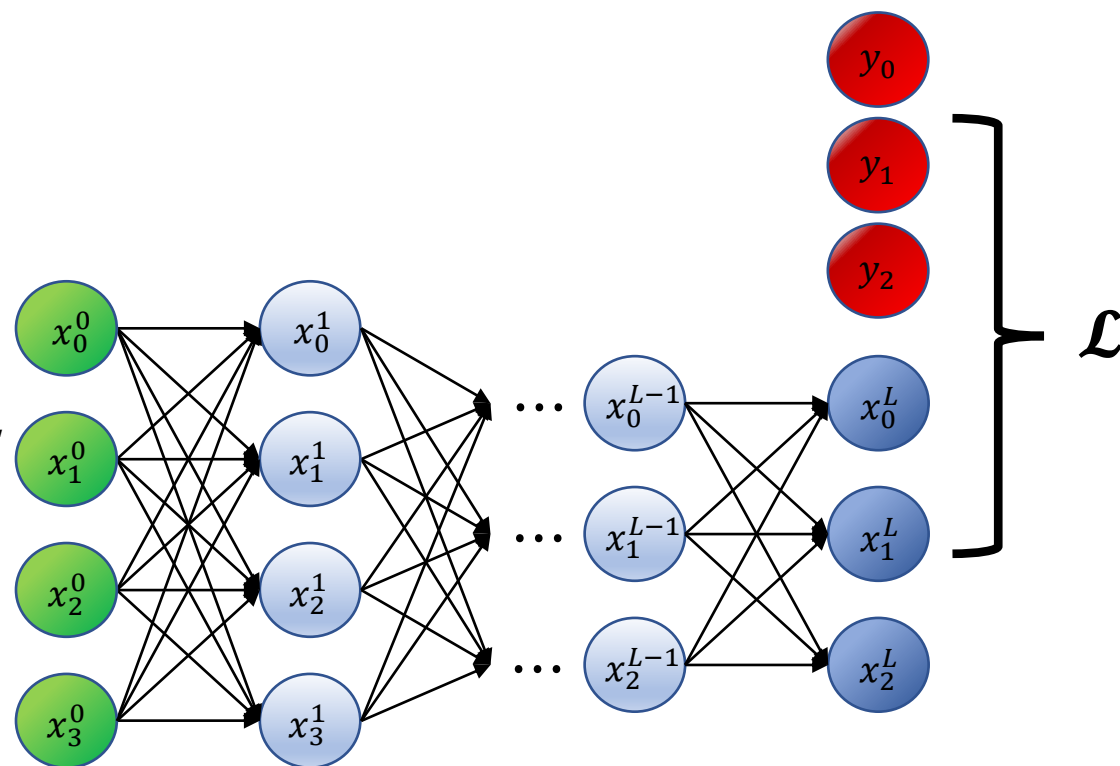
$$\begin{aligned} x_j^l &= \sum_i w_{ij}^l x_i^{l-1} + b_j^l \\ &= \sum_i w_{ij}^l \sigma \left(\sum_k w_{ki}^{l-1} x_k^{l-2} + b_i^{l-1} \right) + b_j^l \end{aligned}$$



Learning by SGD

We need $\frac{\partial \mathcal{L}(\theta; \mathbf{x}, \mathbf{y})}{\partial w_{ij}^l}$, $\frac{\partial \mathcal{L}(\theta; \mathbf{x}, \mathbf{y})}{\partial b_j^l}$ To all l, i, j

$$\begin{aligned}
 x_j^L &= \sum_i W_{ij}^L x_i^{L-1} + b_j^L \\
 &= \sum_i W_{ij}^L \sigma \left(\sum_k W_{ki}^{L-1} x_k^{L-2} + b_i^{L-1} \right) + b_j^L \\
 &= \sum_i W_{ij}^L \sigma \left(\sum_k W_{ki}^{L-1} \sigma \left(\sum_m W_{mk}^{L-2} x_m^{L-3} + b_k^{L-2} \right) + b_i^{L-1} \right) + b_j^L \\
 &\quad \vdots
 \end{aligned}$$



Chain rule reminder



Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$



Chain rule reminder

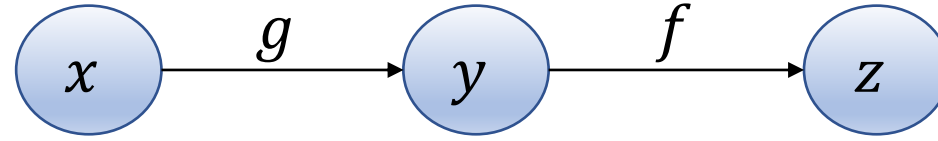
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

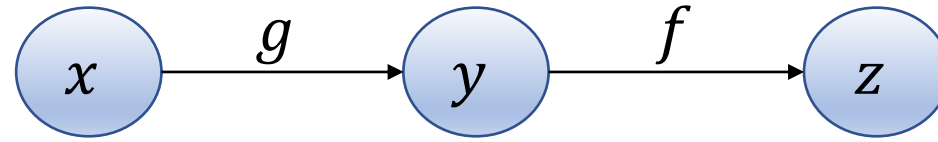


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} =$$

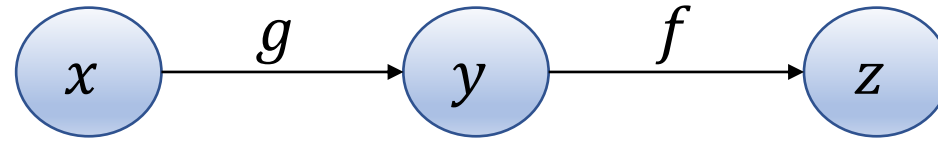


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dx}$$

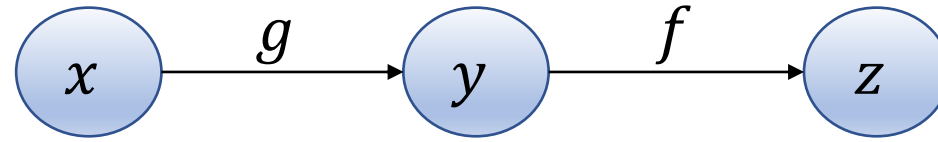


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

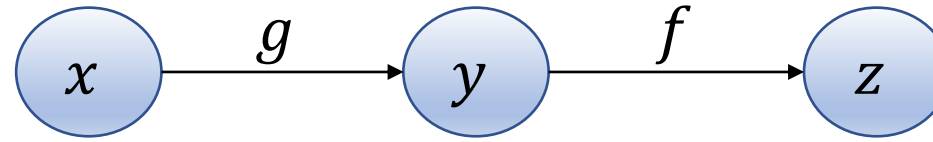


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

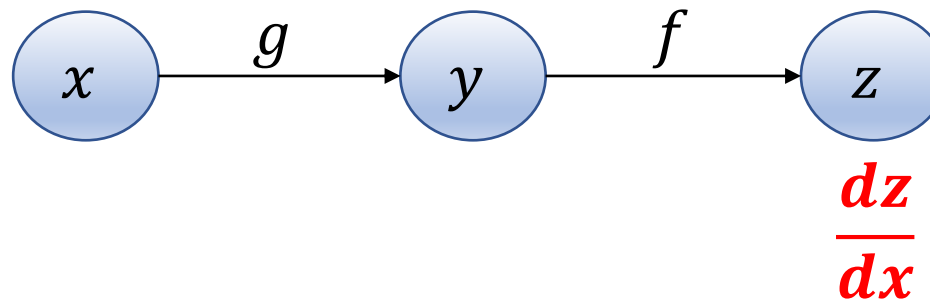


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

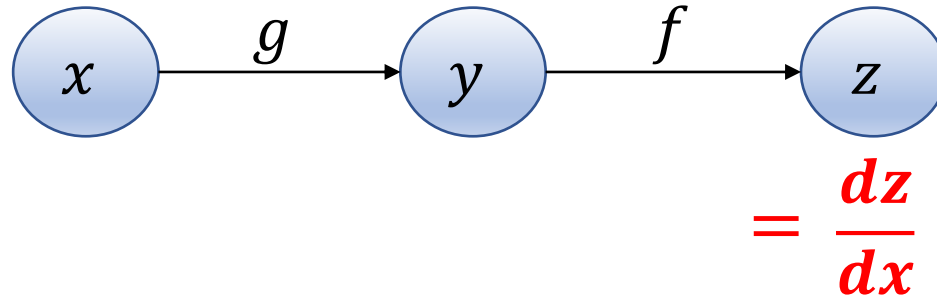


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

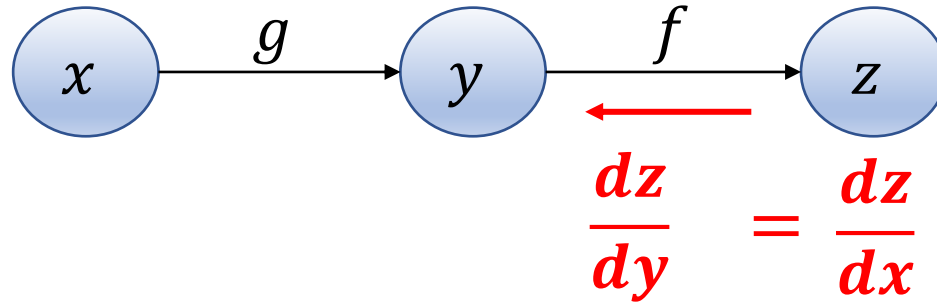


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

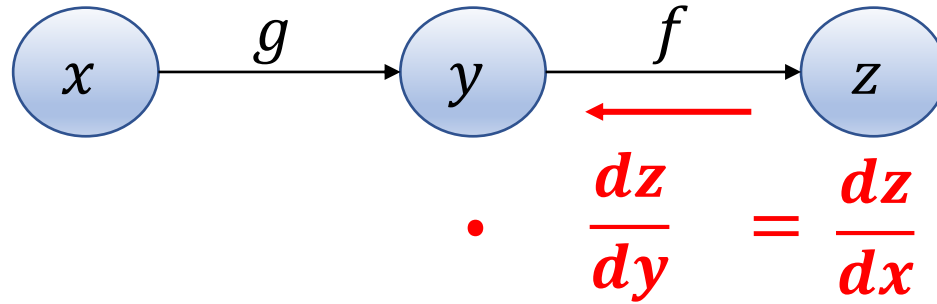


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

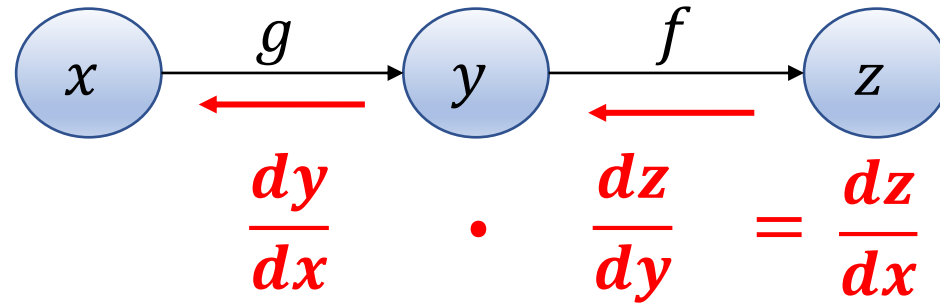


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

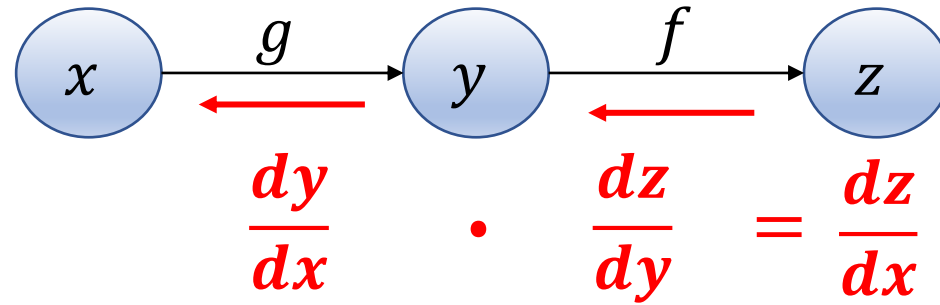


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

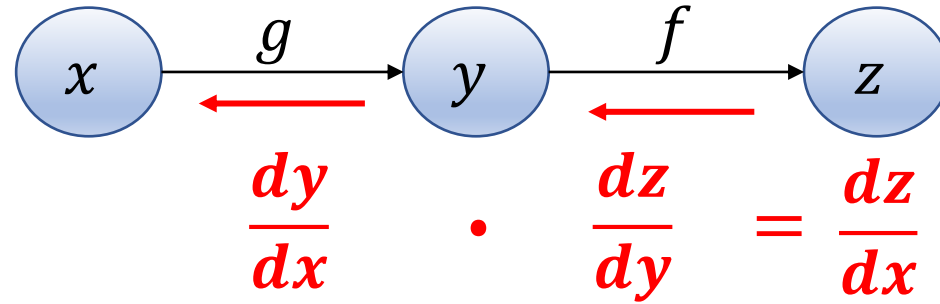


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



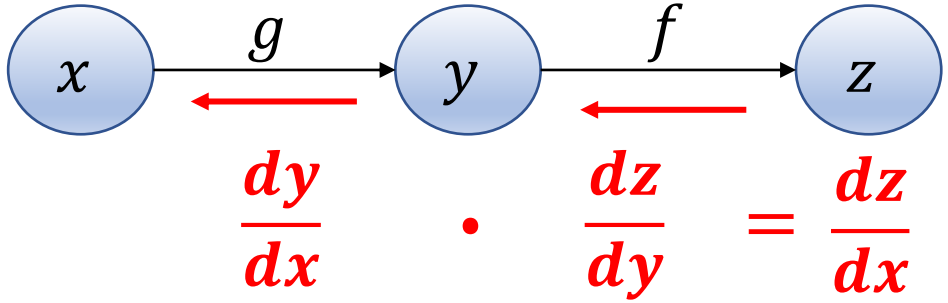
$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

Chain rule reminder

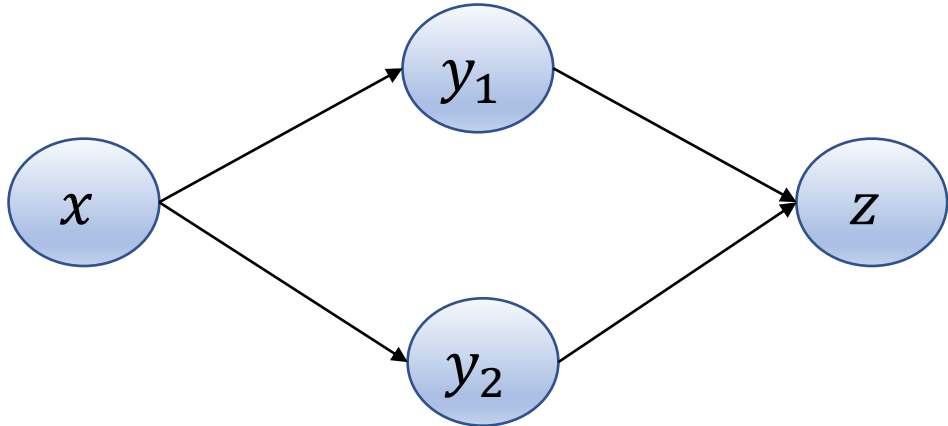
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

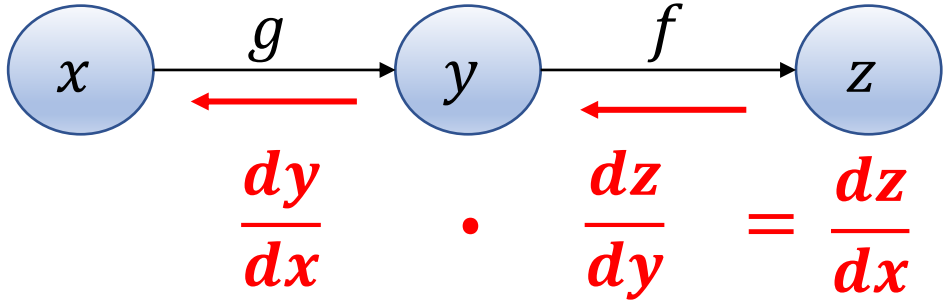


Chain rule reminder

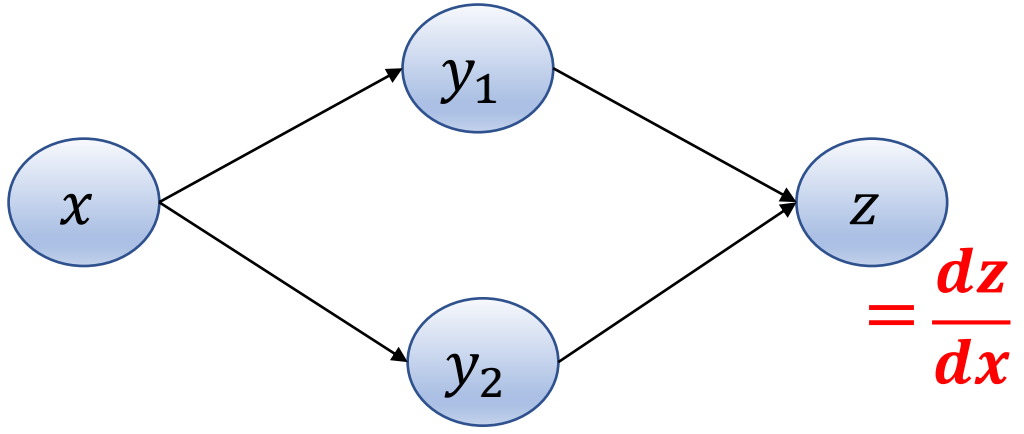
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

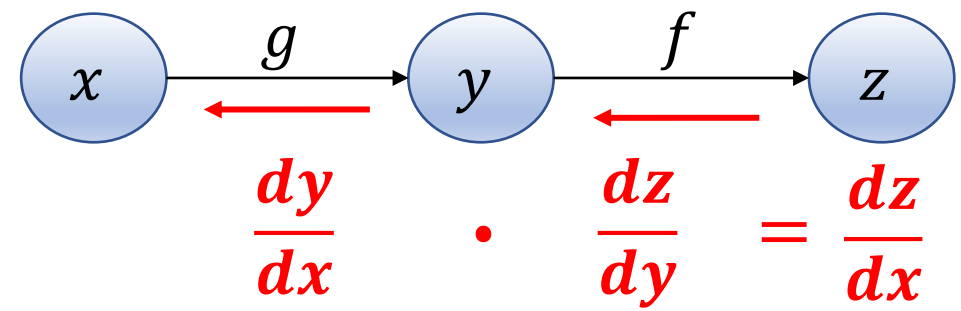


Chain rule reminder

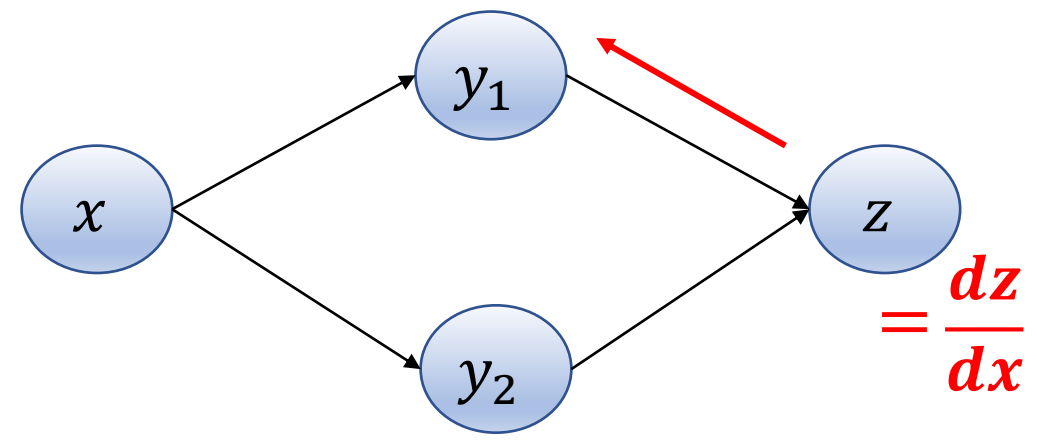
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

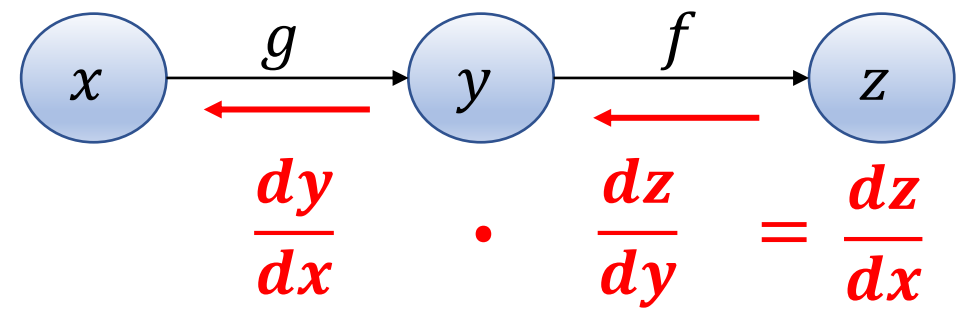


Chain rule reminder

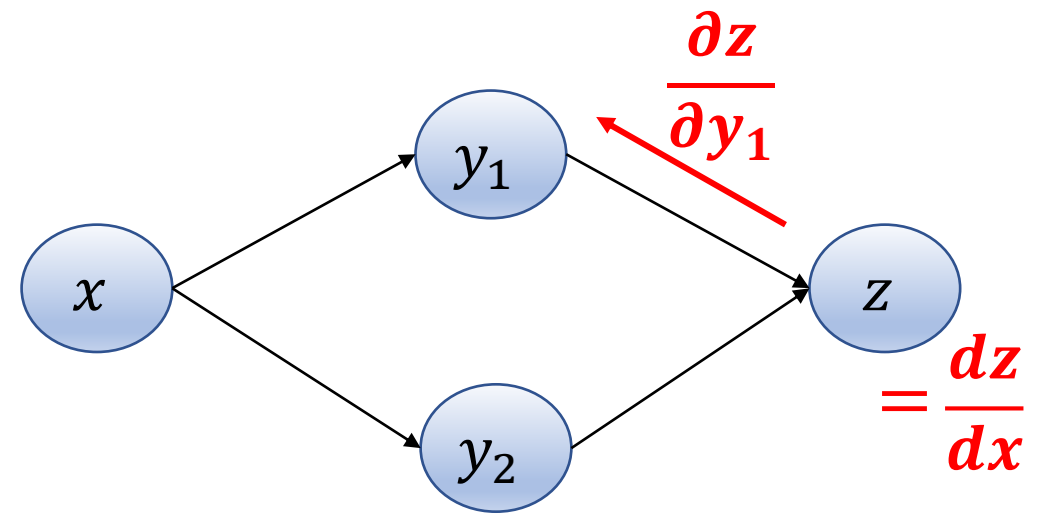
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

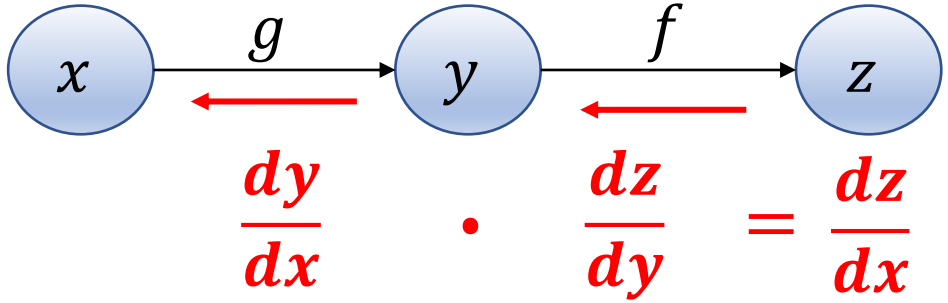


Chain rule reminder

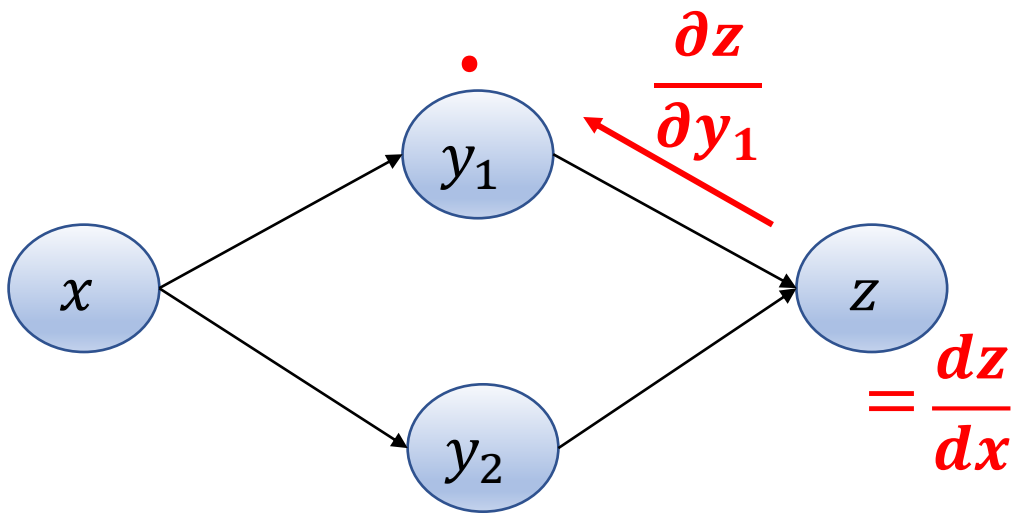
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

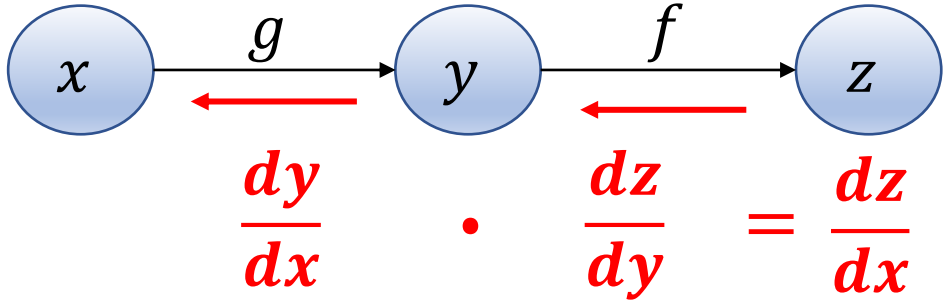


Chain rule reminder

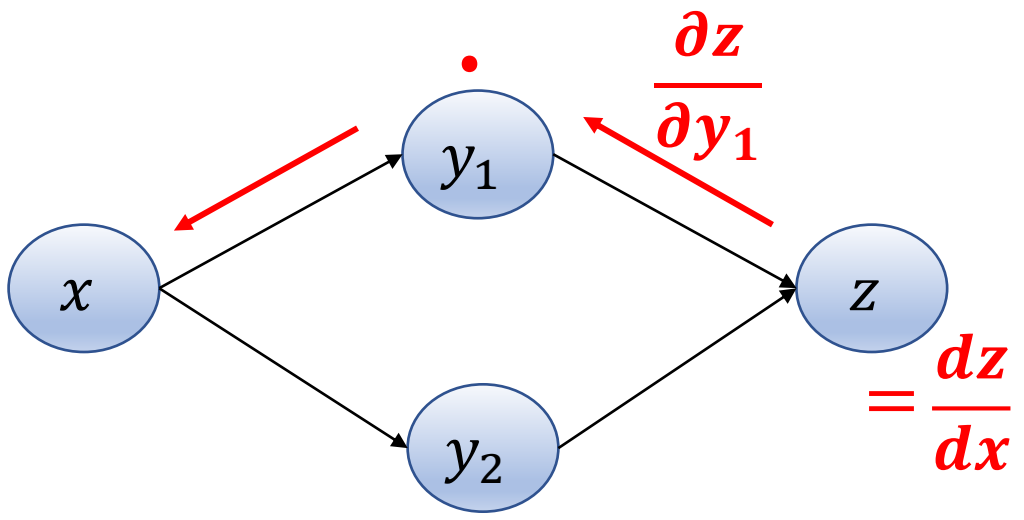
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$



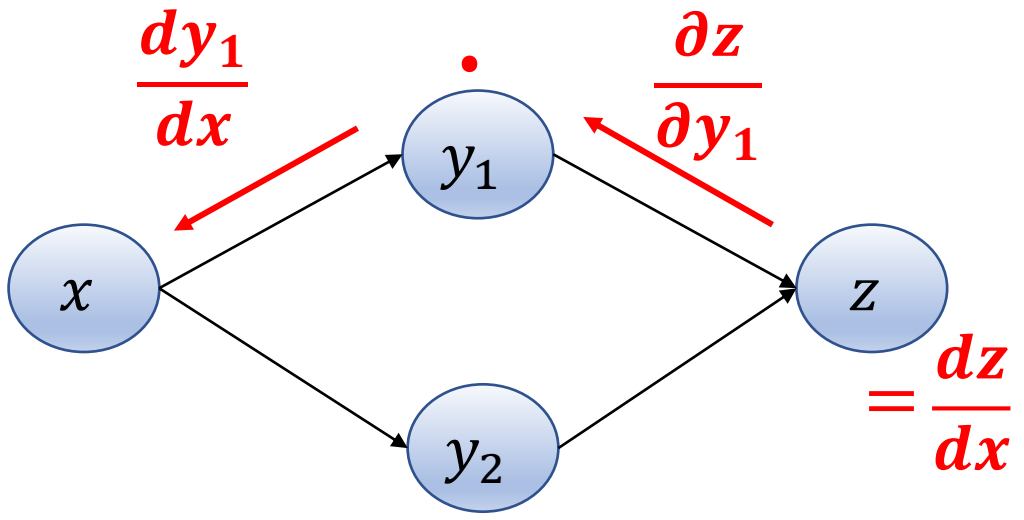
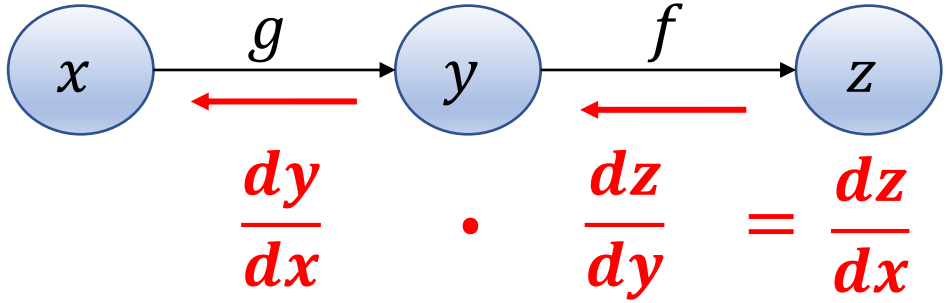
Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

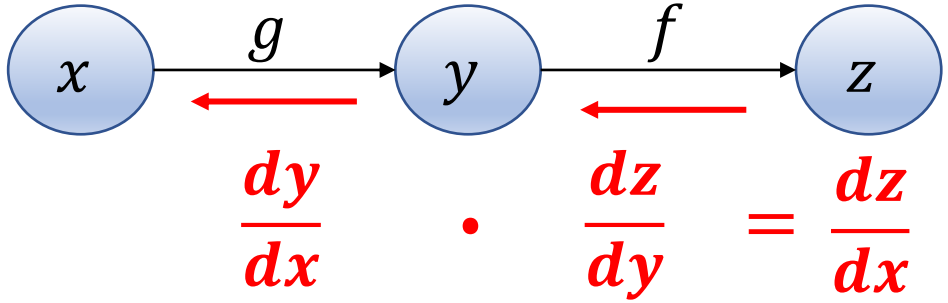


Chain rule reminder

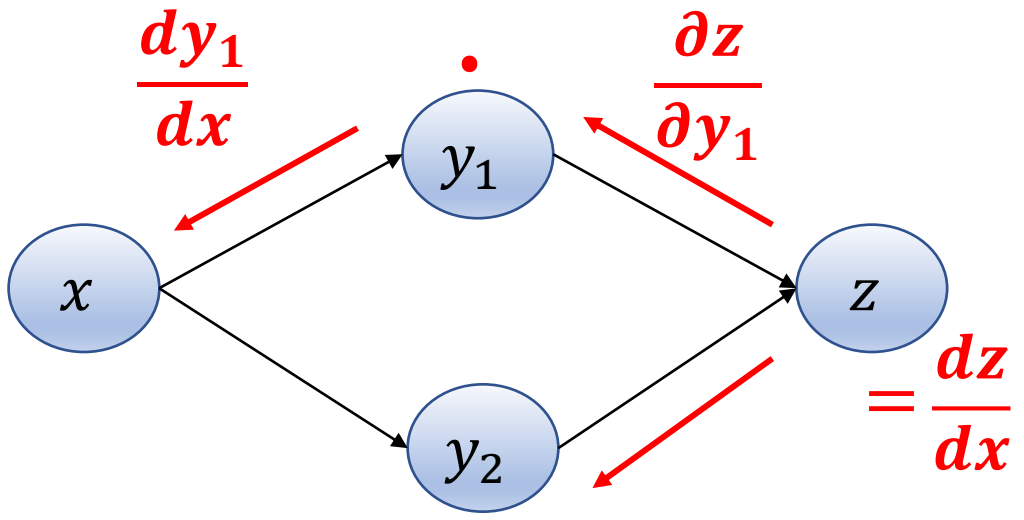
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

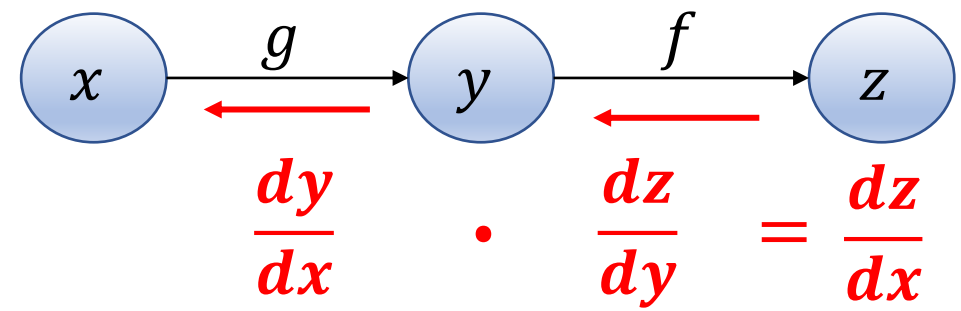


Chain rule reminder

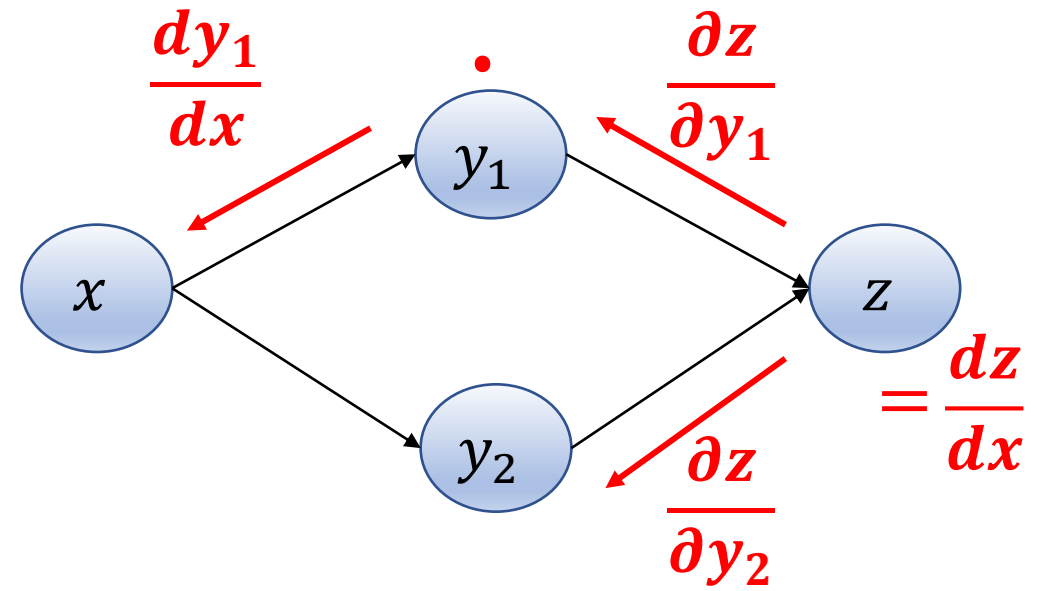
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

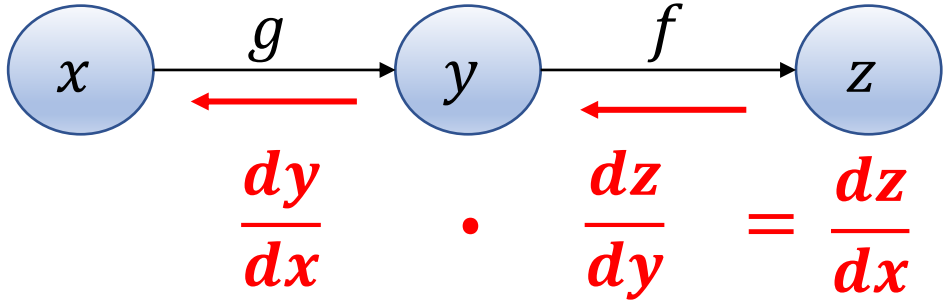


Chain rule reminder

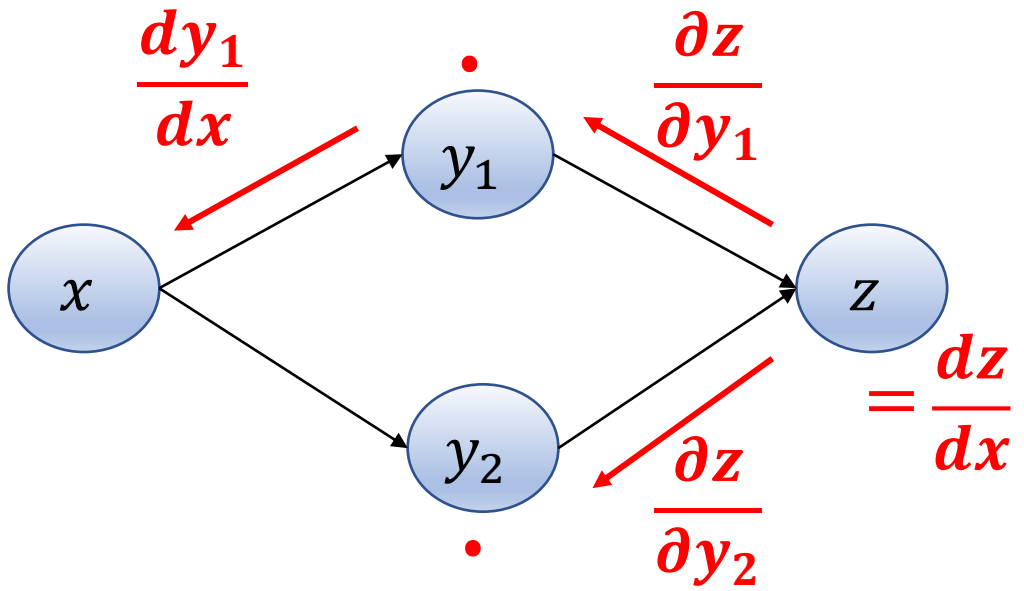
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

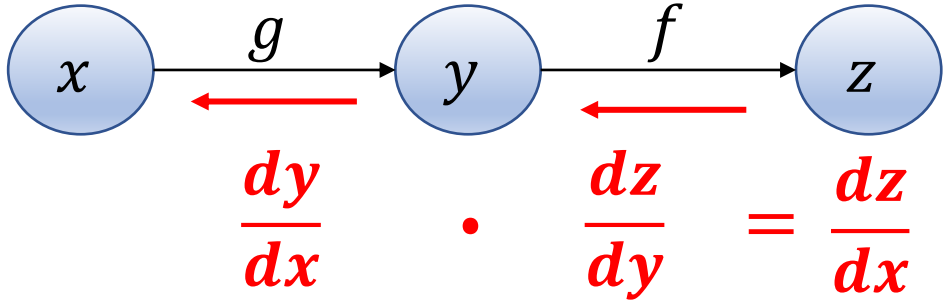


Chain rule reminder

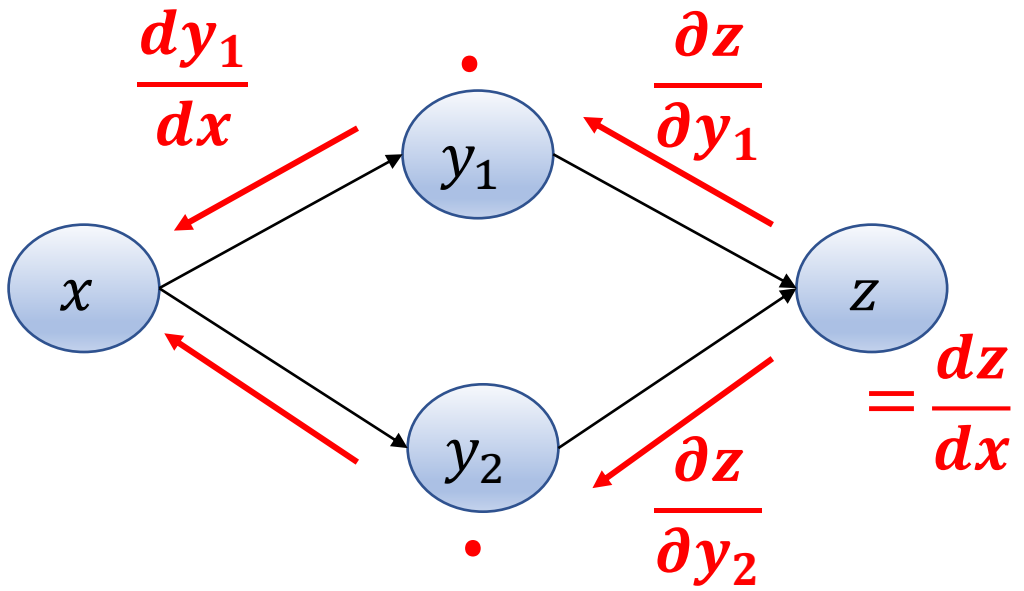
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

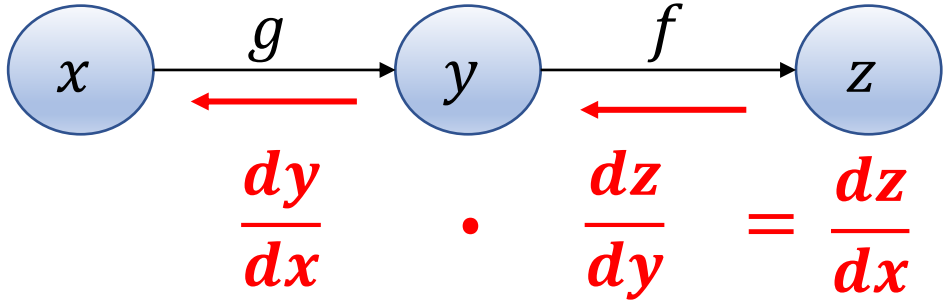


Chain rule reminder

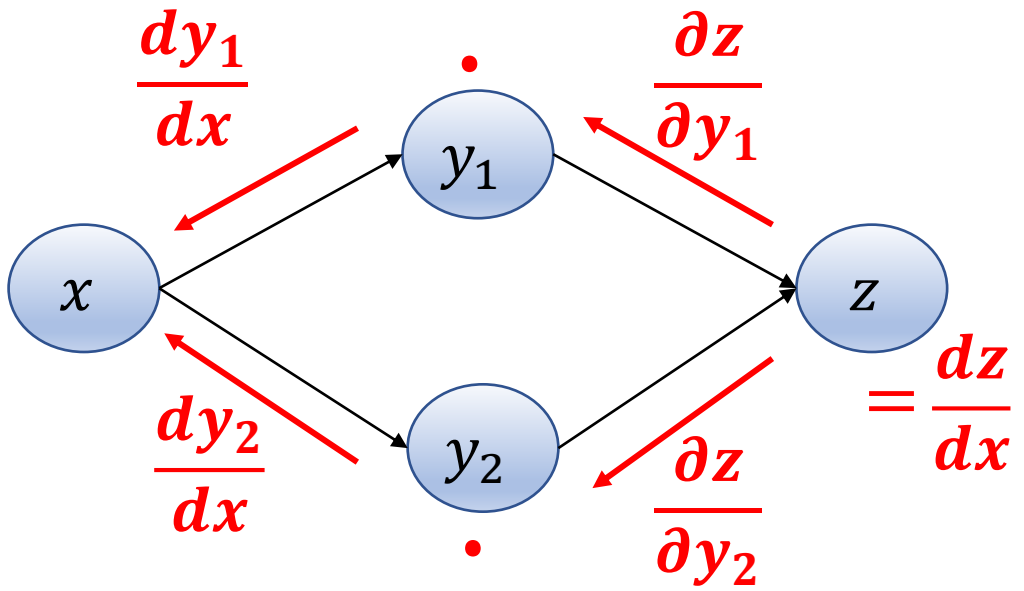
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

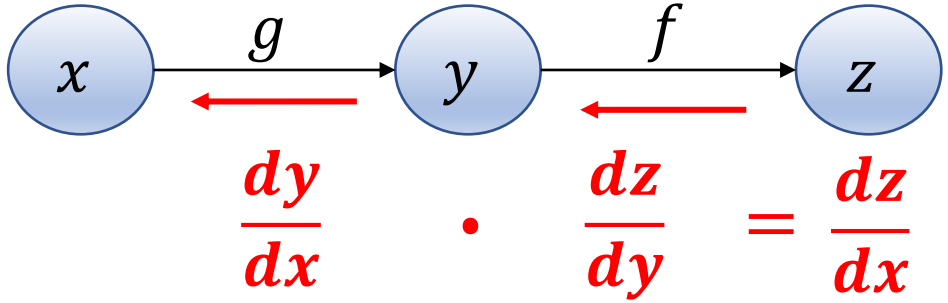


Chain rule reminder

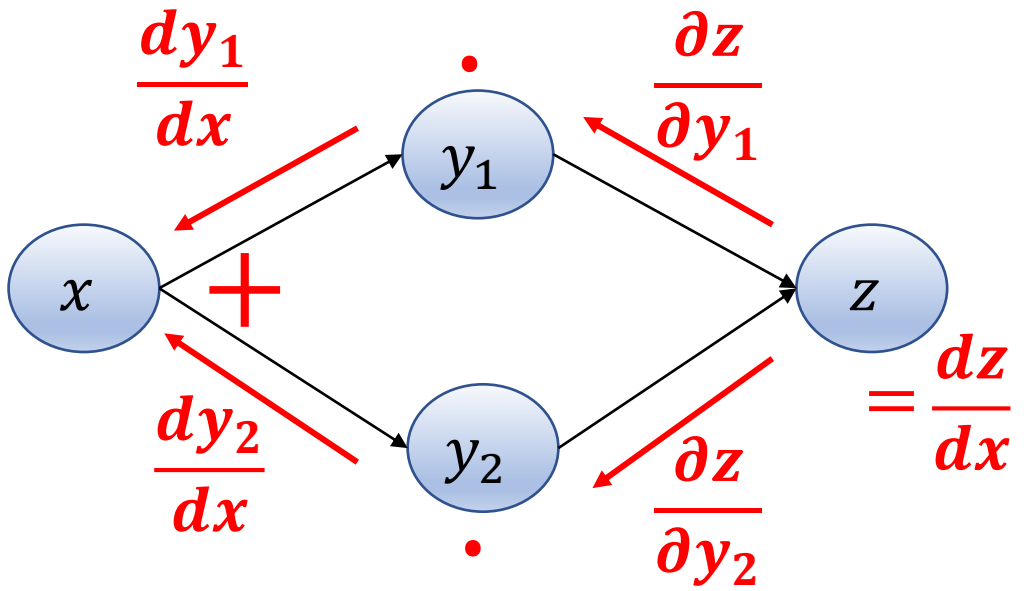
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

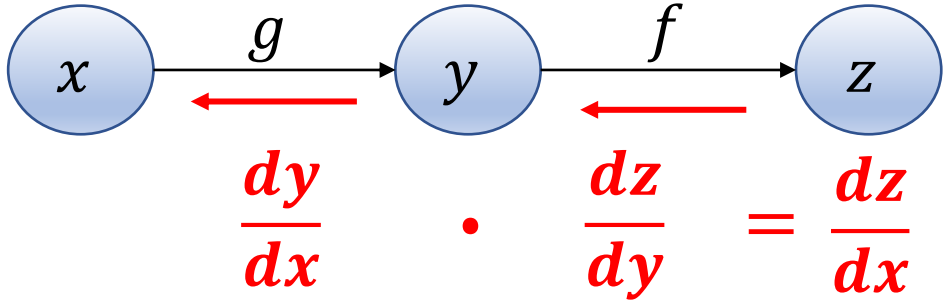


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

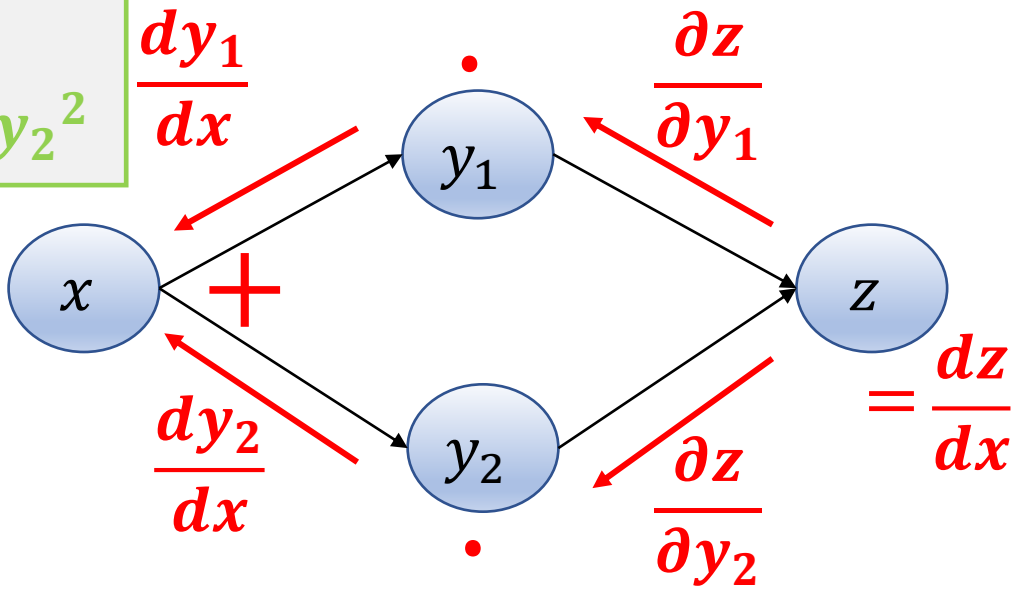


$$y_1 = 2x^2 + e^{2x}$$

$$y_2 = e^x$$

$$z = 3y_1 + y_2^2$$

$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

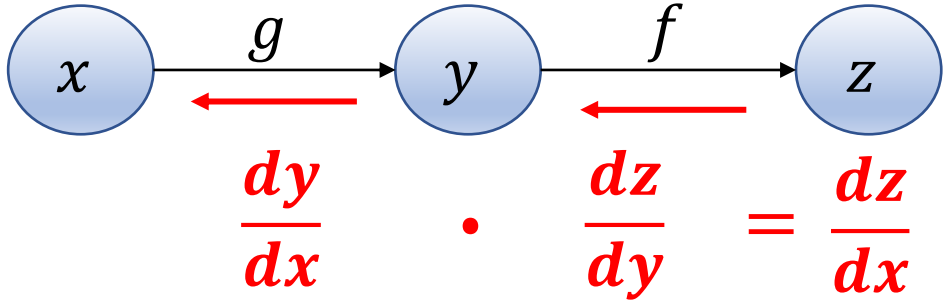


Chain rule reminder

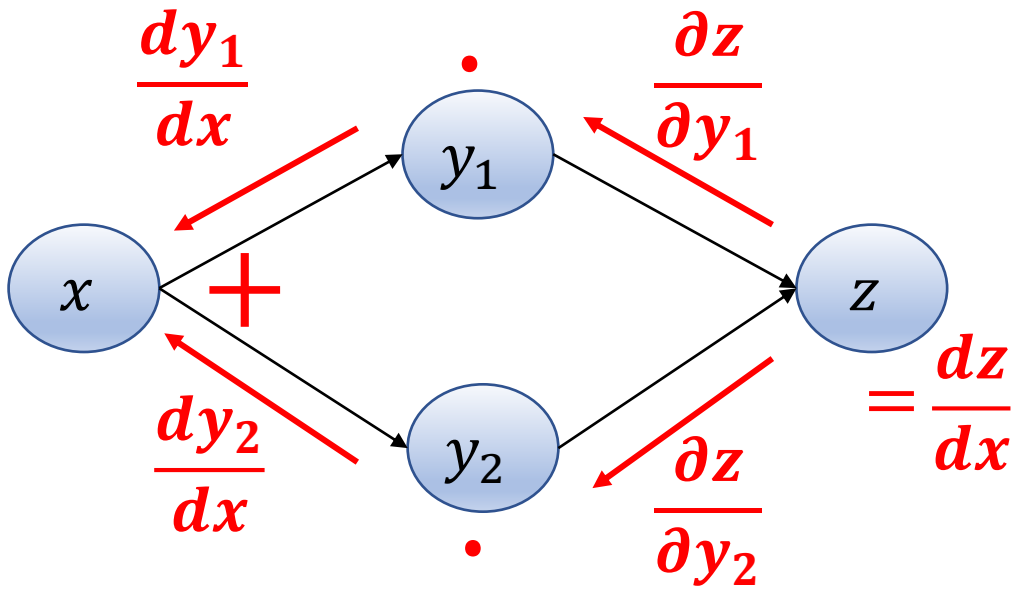
$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$



$$y_1 = 2x^2 + e^{2x}$$

$$y_2 = e^x$$

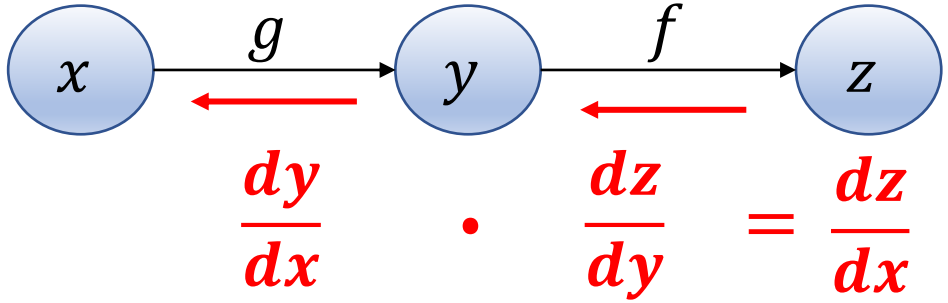
$$z = 3y_1 + y_2^2$$

Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

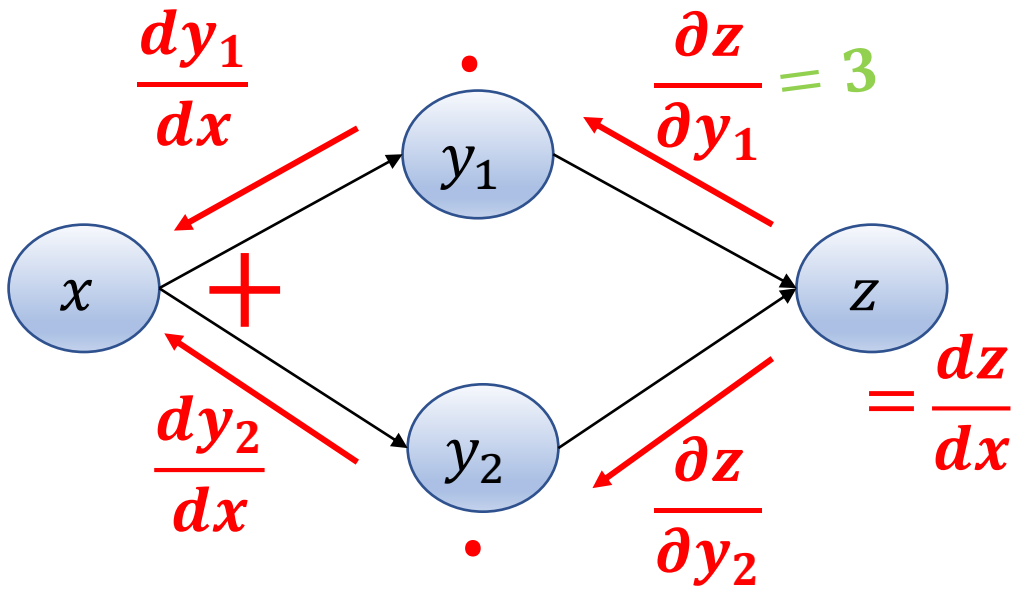


$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

$$y_1 = 2x^2 + e^{2x}$$

$$y_2 = e^x$$

$$z = 3y_1 + y_2^2$$

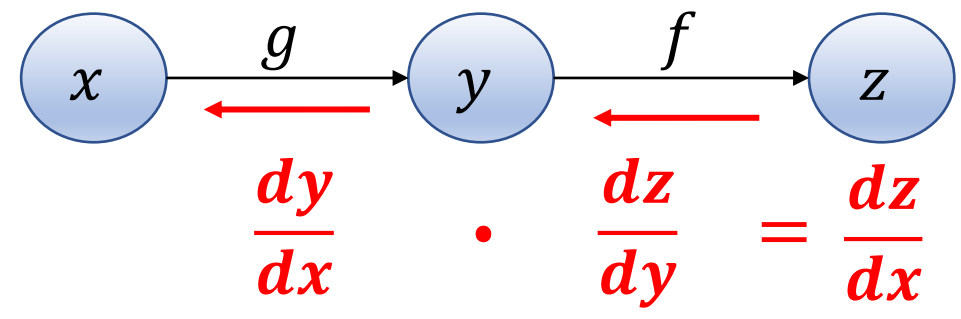


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



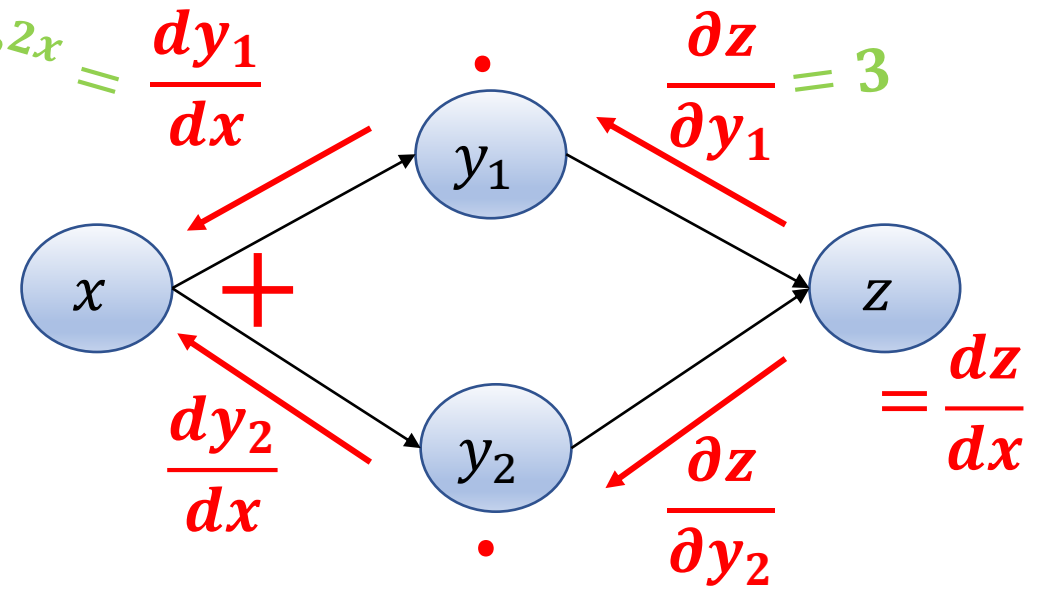
$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

$$y_1 = 2x^2 + e^{2x}$$

$$y_2 = e^x$$

$$z = 3y_1 + y_2^2$$

$$4x + 2e^{2x} = \frac{dy_1}{dx}$$



$$3 \quad 4x + 2e^{2x}$$

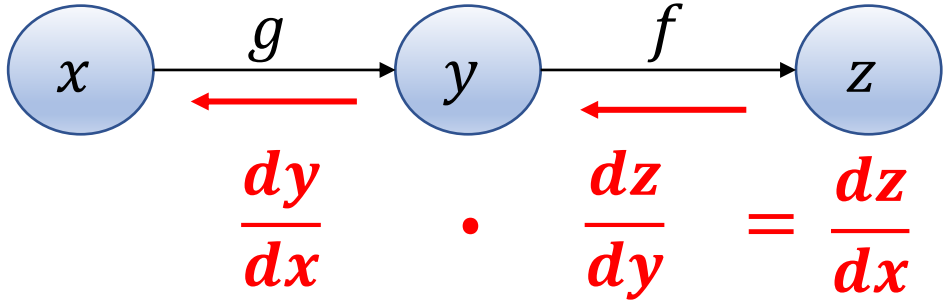


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

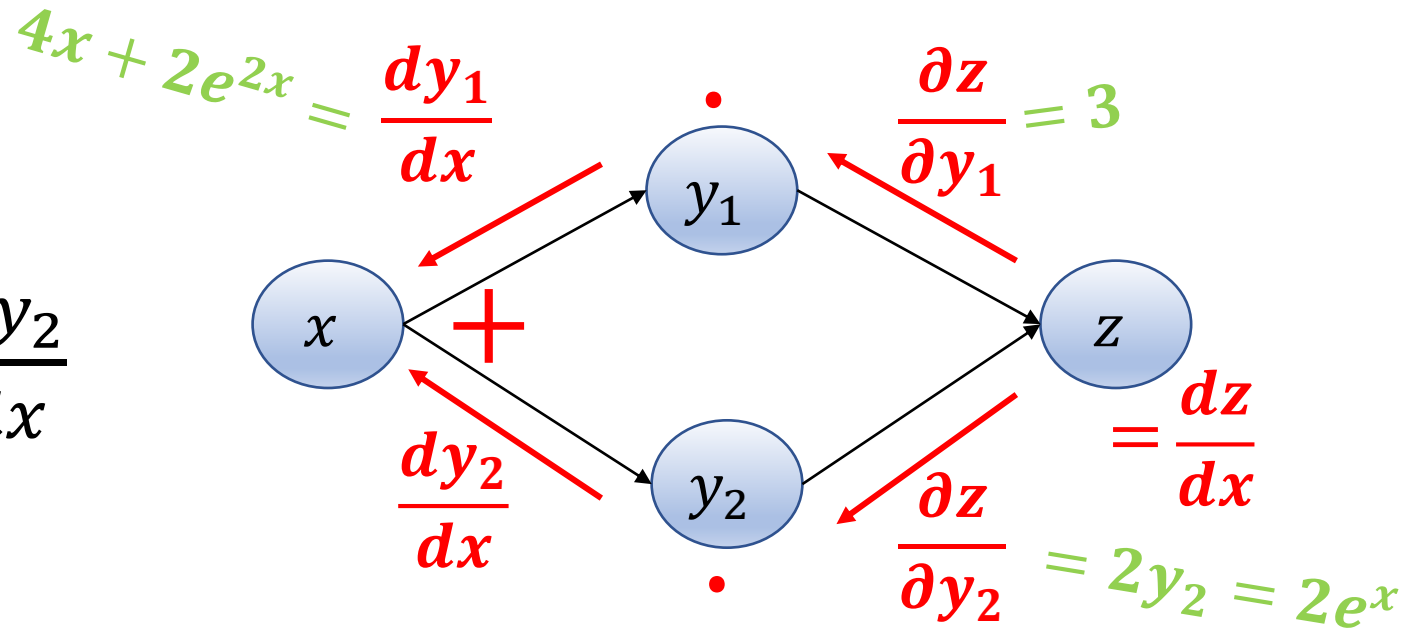


$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

$$y_1 = 2x^2 + e^{2x}$$

$$y_2 = e^x$$

$$z = 3y_1 + y_2^2$$

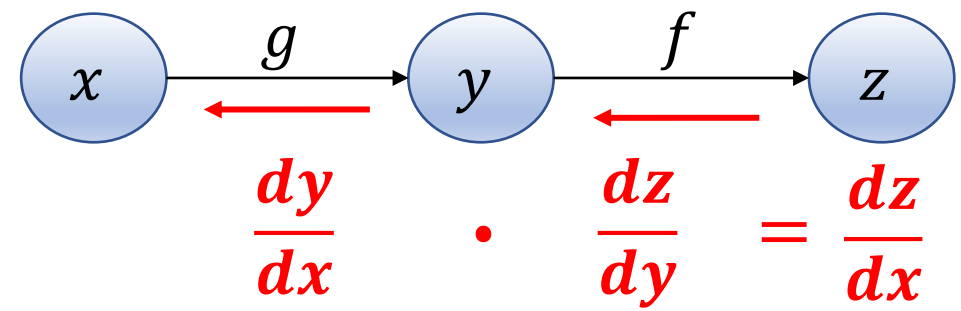


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

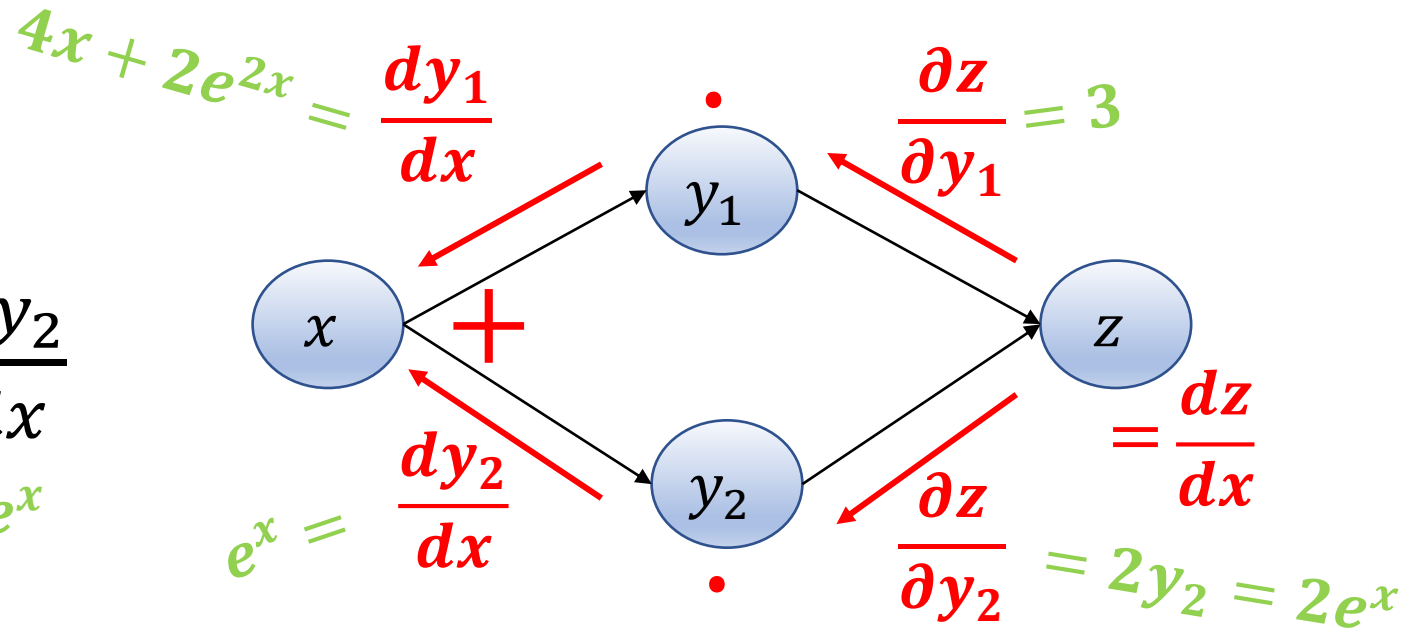


$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

$$y_1 = 2x^2 + e^{2x}$$

$$y_2 = e^x$$

$$z = 3y_1 + y_2^2$$

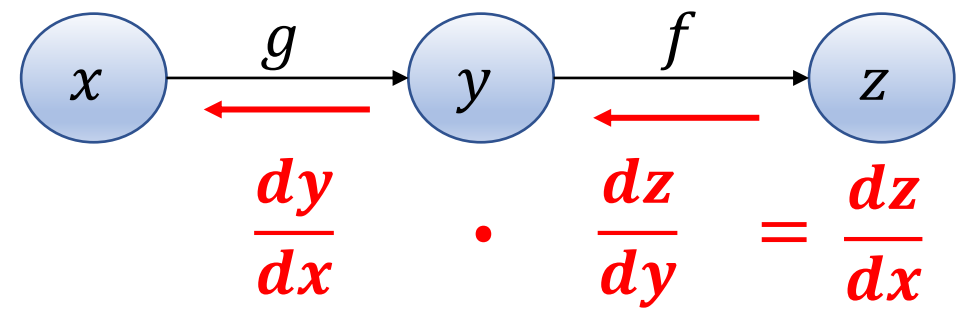


Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

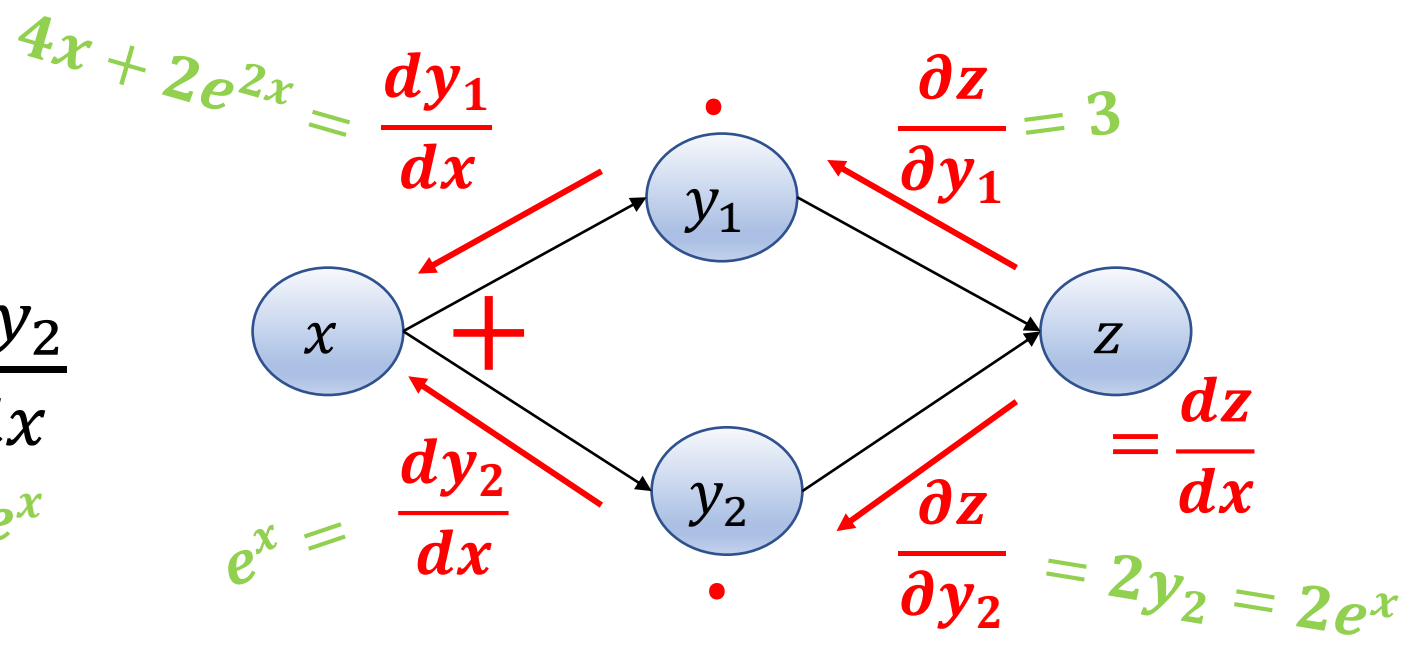


$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

$$y_1 = 2x^2 + e^{2x}$$

$$y_2 = e^x$$

$$z = 3y_1 + y_2^2$$



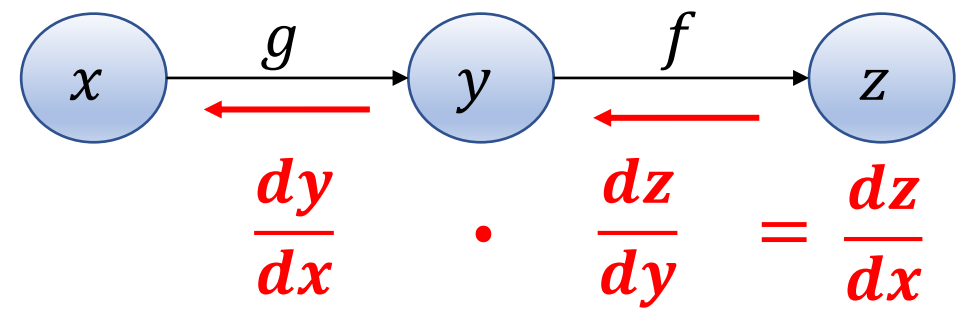
$$\frac{dz(y_1, y_2)}{dx} = 12x + 8e^{2x}$$

Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$

$$y = f(x), \quad z = g(y)$$

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$



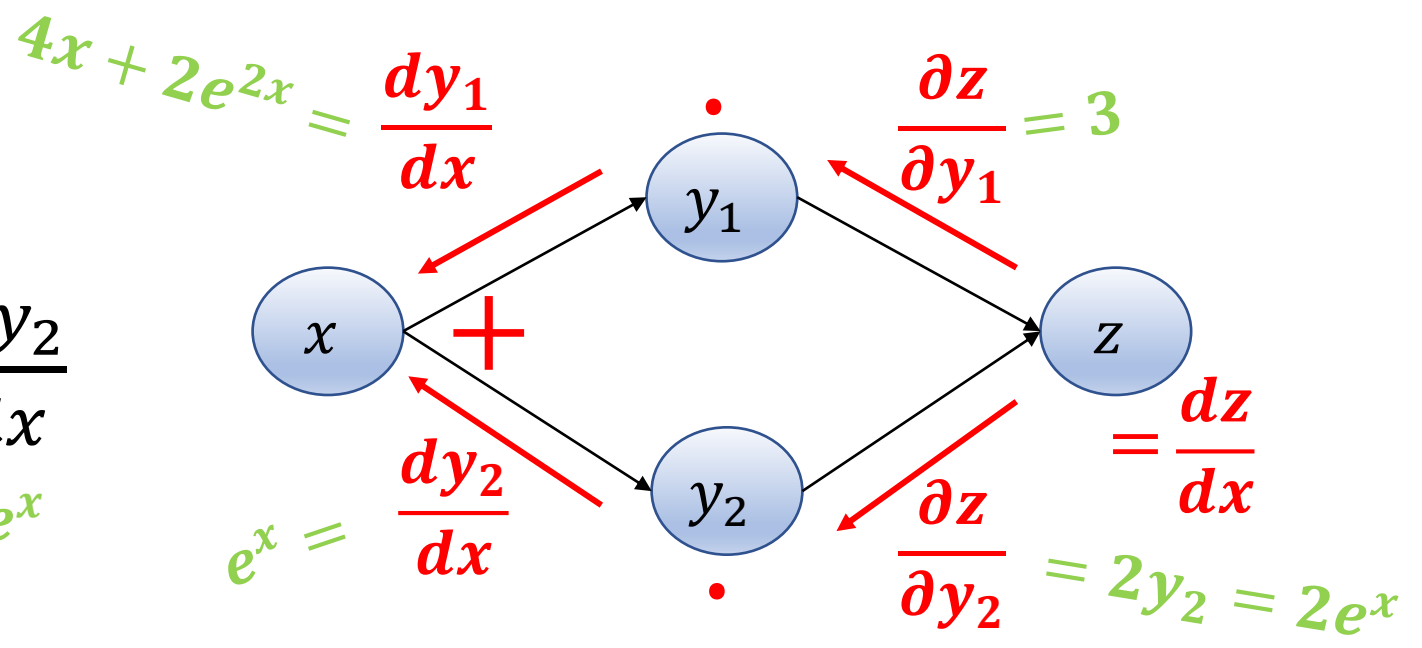
$$\frac{dz(y_1, y_2)}{dx} = \frac{\partial z}{\partial y_1} \frac{dy_1}{dx} + \frac{\partial z}{\partial y_2} \frac{dy_2}{dx}$$

$$y_1 = 2x^2 + e^{2x}$$

$$y_2 = e^x$$

$$z = 3y_1 + y_2^2$$

$$= 6x^2 + 4e^{2x}$$

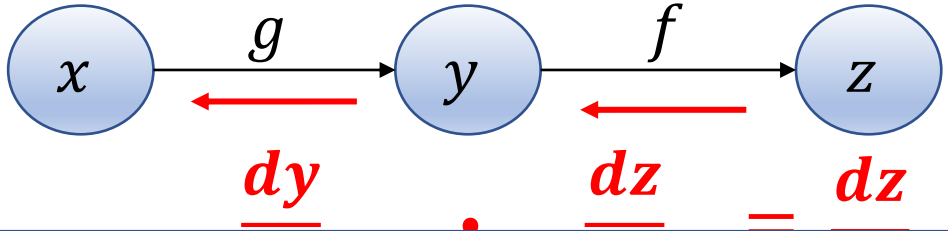


$$\frac{dz(y_1, y_2)}{dx} = 12x + 8e^{2x}$$



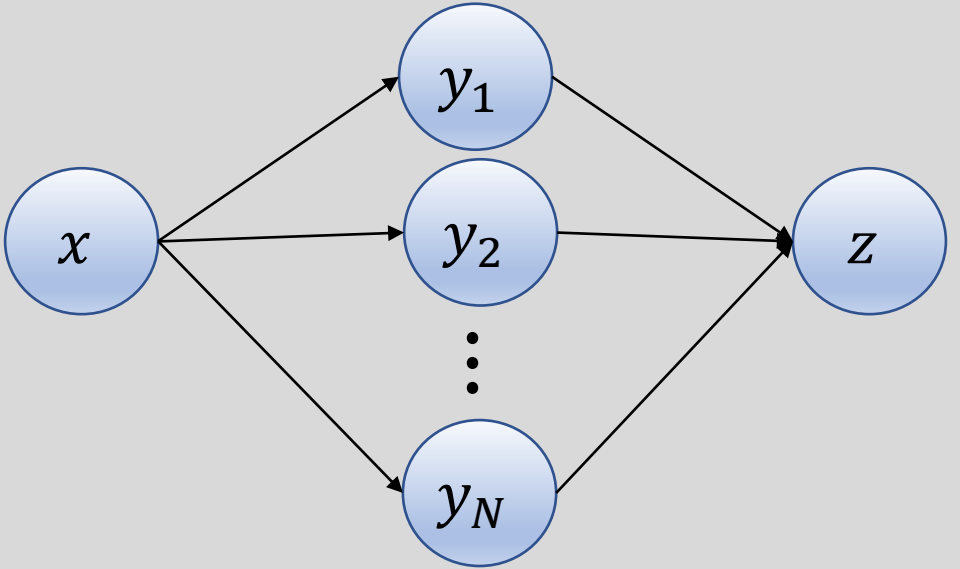
Chain rule reminder

$$f(g(x))' = f'(g(x)) \cdot g'(x)$$



Conclusion:

$$\frac{df(y_1, y_2 \dots y_N)}{dx} = \sum_n \frac{\partial f}{\partial y_n} \frac{dy_n}{dx}$$



$$z = 3y_1 + y_2^2 = 6x^2 + 4e^{2x}$$

$$\frac{dz(y_1, y_2)}{dx} = 12x + 8e^{2x}$$



Back Propagation - preliminaries

$$x_j^l = \sigma \left(\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j \right)$$

Back Propagation - preliminaries

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1}}_{z_j^l} + b_j \right)$$



Back Propagation - preliminaries

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1}}_{z_j^l} + b_j \right)$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^l}$$

Back Propagation - preliminaries

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1}}_{z_j^l} + b_j \right)$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^l} = \frac{\partial \mathcal{L}}{\partial x_j^l} \cdot \frac{\partial x_j^l}{\partial w_{ij}^l}$$

Back Propagation - preliminaries

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^l} = \frac{\partial \mathcal{L}}{\partial x_j^l} \cdot \frac{\partial x_j^l}{\partial w_{ij}^l}$$

Back Propagation - preliminaries

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^l} = \frac{\partial \mathcal{L}}{\partial x_j^l} \cdot \frac{\partial x_j^l}{\partial w_{ij}^l}$$

Easy!
 $x_i^{l-1} \cdot \sigma'(z_j^l)$

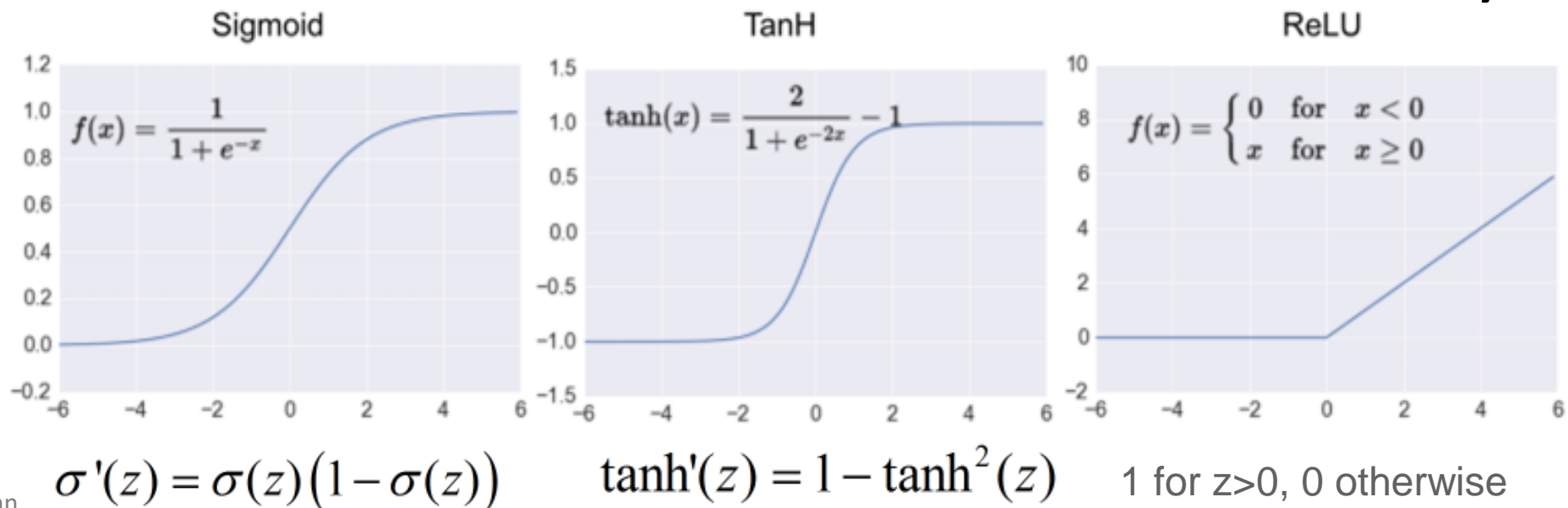
Back Propagation - preliminaries

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^l} = \frac{\partial \mathcal{L}}{\partial x_j^l} \cdot \frac{\partial x_j^l}{\partial w_{ij}^l}$$

Easy!
 $x_i^{l-1} \cdot \sigma'(z_j^l)$

Derivatives of common activations are easy!



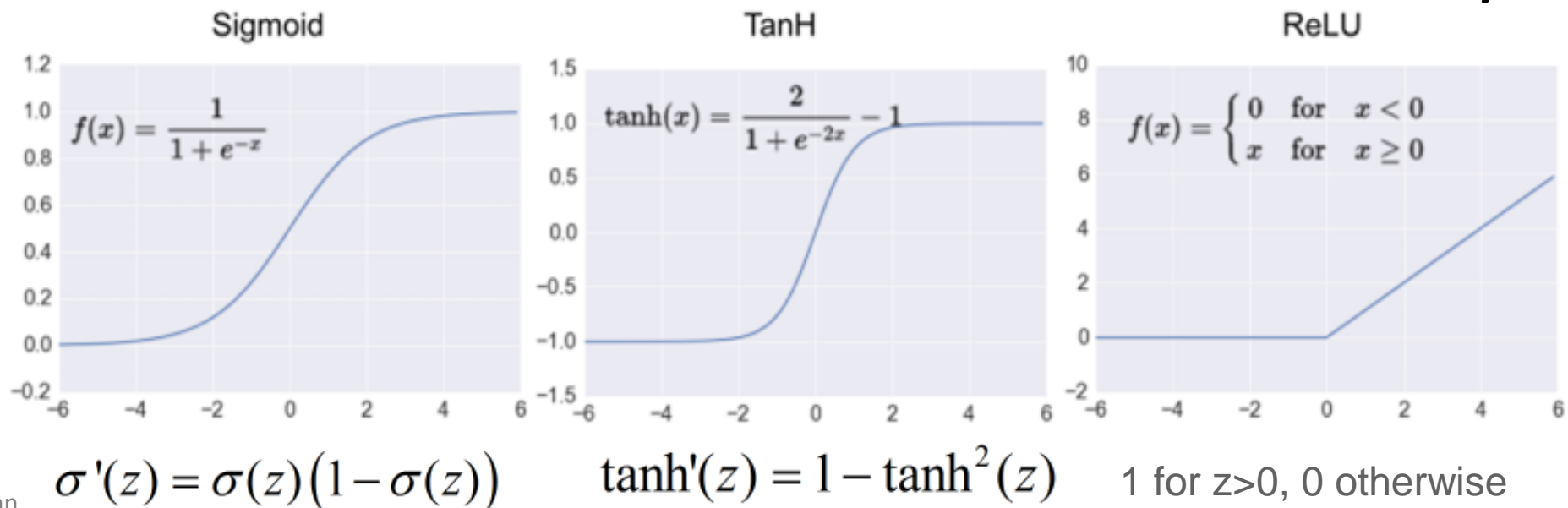
Back Propagation - preliminaries

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^l} = \frac{\partial \mathcal{L}}{\partial x_j^l} \cdot \frac{\partial x_j^l}{\partial w_{ij}^l}$$

Easy!
 $x_i^{l-1} \cdot \sigma'(z_j^l)$

Derivatives of common activations are easy!



Back Propagation - preliminaries

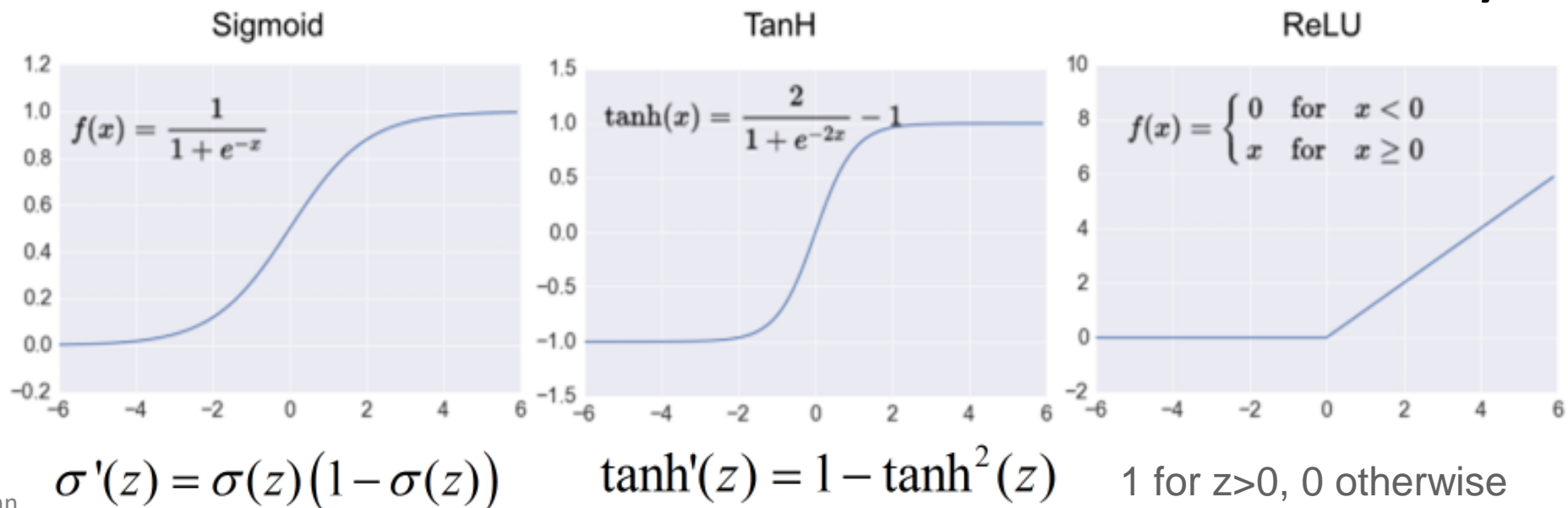
$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

$$\frac{\partial \mathcal{L}}{\partial w_{ij}^l} = \frac{\partial \mathcal{L}}{\partial x_j^l} \cdot \frac{\partial x_j^l}{\partial w_{ij}^l}$$

$\triangleq g_j^l$
Obtained by
backprop

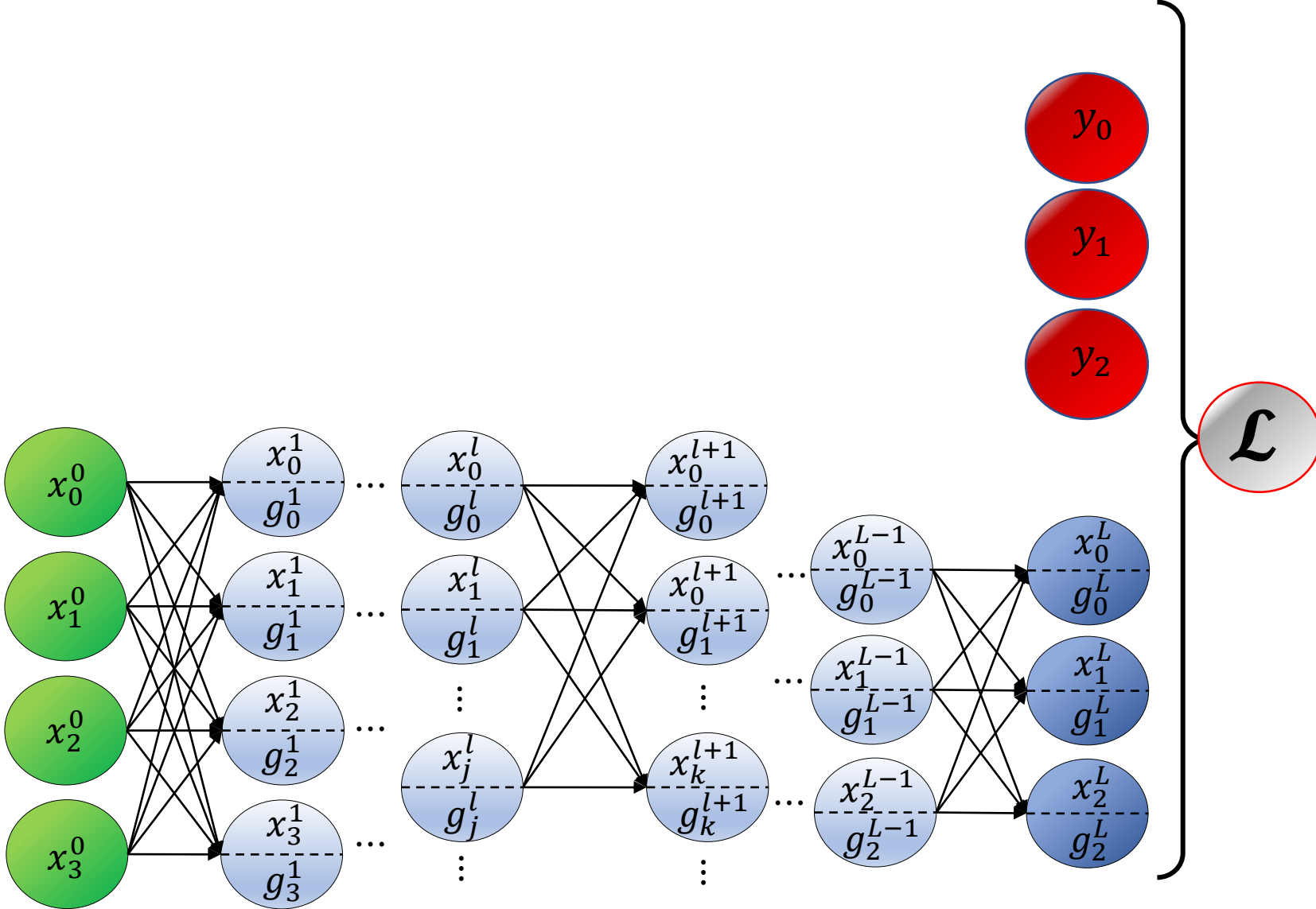
Easy!
 $x_i^{l-1} \cdot \sigma'(z_j^l)$

Derivatives of common activations are easy!



Back Propagation

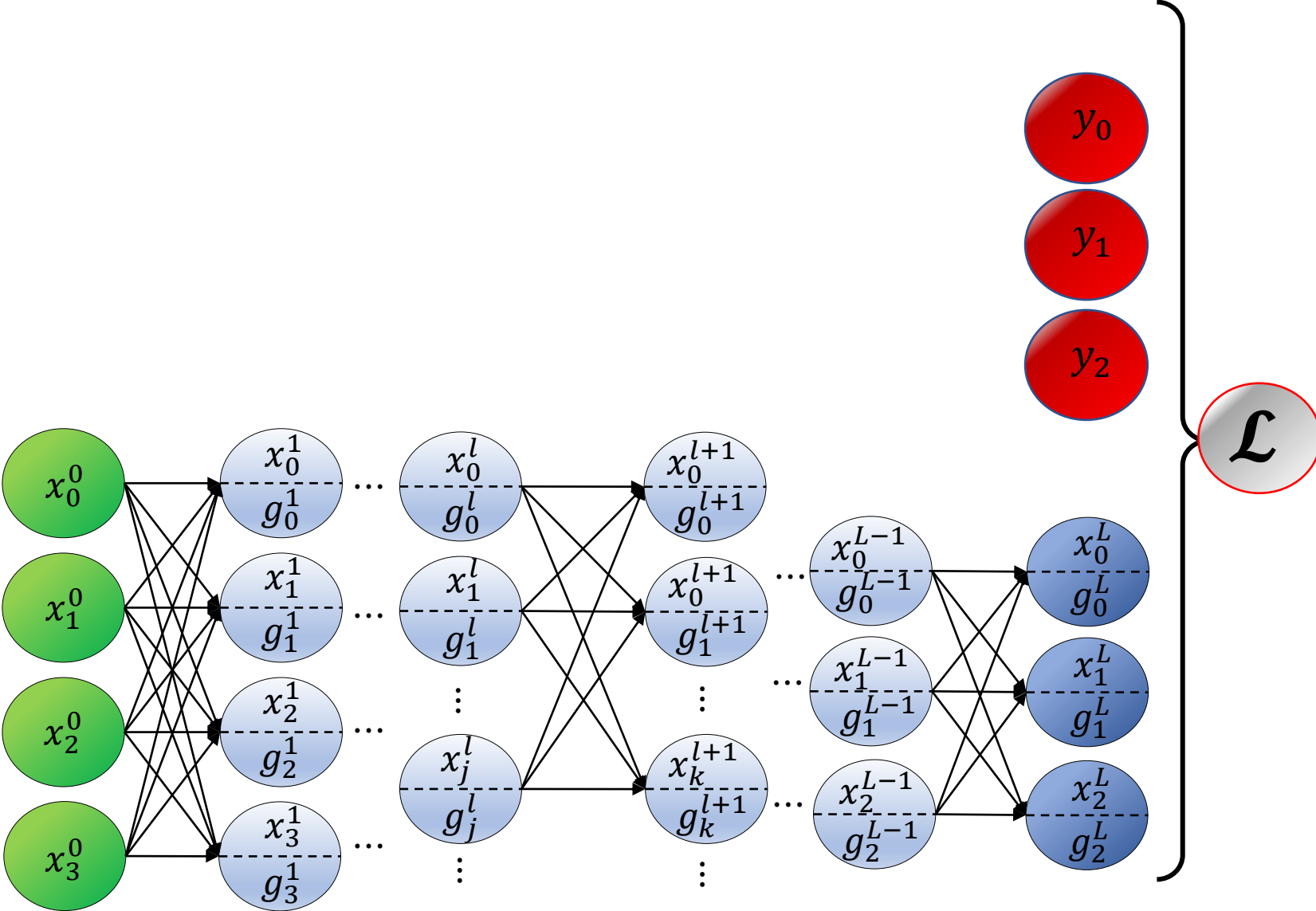
$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1}}_{z_j^l} + b_j \right)$$



Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

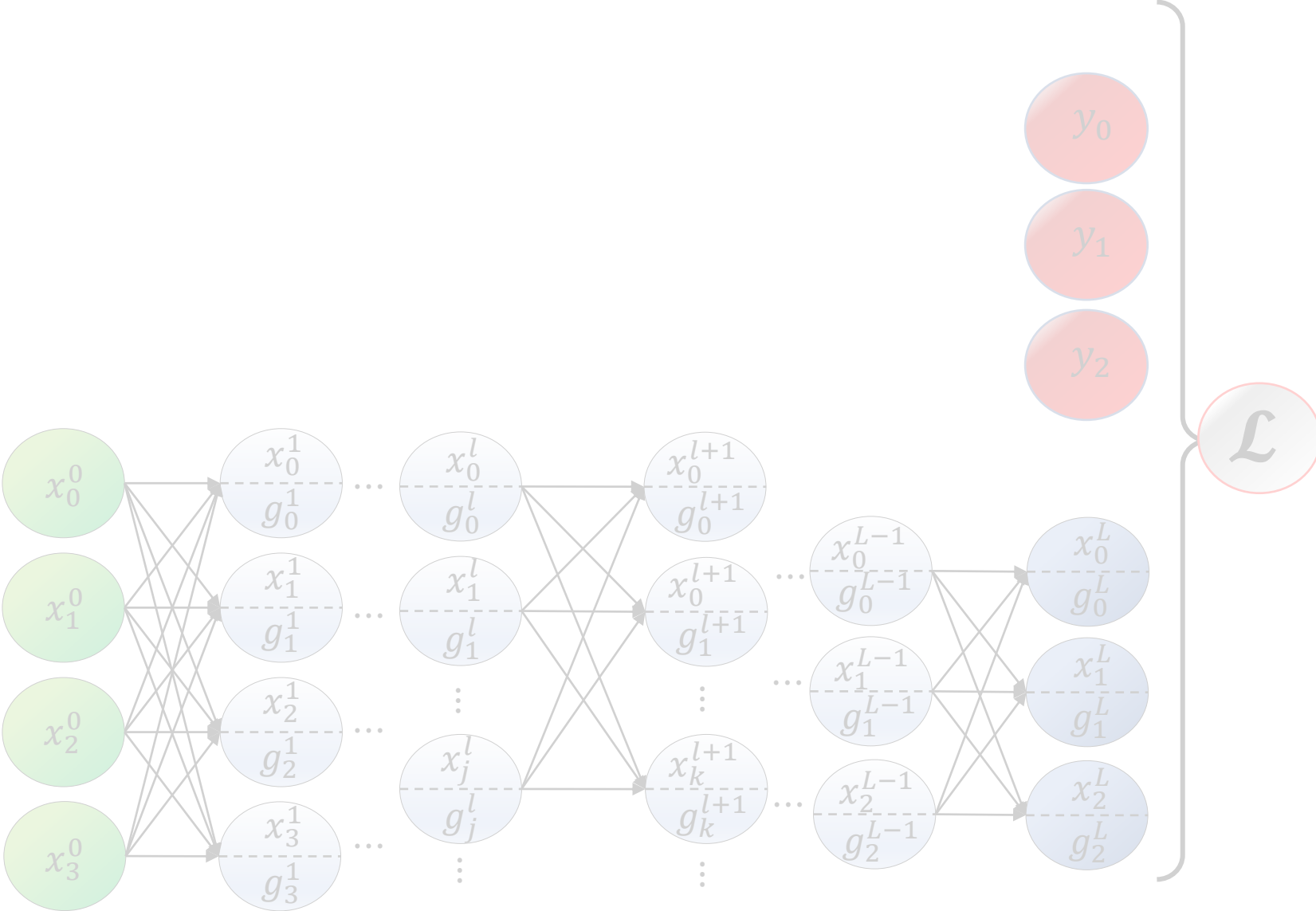
$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1}}_{z_j^l} + b_j \right)$$



Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

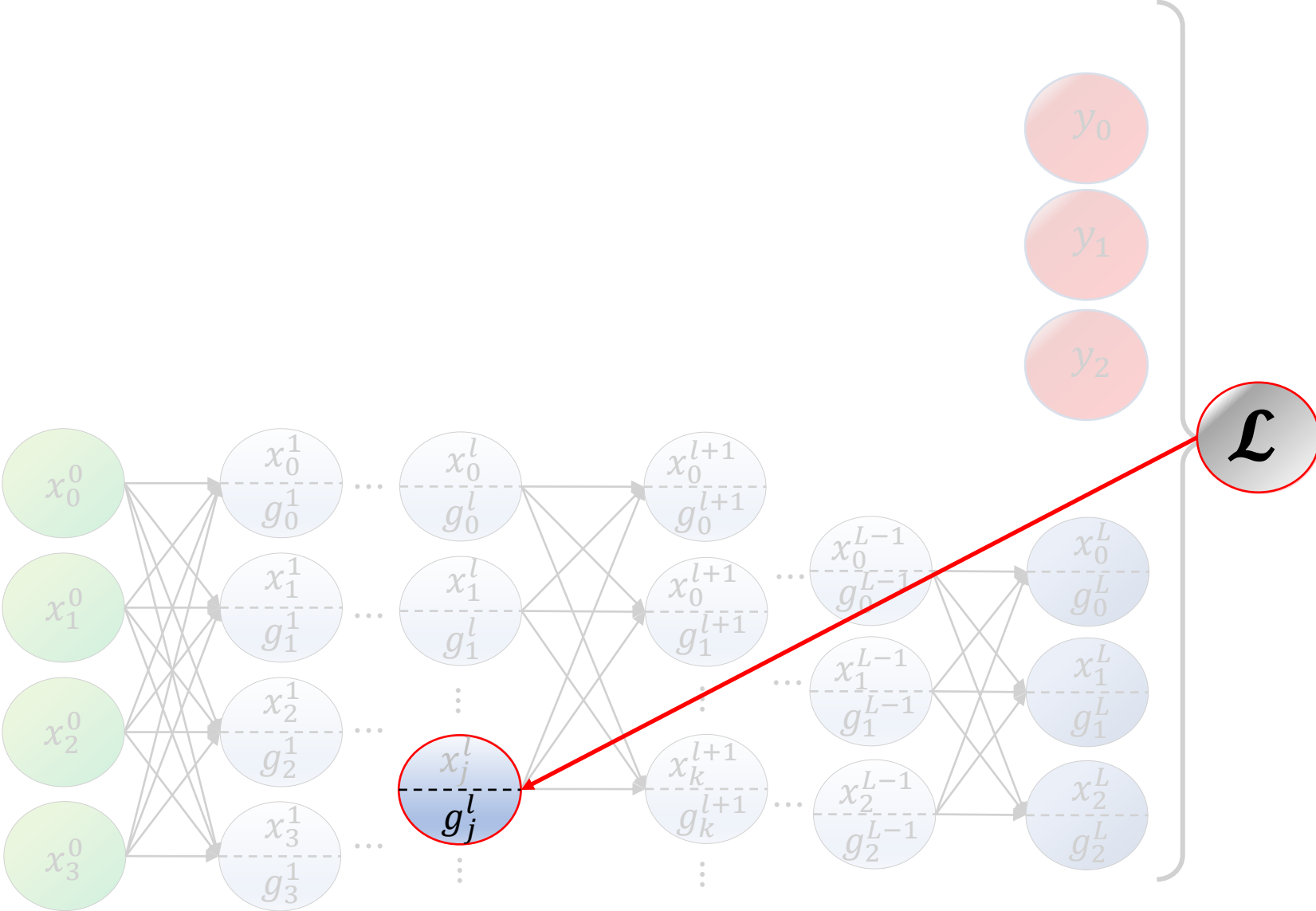
$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1}}_{z_j^l} + b_j \right)$$



Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

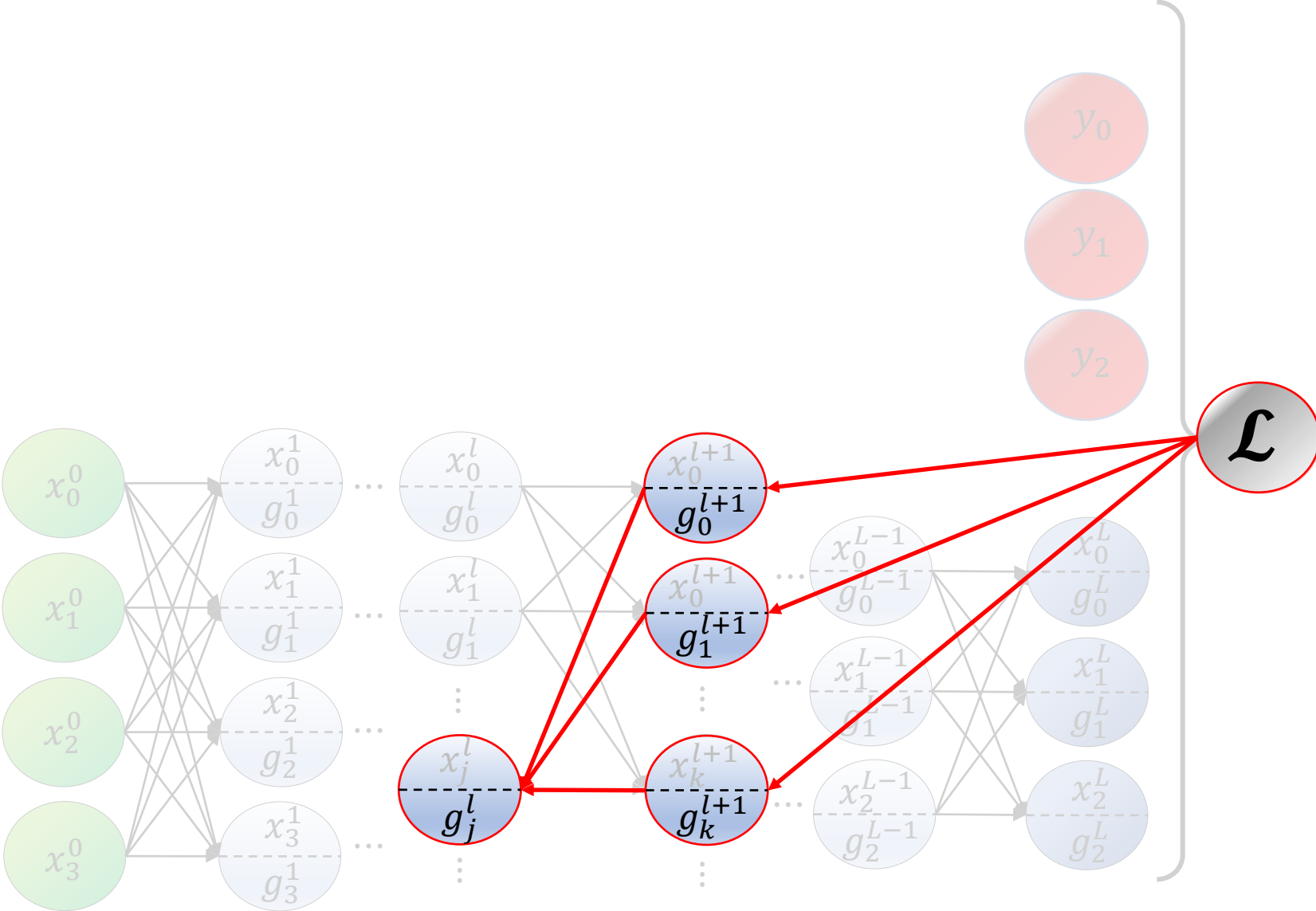
$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1}}_{z_j^l} + b_j \right)$$



Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1}}_{z_j^l} + b_j \right)$$

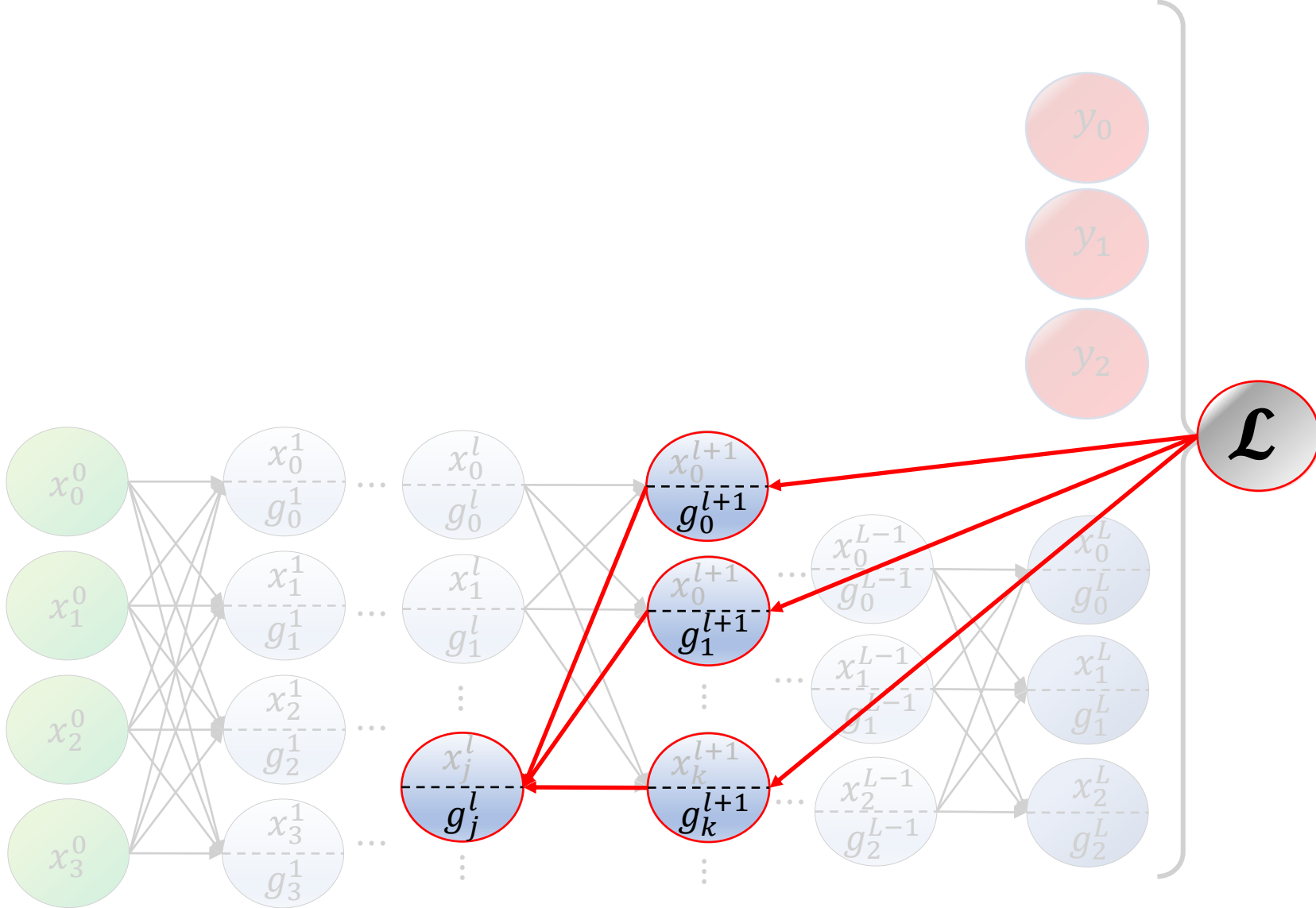


Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

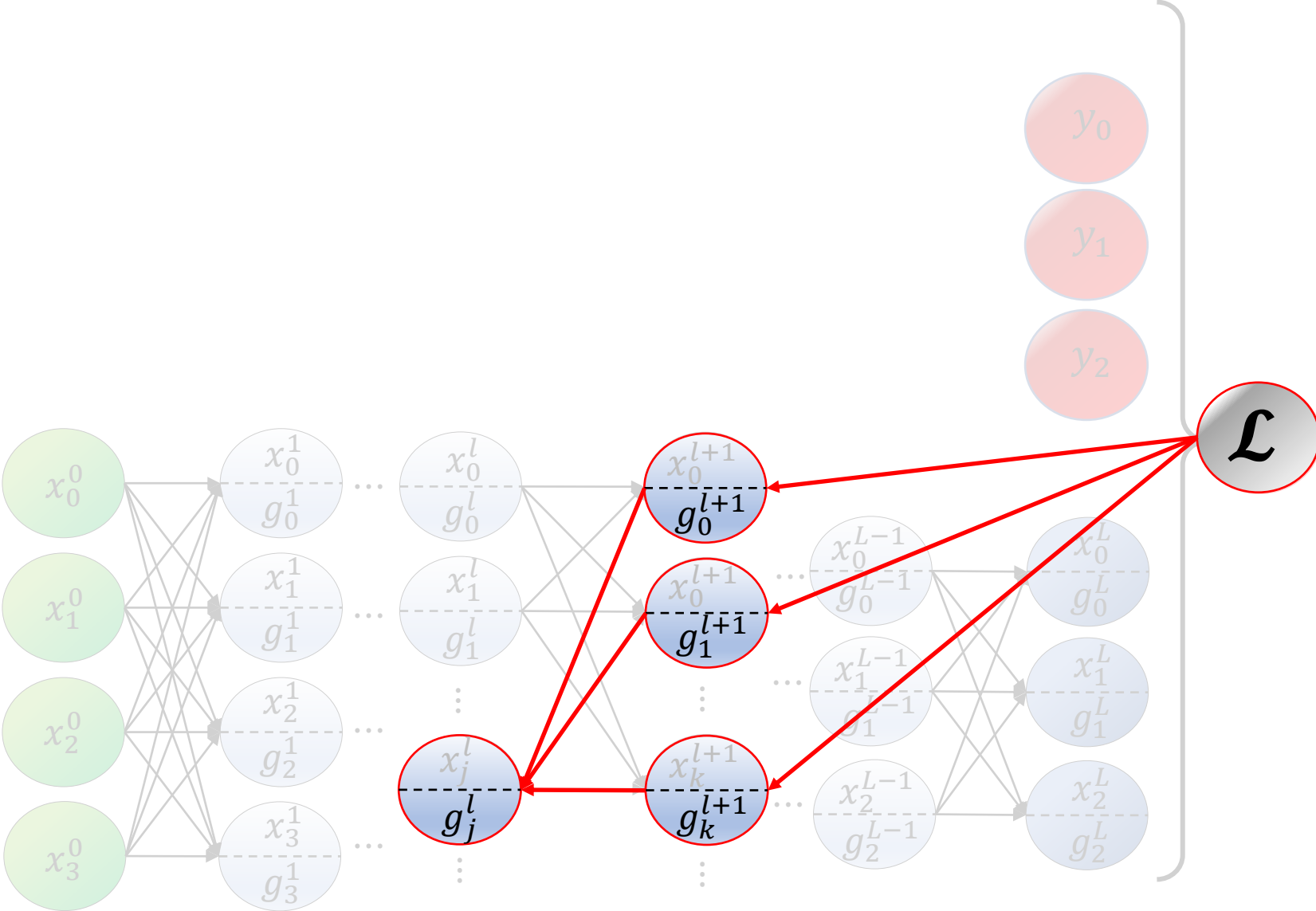


Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

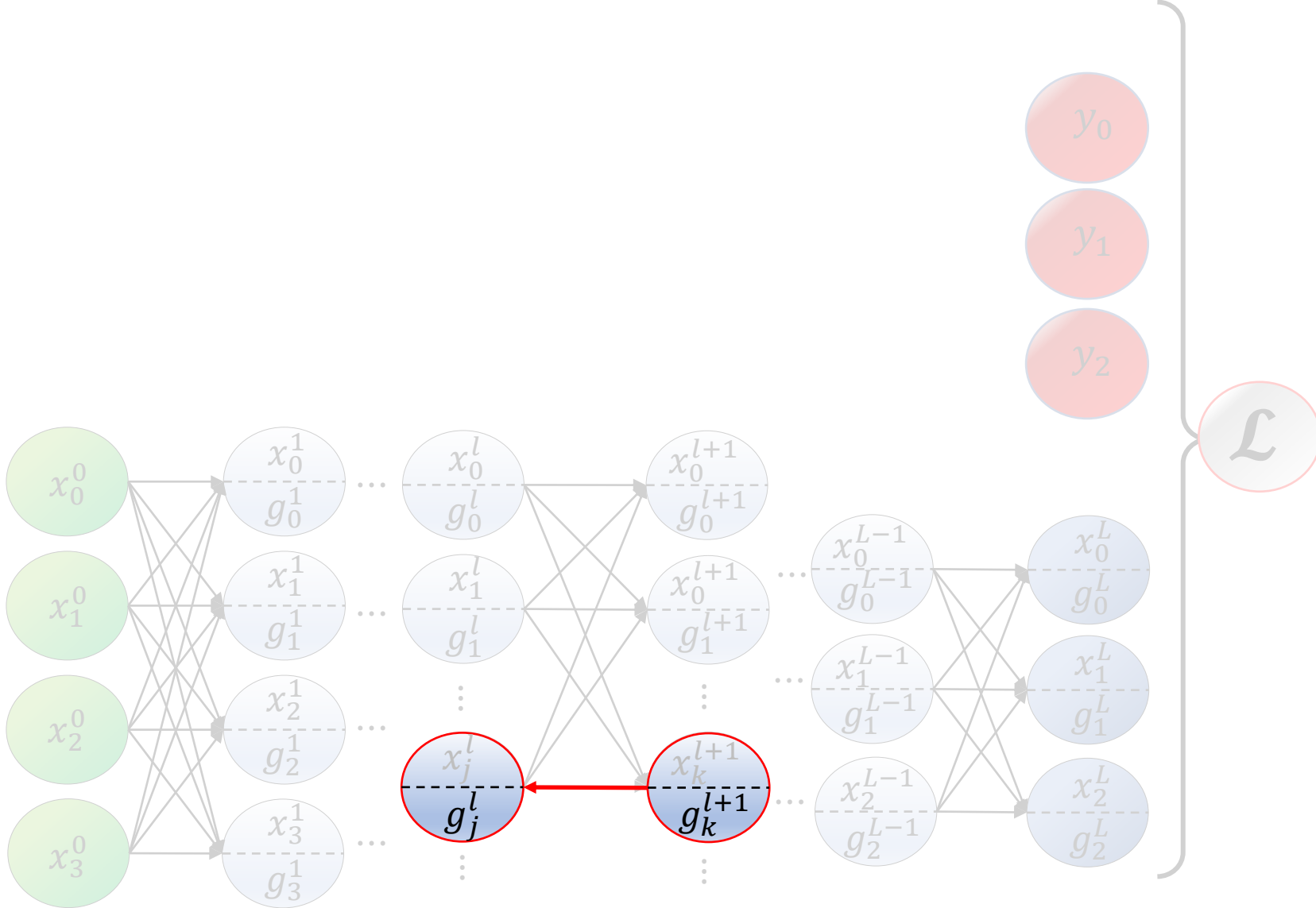


Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

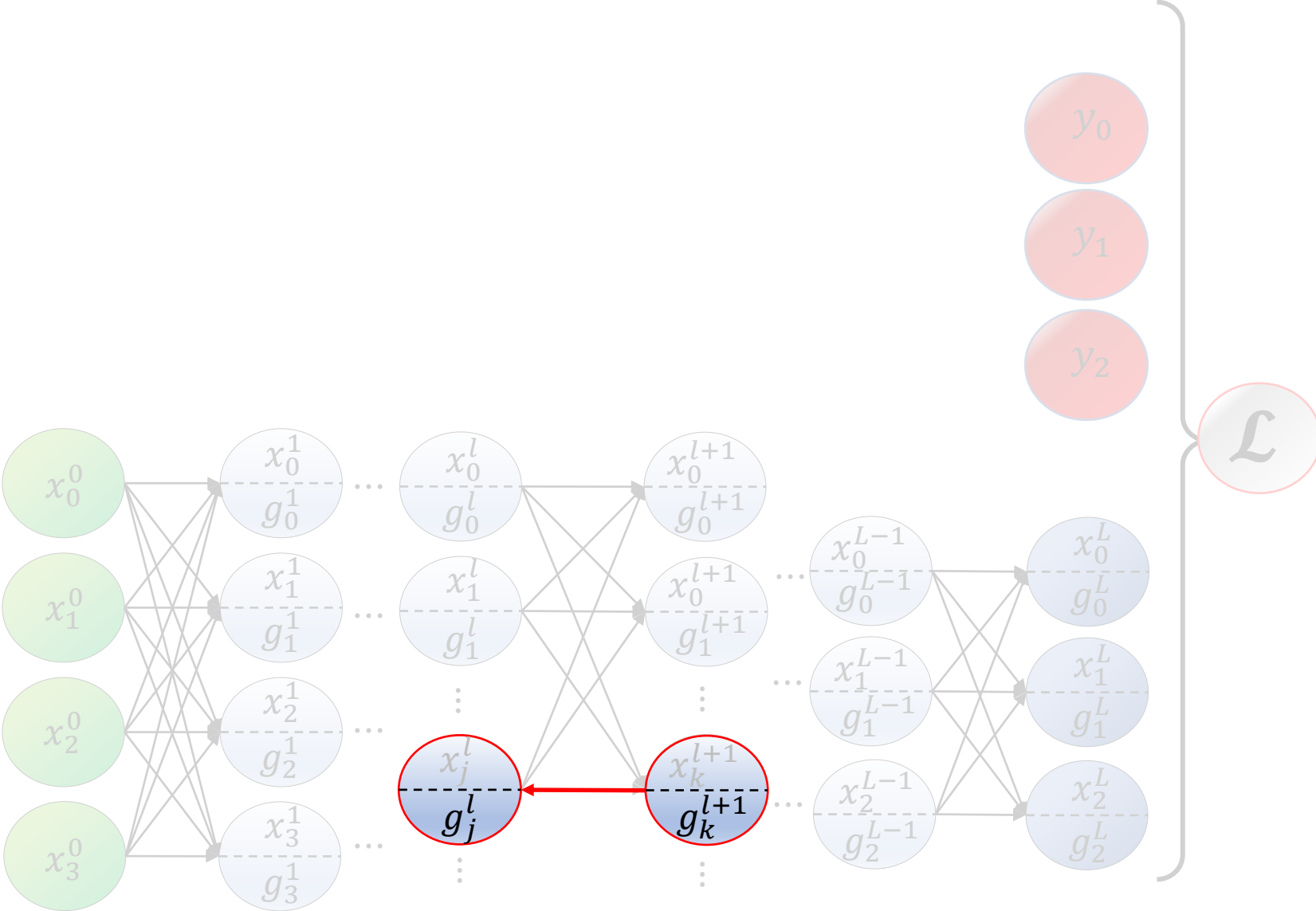


Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \underbrace{\frac{\partial x_k^{l+1}}{\partial x_j^l}}_{w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})}$$

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

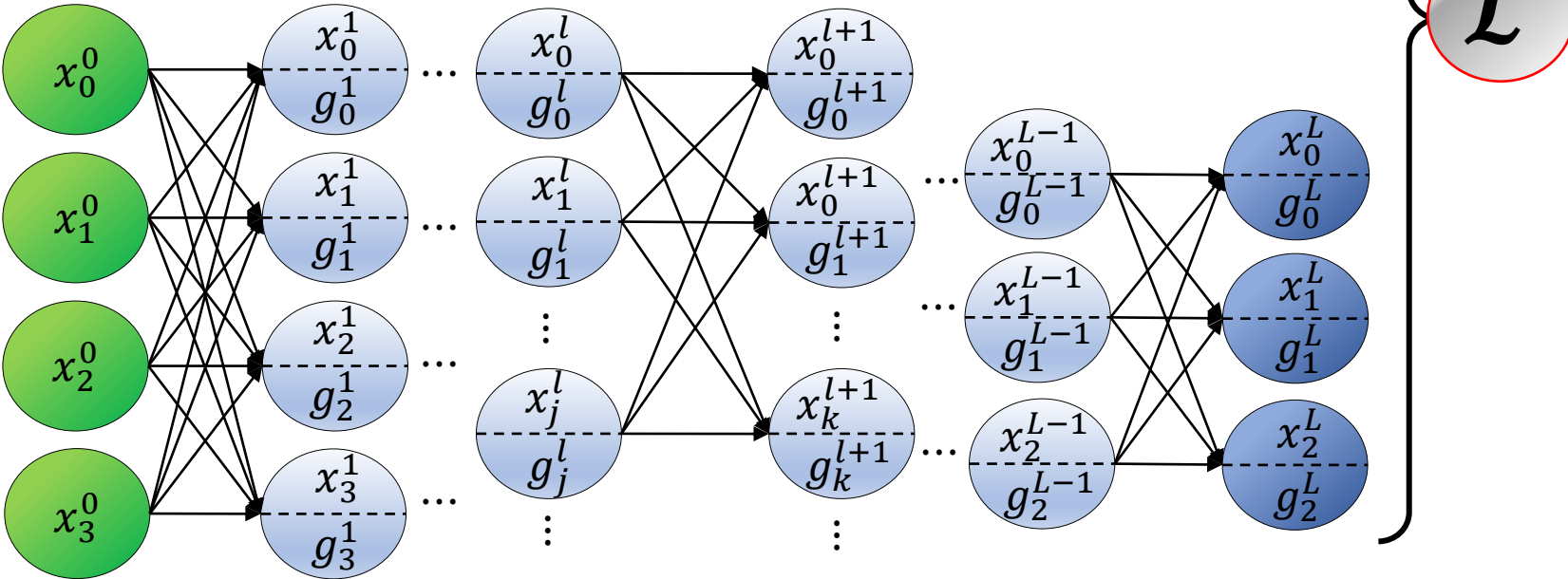


Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \underbrace{\frac{\partial x_k^{l+1}}{\partial x_j^l}}_{w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})}$$

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$



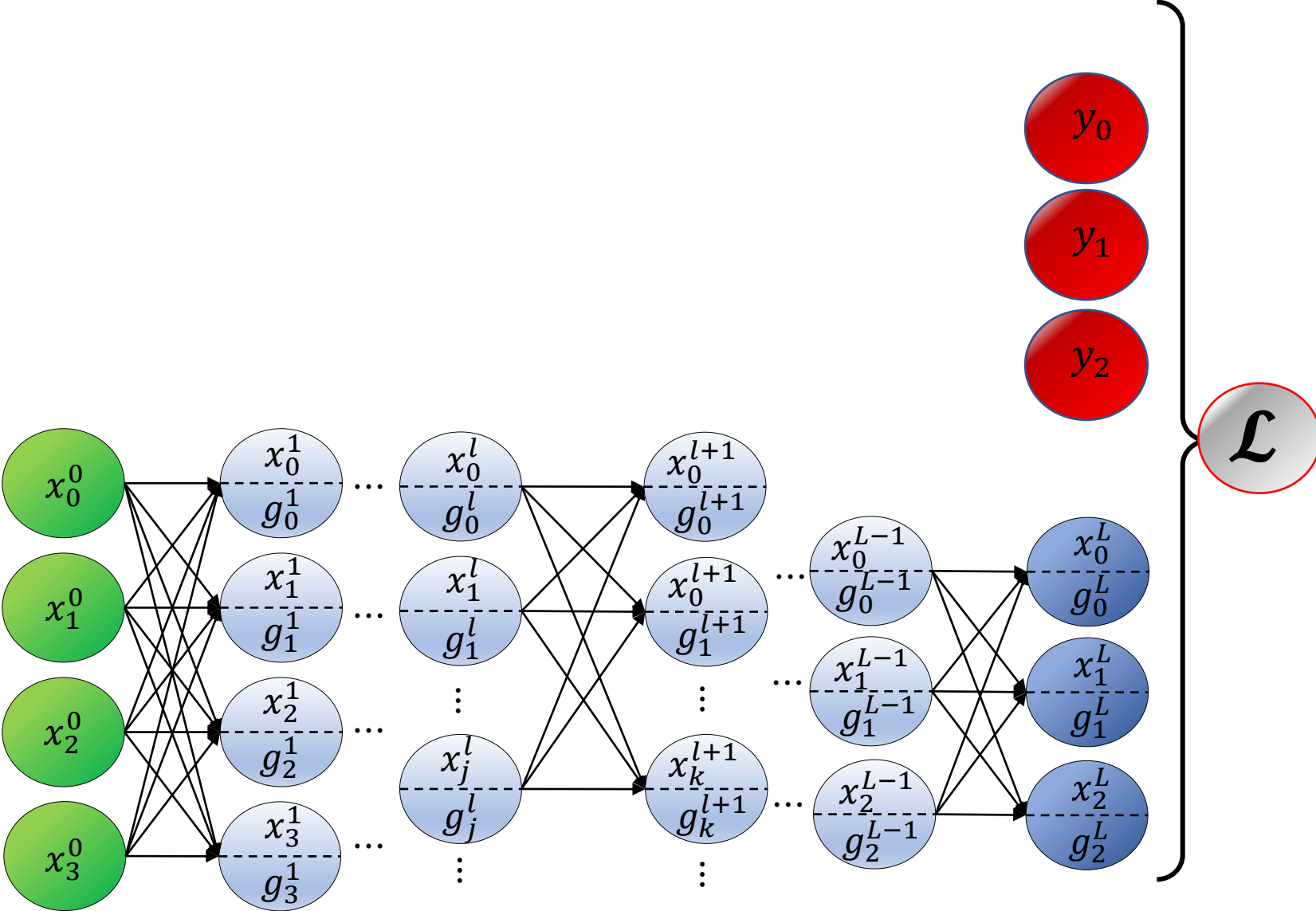
Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$\underbrace{\hspace{10em}}_{w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})}$

$$x_j^l = \sigma \left(\underbrace{\sum_i w_{ij}^l \cdot x_i^{l-1} + b_j}_{z_j^l} \right)$$

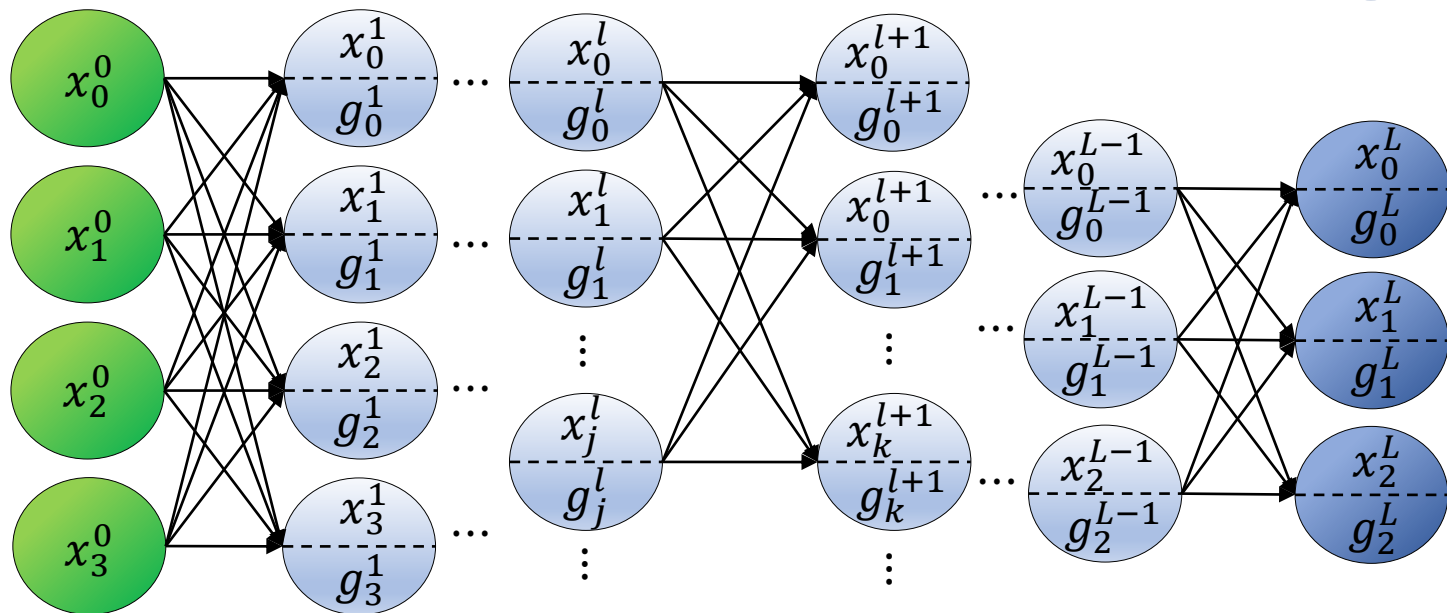
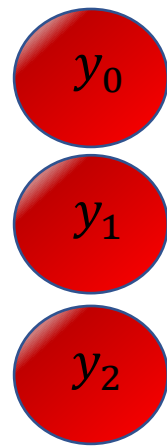
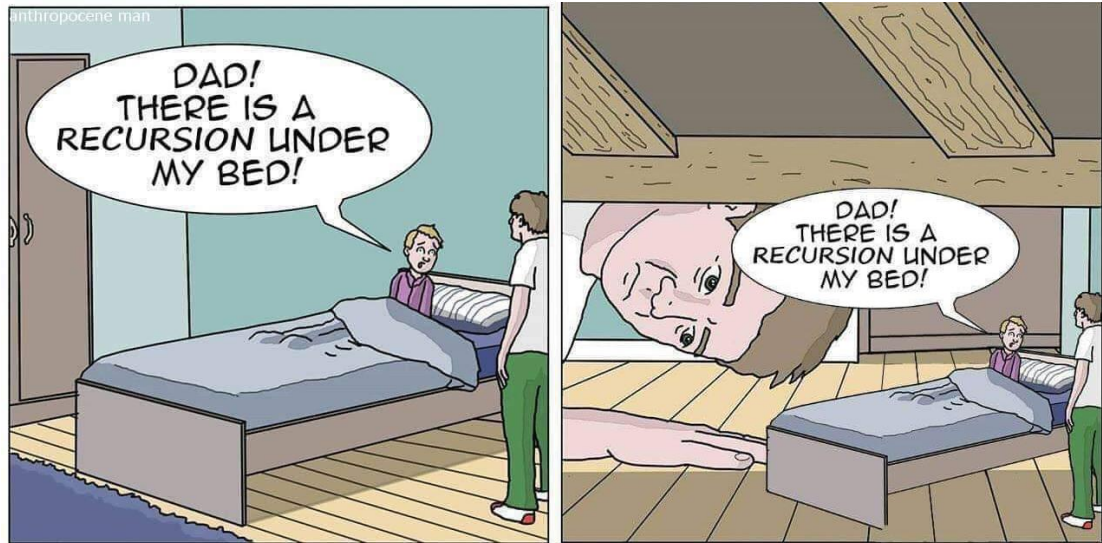


Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$

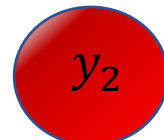
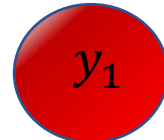
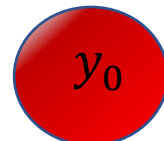
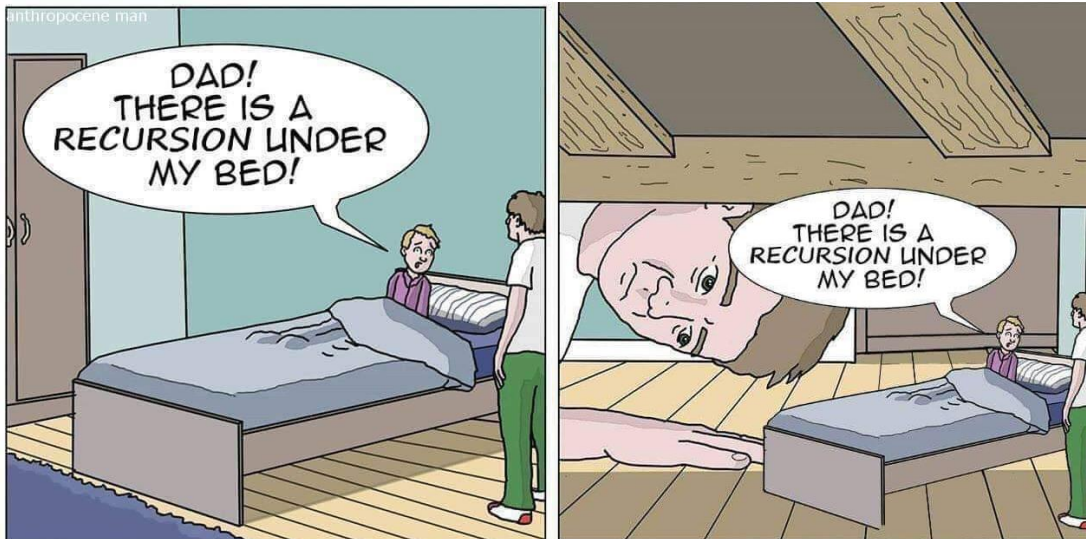


Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

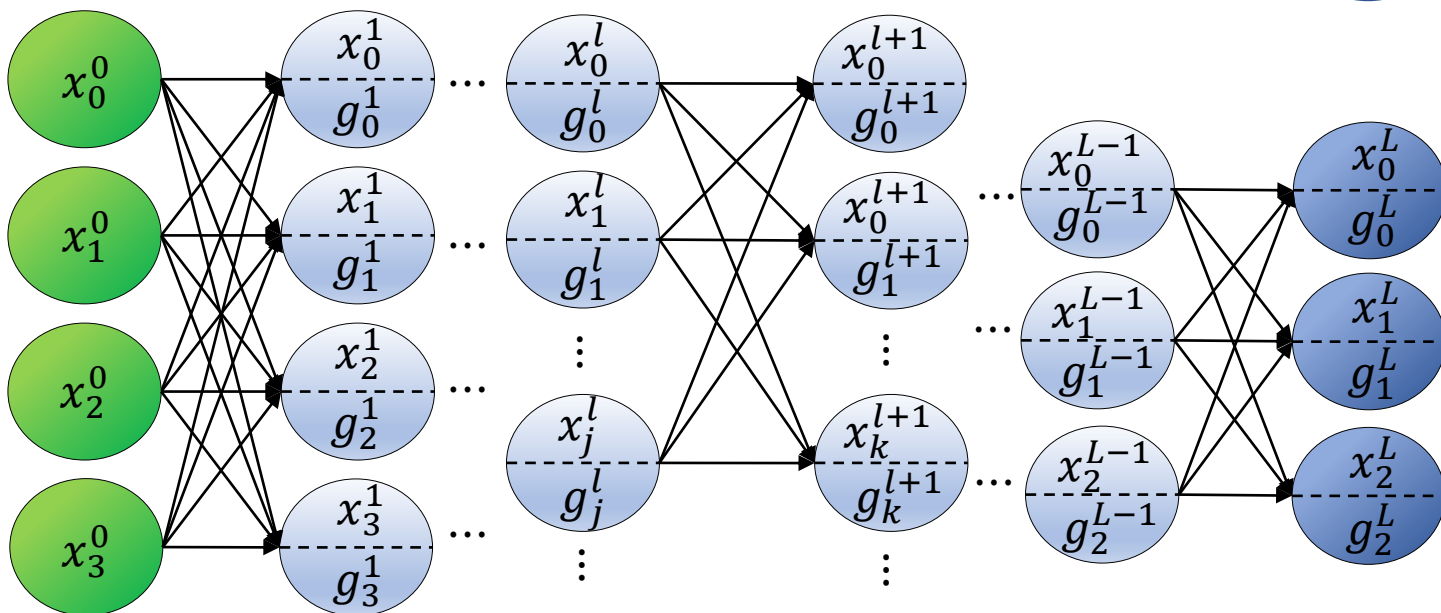
$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$



Stopping criterion:

$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$

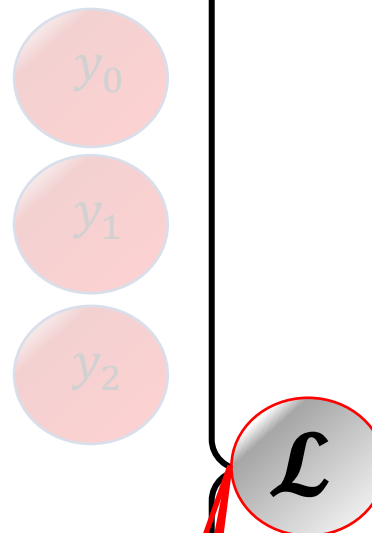
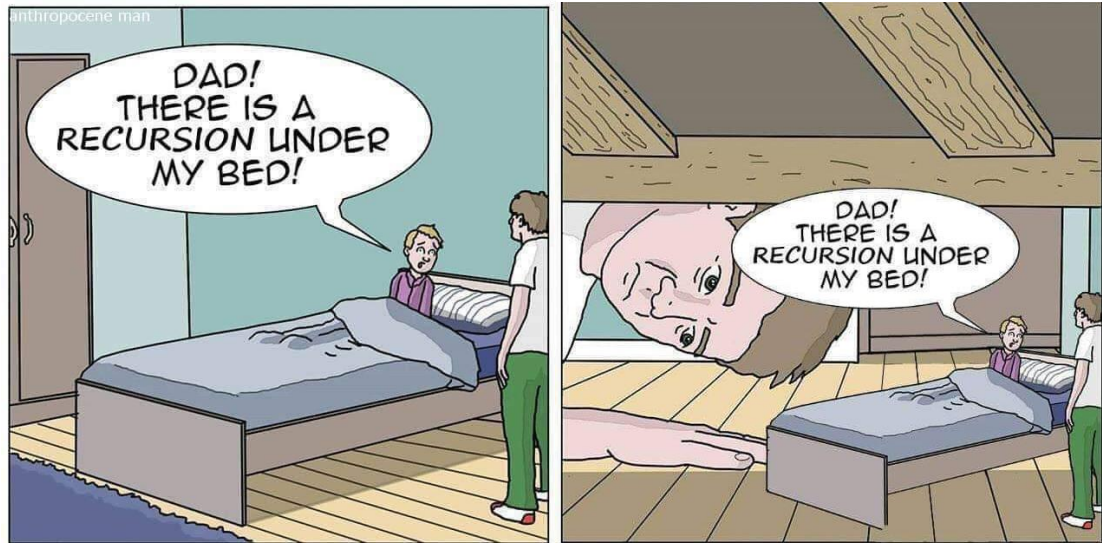


Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

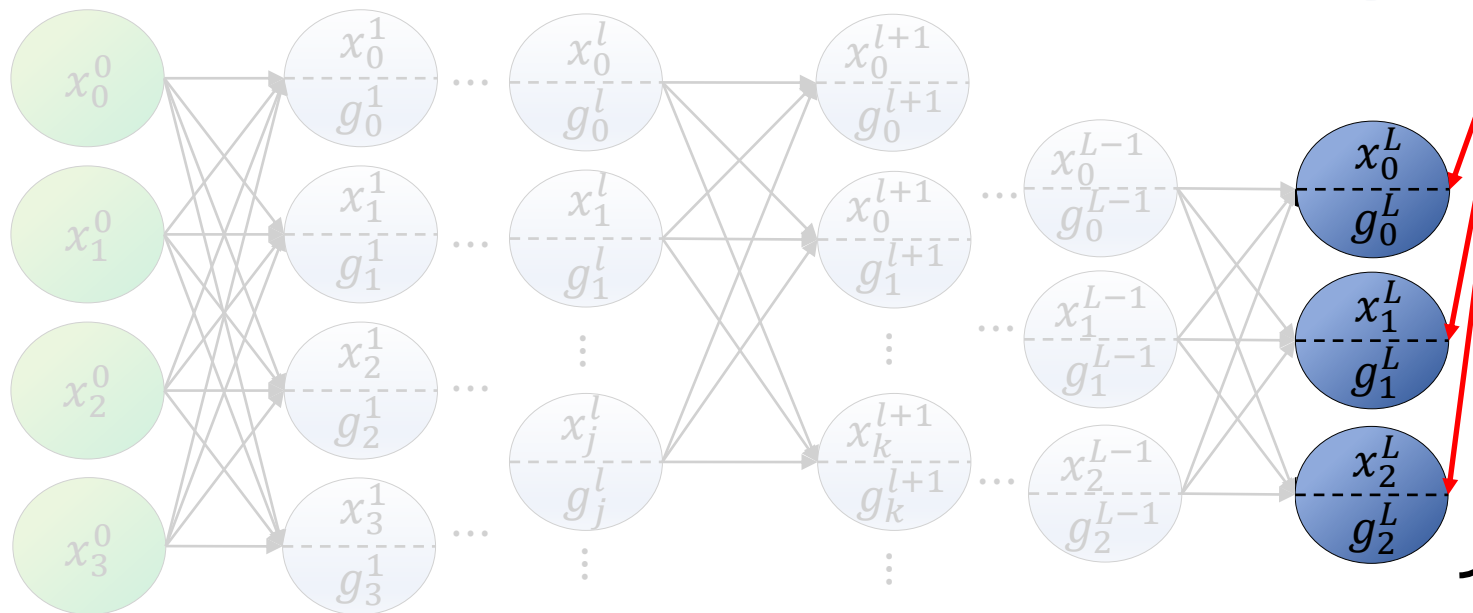
$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$



Stopping criterion:

$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$

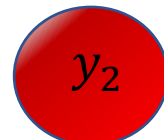
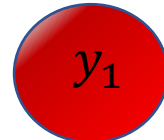
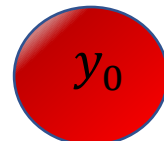
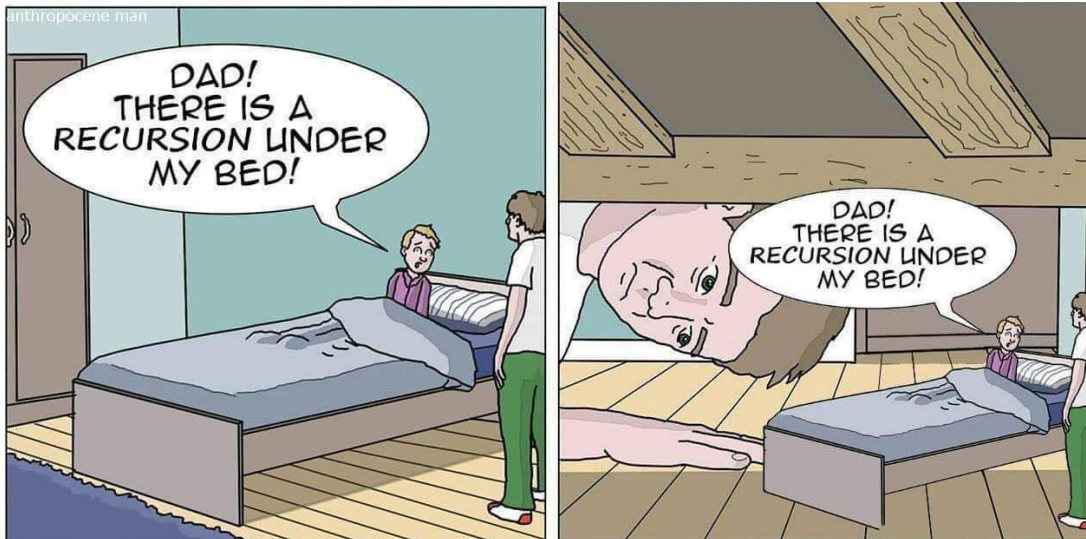


Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

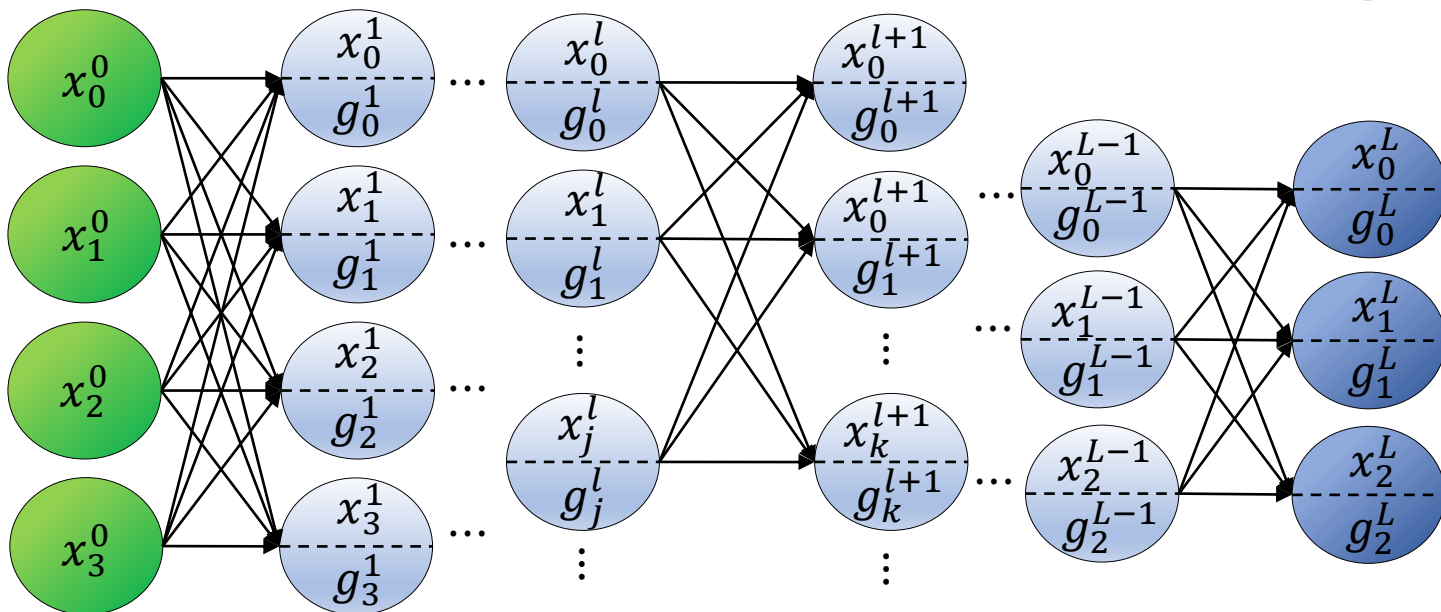
$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$



Stopping criterion:

$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$



Back Propagation

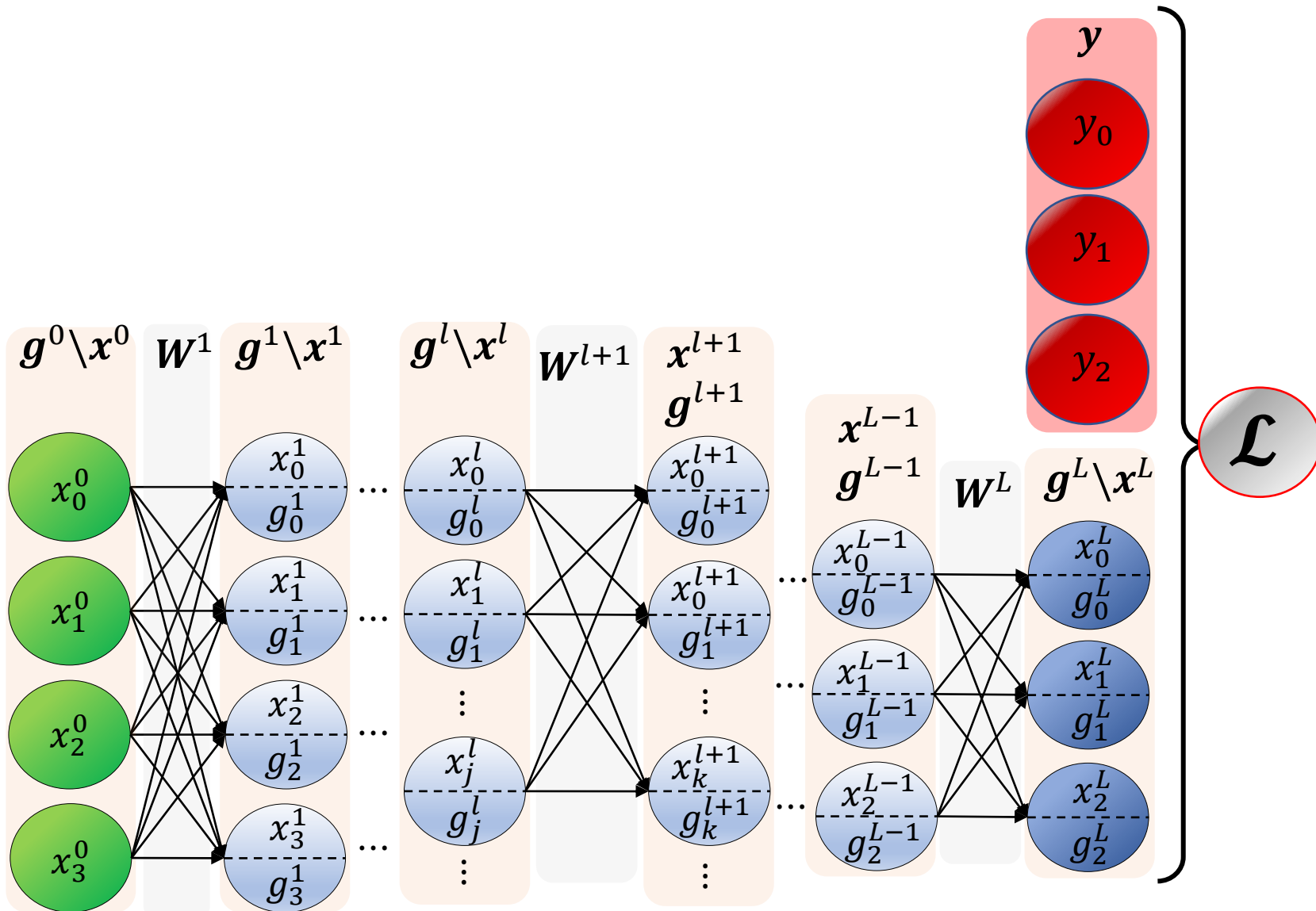
$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$

Stopping criterion:

$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$



Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

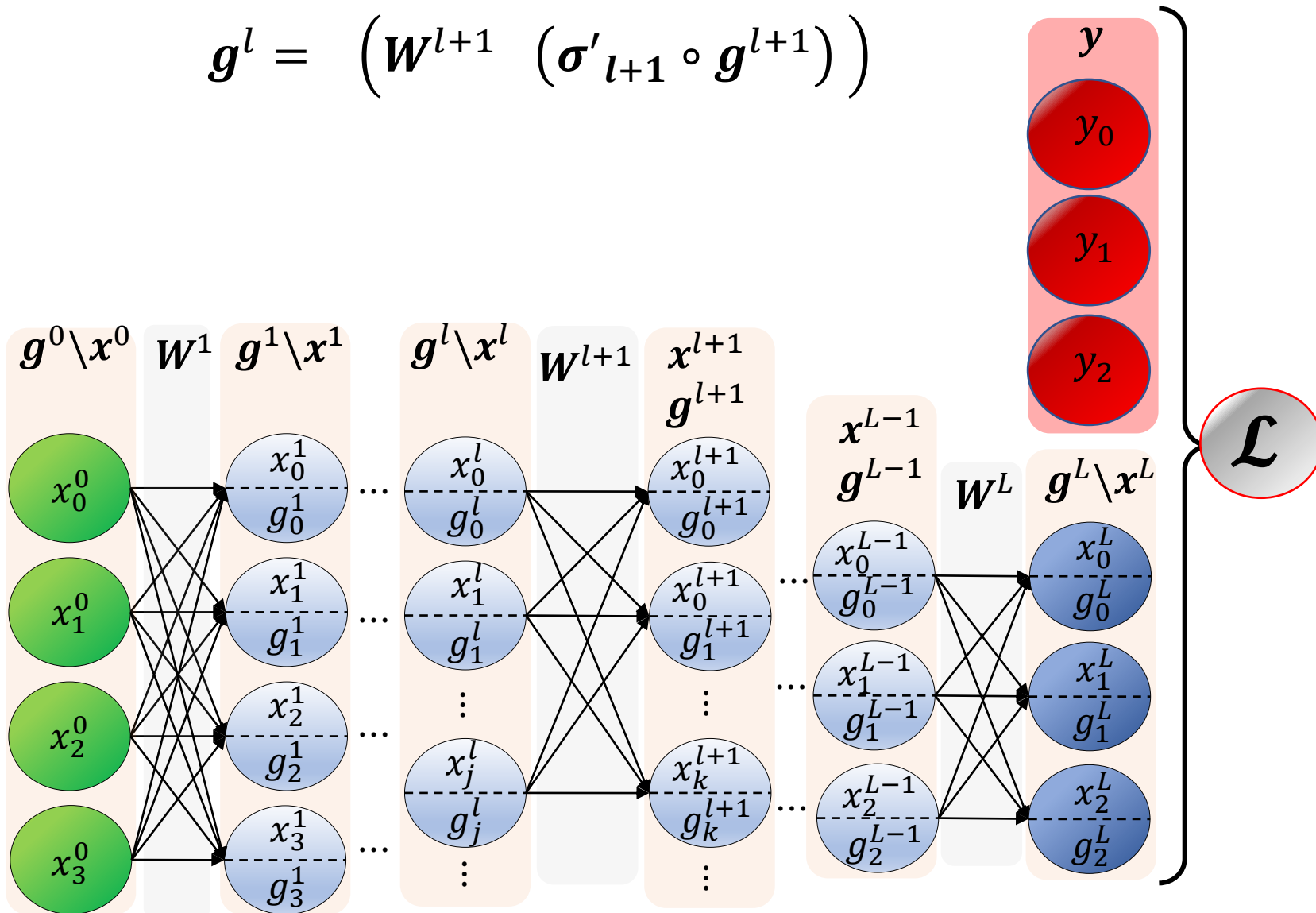
$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$

$$g^l = \left(W^{l+1} \left(\sigma'_{l+1} \circ g^{l+1} \right) \right)$$

Stopping criterion:

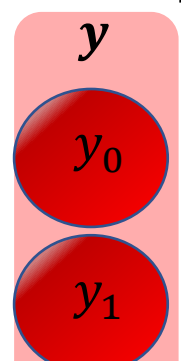
$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$



Back Propagation

$$\begin{aligned}
 g_j^l &\triangleq \frac{\partial \mathcal{L}}{\partial x_j^l} \\
 &= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l} \\
 &= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})
 \end{aligned}$$

$$\begin{aligned}
 x^l &= \sigma(W^{lT} x^{l-1} + b) \\
 g^l &= (W^{l+1} (\sigma'_{l+1} \circ g^{l+1}))
 \end{aligned}$$



Stopping criterion:

$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$



Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$

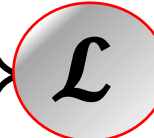
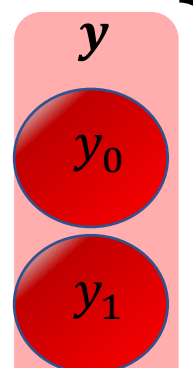
Stopping criterion:

$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$

$$x^l = \sigma(W^{lT} x^{l-1} + b)$$

$$g^l = (W^{l+1} (\sigma'_{l+1} \circ g^{l+1}))$$

Linear activation

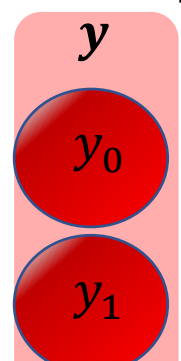


Back Propagation

$$\begin{aligned}
 g_j^l &\triangleq \frac{\partial \mathcal{L}}{\partial x_j^l} \\
 &= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l} \\
 &= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})
 \end{aligned}$$

$$\begin{aligned}
 x^l &= \sigma(W^{lT} x^{l-1} + b) \\
 g^l &= (W^{l+1} (\sigma'_{l+1} \circ g^{l+1}))
 \end{aligned}$$

Linear activation Not Transposed



Stopping criterion:

$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$



Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$

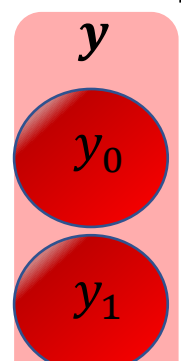
$$x^l = \sigma(W^{lT} x^{l-1} + b)$$

$$g^l = (W^{l+1} (\sigma'_{l+1} \circ g^{l+1}))$$

Linear activation

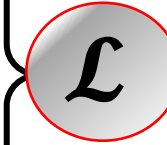
Not Transposed

Multiply by activation derivative



Stopping criterion:

$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$



Back Propagation

$$g_j^l \triangleq \frac{\partial \mathcal{L}}{\partial x_j^l}$$

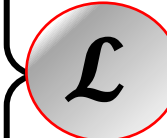
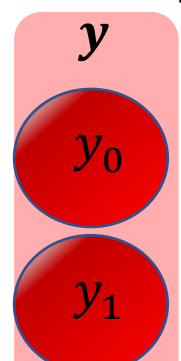
$$= \sum_k \frac{\partial \mathcal{L}}{\partial x_k^{l+1}} \cdot \frac{\partial x_k^{l+1}}{\partial x_j^l}$$

$$= \sum_k g_k^{l+1} \cdot w_{jk}^{l+1} \cdot \sigma'(z_k^{l+1})$$

$$x^l = \sigma(W^{lT} x^{l-1} + b)$$

$$g^l = (W^{l+1} (\sigma'_{l+1} \circ g^{l+1}))$$

Linear activation Not Transposed Multiply by activation derivative No bias

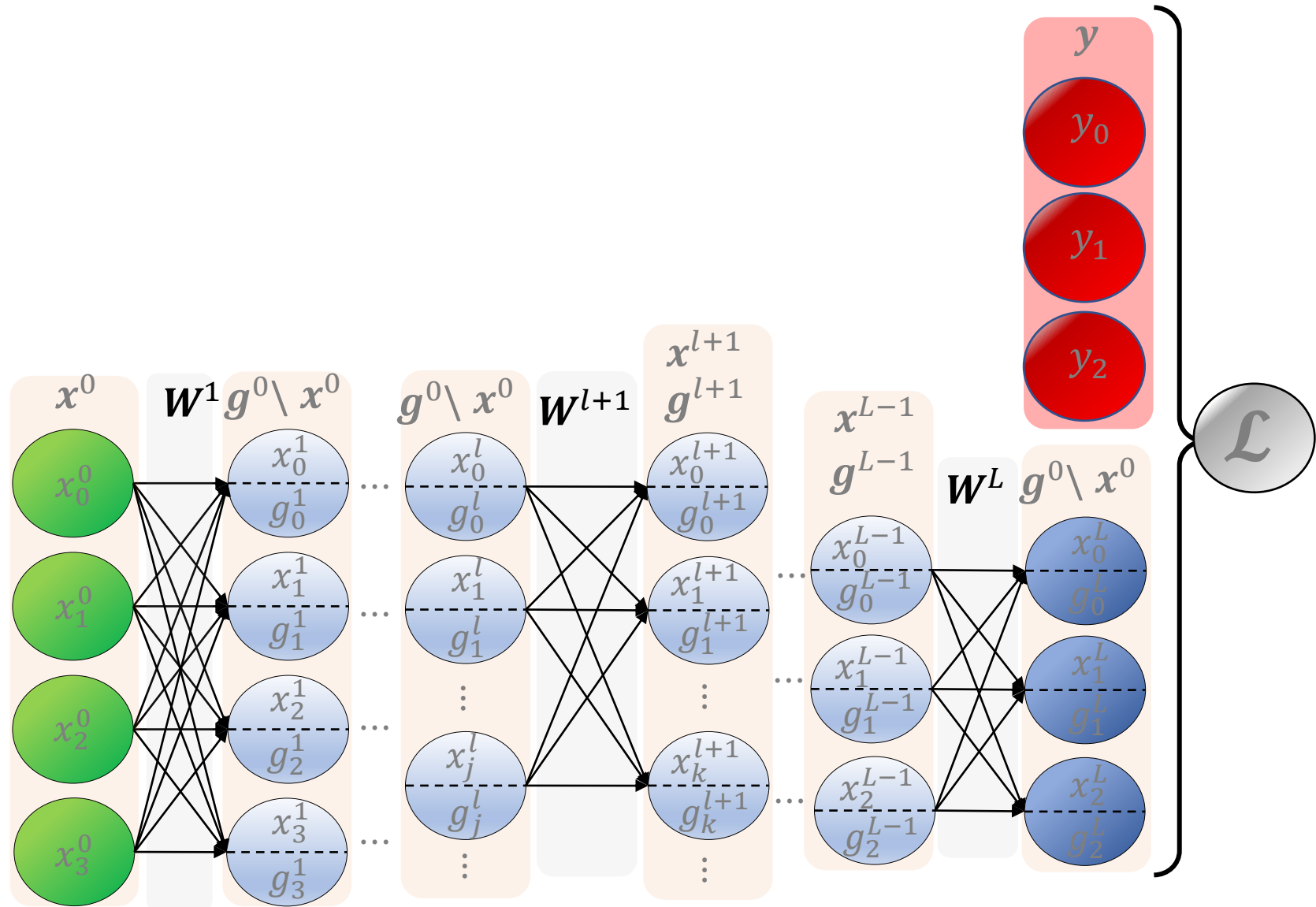


Stopping criterion:

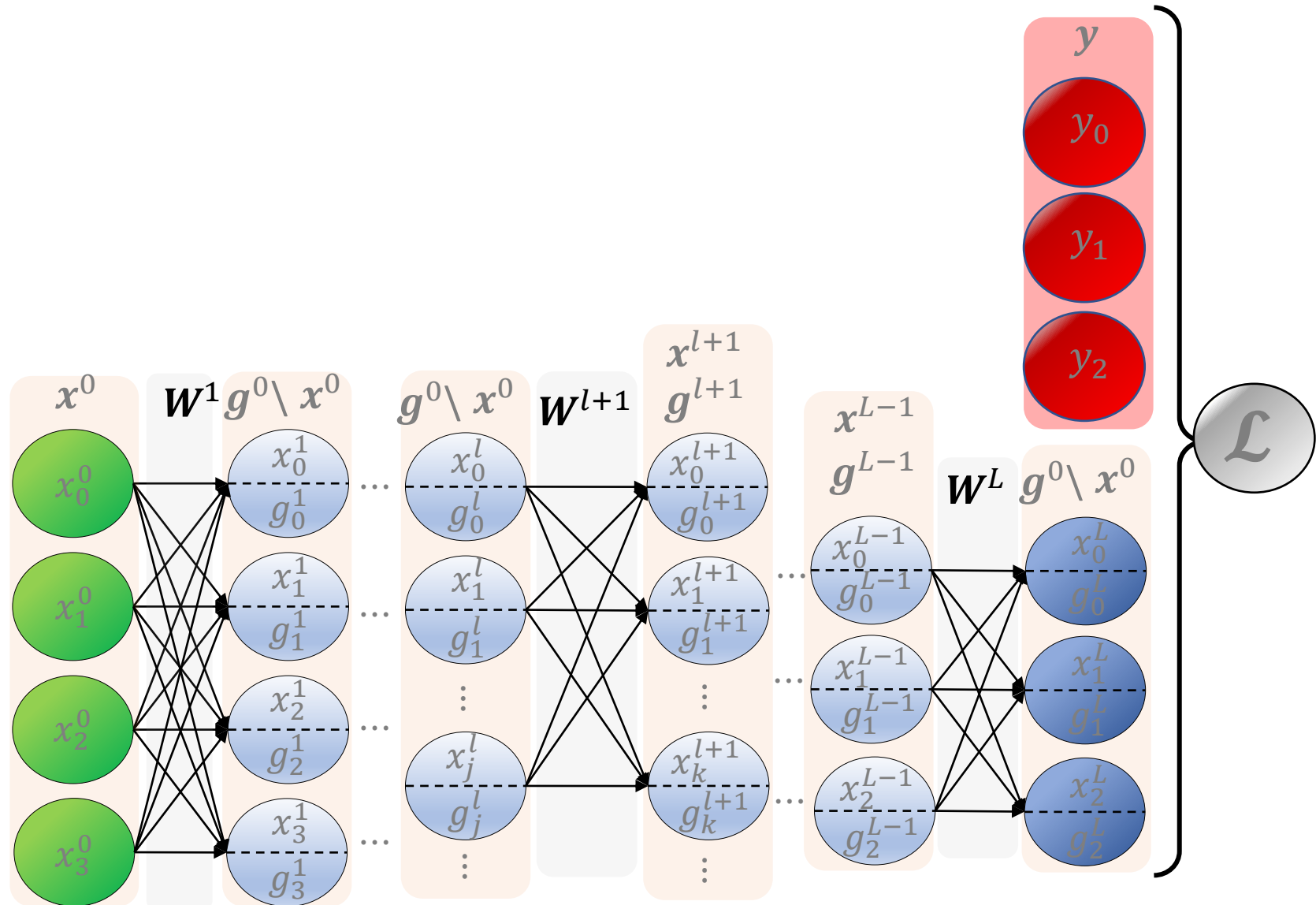
$$g_j^L = \frac{\partial \mathcal{L}}{\partial x_j^L}$$



Back Propagation

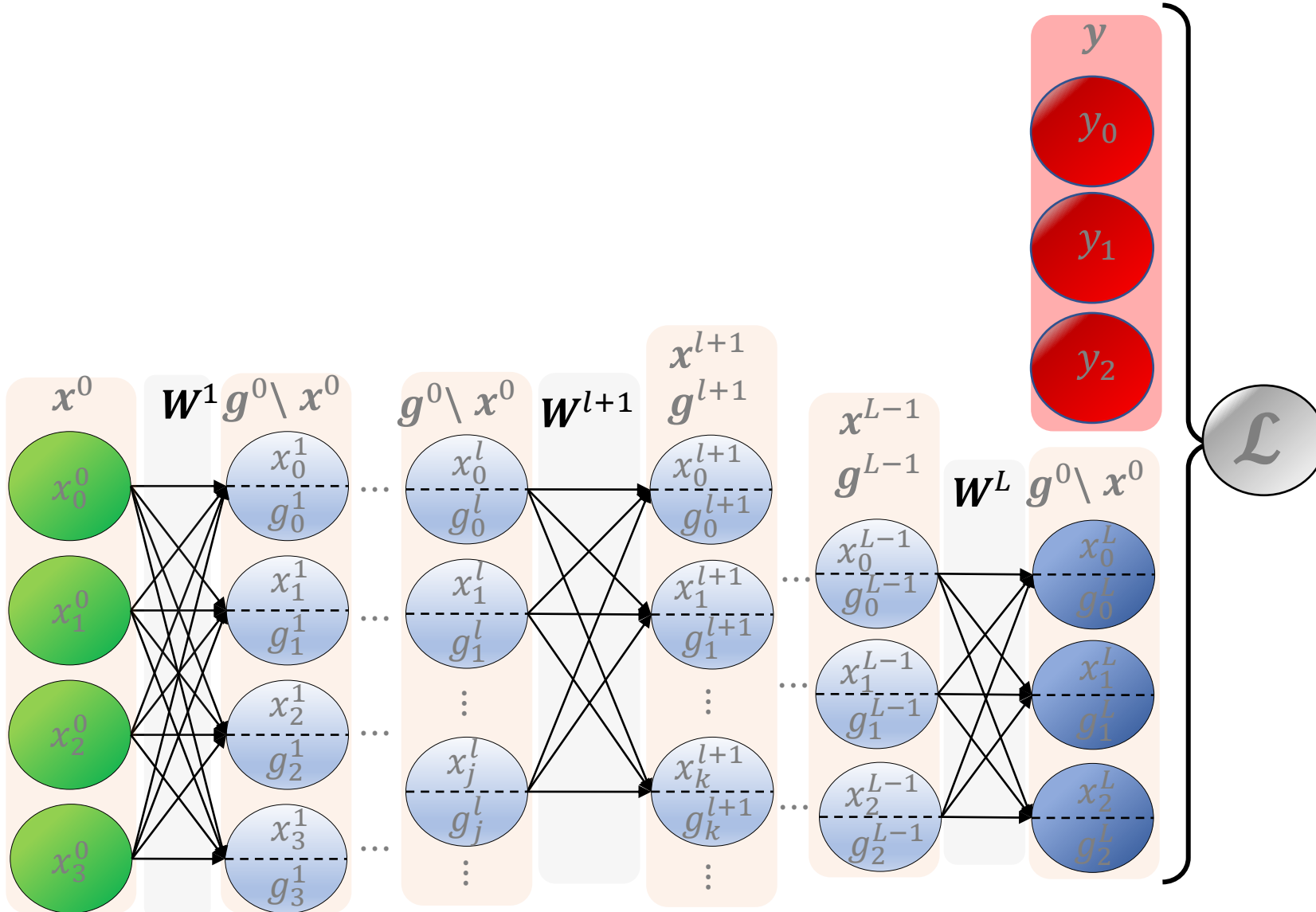


Back Propagation



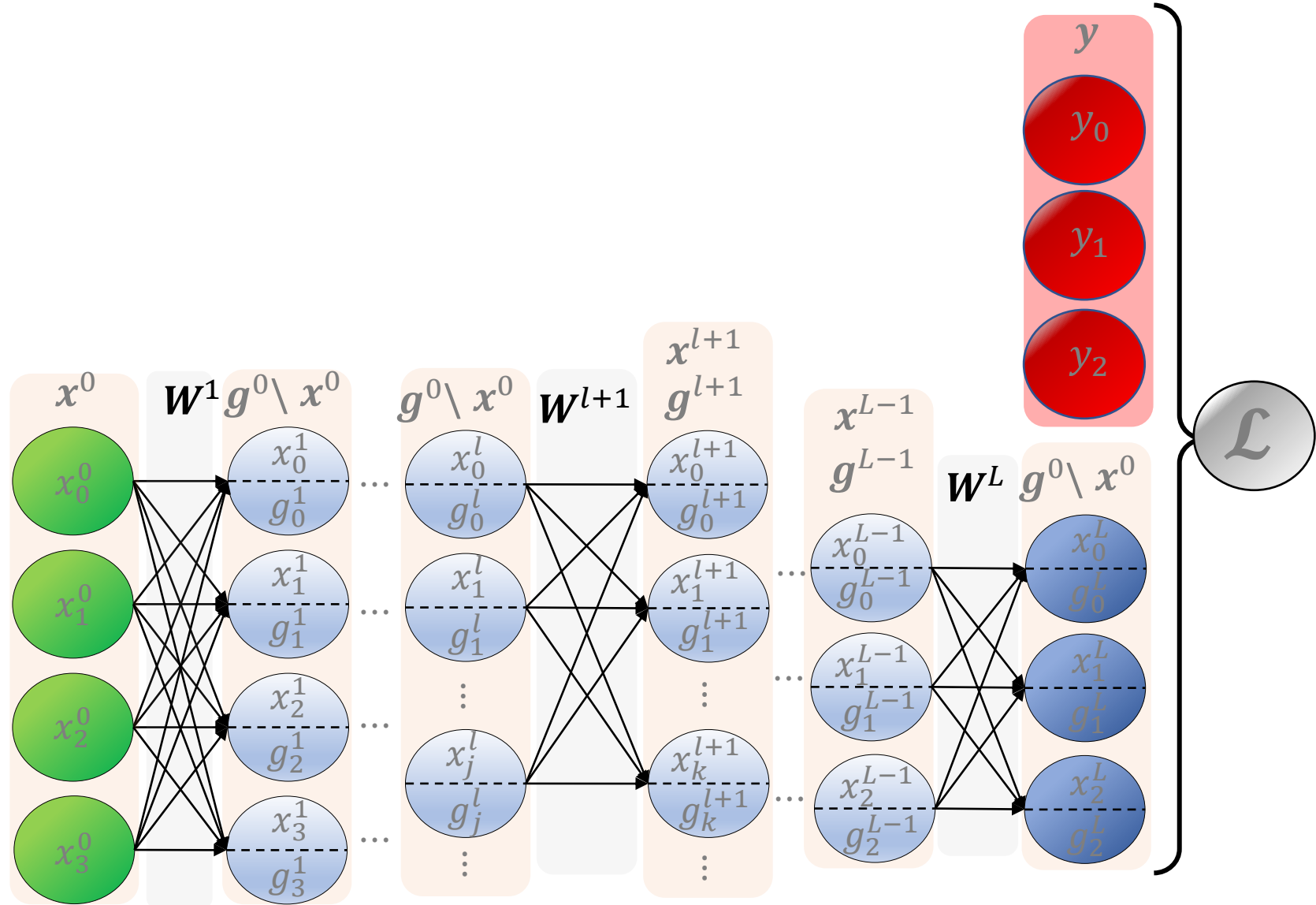
- Initialize weights

Back Propagation



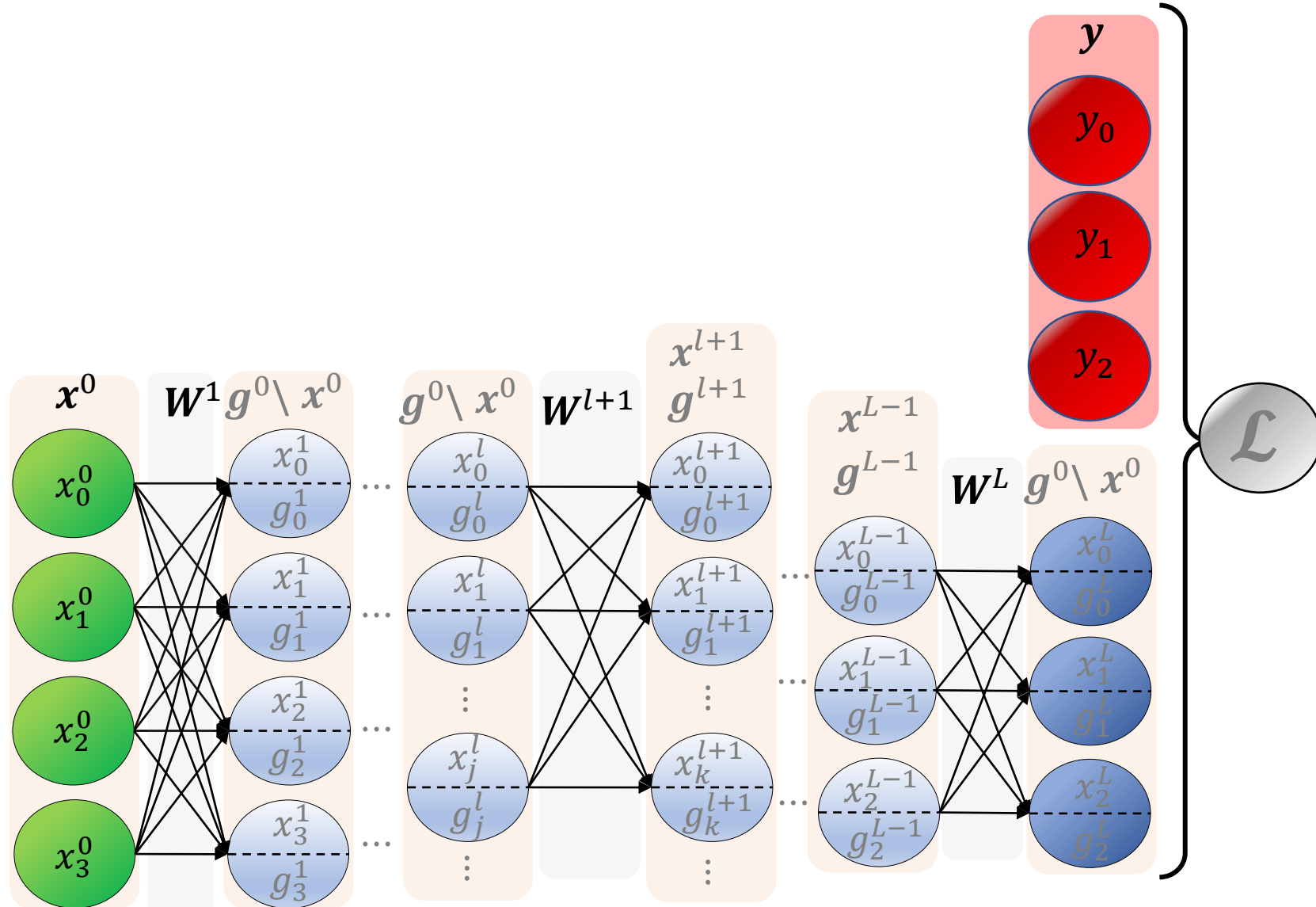
- Initialize weights
- Repeat until convergence:

Back Propagation

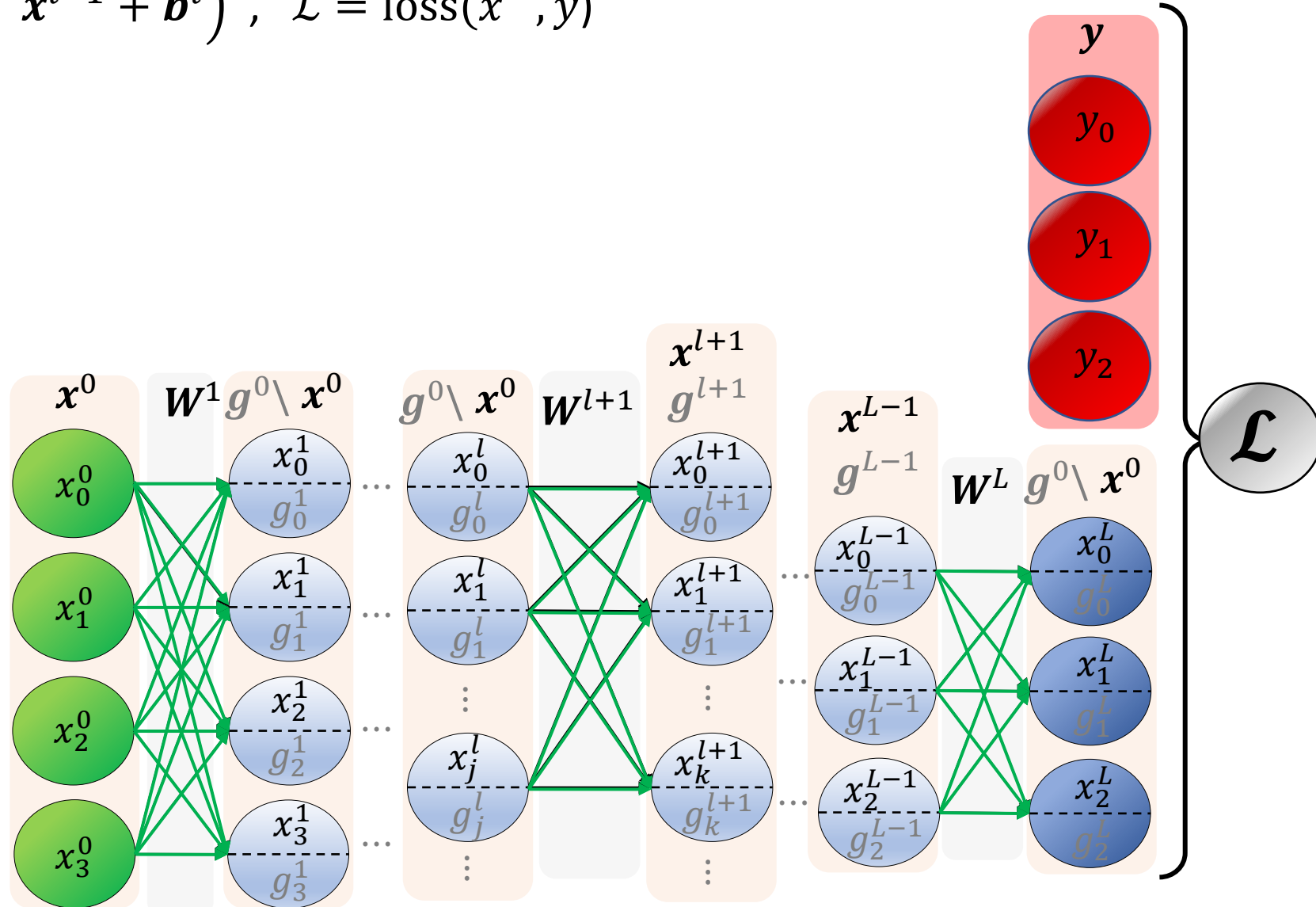


- Initialize weights
- Repeat until convergence:
 1. Sample a batch from the data: $\{(x_i, y_i) \dots\}$

Back Propagation



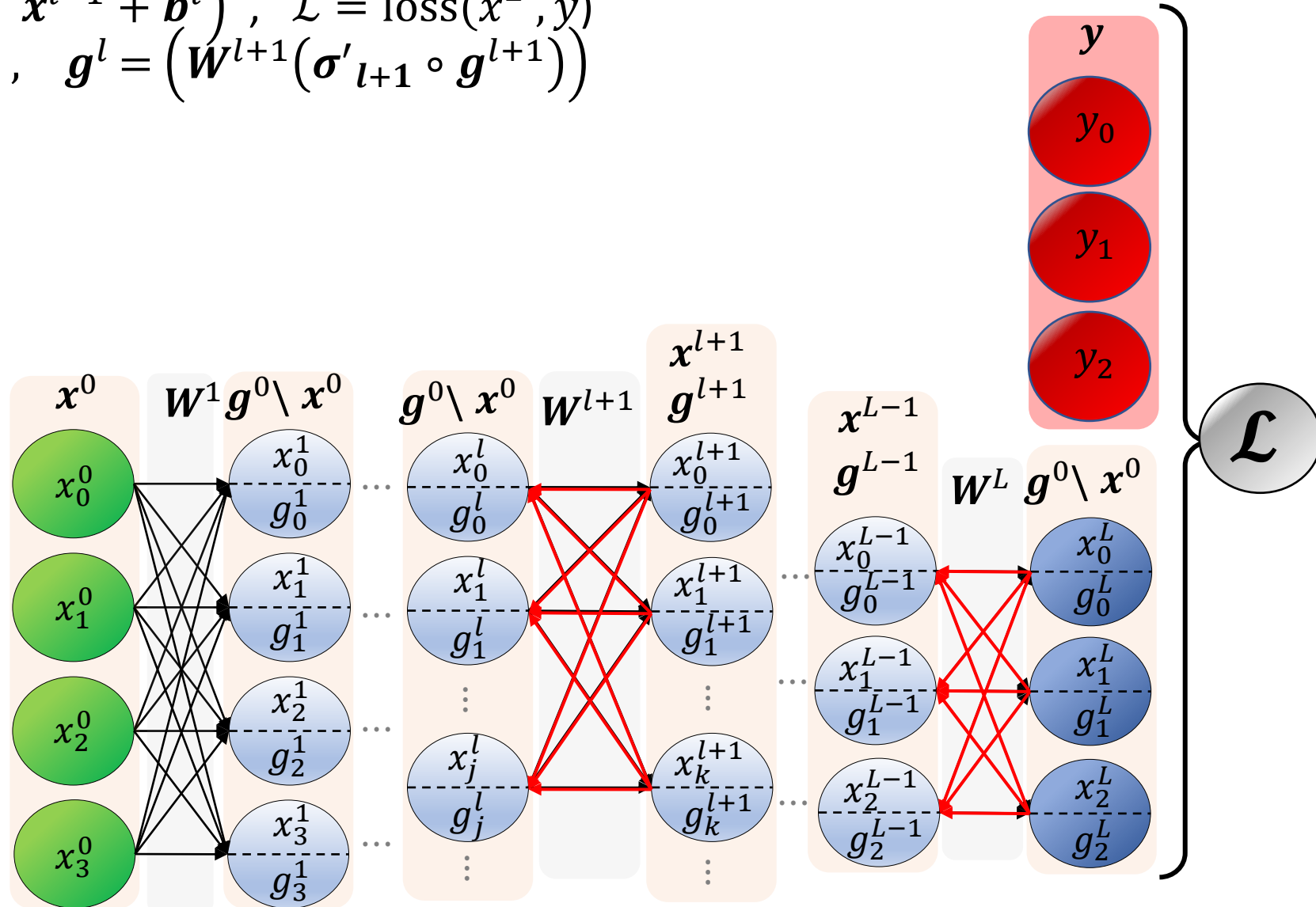
- Initialize weights
- Repeat until convergence:
 1. Sample a batch from the data: $\{(\mathbf{x}_i, \mathbf{y}_i) \dots\}$
 2. Forward pass: $\mathbf{x}^l = \sigma(\mathbf{W}^{lT} \mathbf{x}^{l-1} + \mathbf{b}^l)$, $\mathcal{L} = \text{loss}(\mathbf{x}^L, \mathbf{y})$



Back Propagation

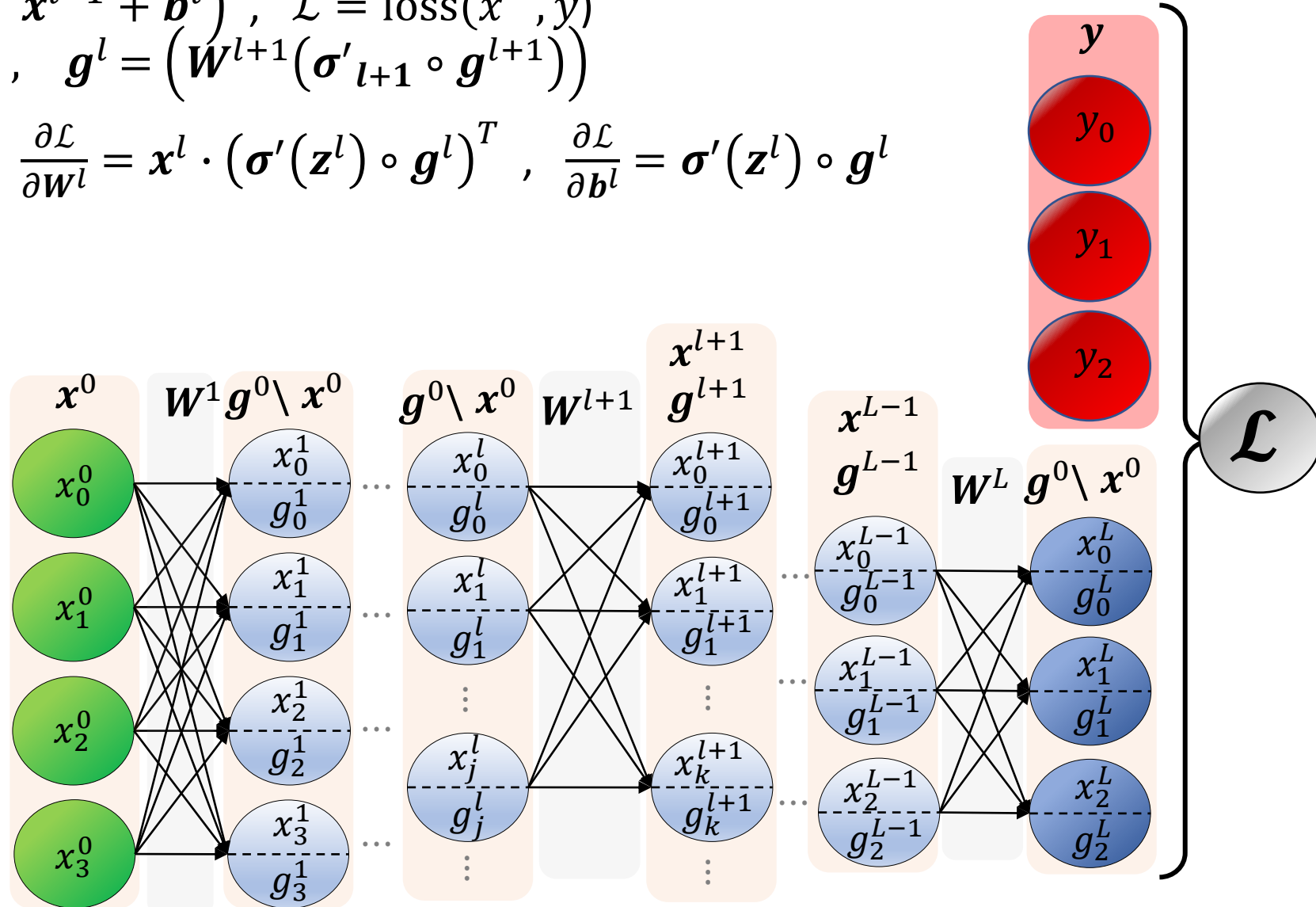
- Initialize weights
- Repeat until convergence:
 1. Sample a batch from the data: $\{(x_i, y_i) \dots\}$
 2. Forward pass: $x^l = \sigma(W^{lT} x^{l-1} + b^l)$, $\mathcal{L} = \text{loss}(x^L, y)$
 3. Backward pass: $g^L = \frac{\partial \mathcal{L}}{\partial x^L}$, $g^l = (W^{l+1}(\sigma'_{l+1} \circ g^{l+1}))$

Back Propagation



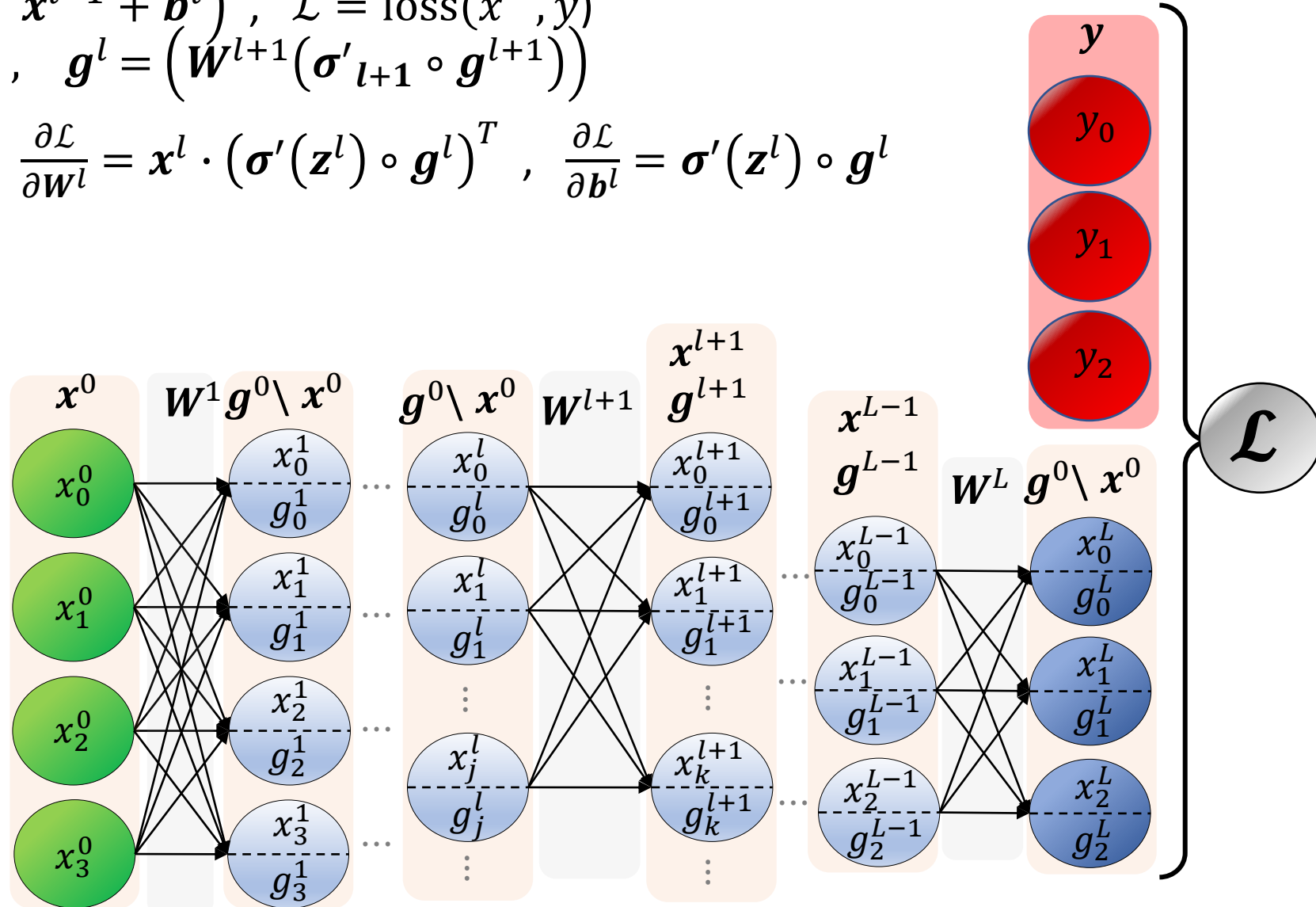
- Initialize weights
- Repeat until convergence:
 1. Sample a batch from the data: $\{(x_i, y_i) \dots\}$
 2. Forward pass: $x^l = \sigma(W^{lT} x^{l-1} + b^l)$, $\mathcal{L} = \text{loss}(x^L, y)$
 3. Backward pass: $g^L = \frac{\partial \mathcal{L}}{\partial x^L}$, $g^l = (W^{l+1}(\sigma'_{l+1} \circ g^{l+1}))$
 4. Calculate weights gradient: $\frac{\partial \mathcal{L}}{\partial W^l} = x^l \cdot (\sigma'(z^l) \circ g^l)^T$, $\frac{\partial \mathcal{L}}{\partial b^l} = \sigma'(z^l) \circ g^l$

Back Propagation



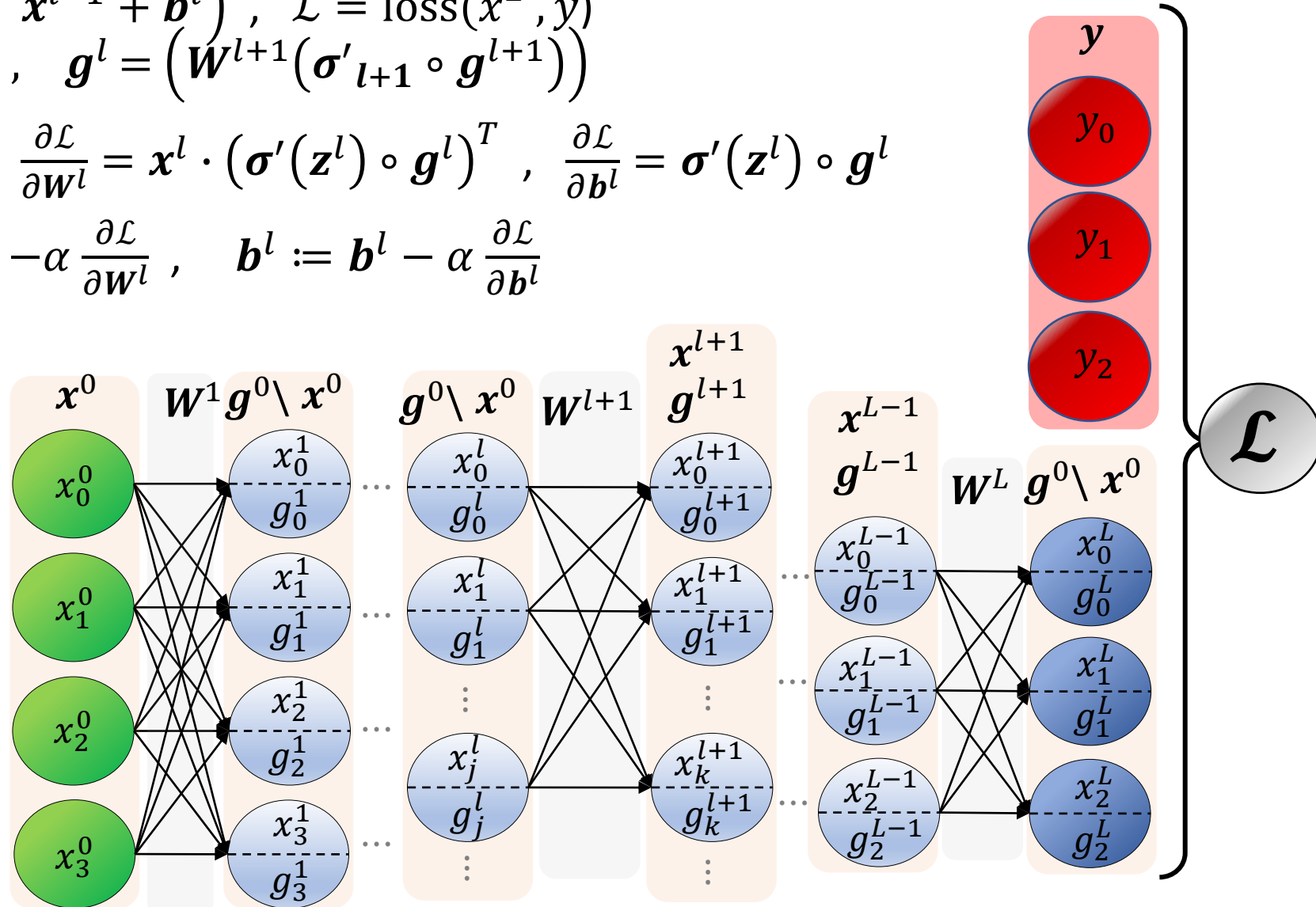
- Initialize weights
- Repeat until convergence:
 1. Sample a batch from the data: $\{(x_i, y_i) \dots\}$
 2. Forward pass: $x^l = \sigma(W^{lT} x^{l-1} + b^l)$, $\mathcal{L} = \text{loss}(x^L, y)$
 3. Backward pass: $g^L = \frac{\partial \mathcal{L}}{\partial x^L}$, $g^l = (W^{l+1}(\sigma'_{l+1} \circ g^{l+1}))$
 4. Calculate weights gradient: $\frac{\partial \mathcal{L}}{\partial W^l} = x^l \cdot (\sigma'(z^l) \circ g^l)^T$, $\frac{\partial \mathcal{L}}{\partial b^l} = \sigma'(z^l) \circ g^l$

Back Propagation



- Initialize weights
- Repeat until convergence:
 1. Sample a batch from the data: $\{(x_i, y_i) \dots\}$
 2. Forward pass: $x^l = \sigma(W^{lT} x^{l-1} + b^l)$, $\mathcal{L} = \text{loss}(x^L, y)$
 3. Backward pass: $g^L = \frac{\partial \mathcal{L}}{\partial x^L}$, $g^l = (W^{l+1}(\sigma'_{l+1} \circ g^{l+1}))$
 4. Calculate weights gradient: $\frac{\partial \mathcal{L}}{\partial W^l} = x^l \cdot (\sigma'(z^l) \circ g^l)^T$, $\frac{\partial \mathcal{L}}{\partial b^l} = \sigma'(z^l) \circ g^l$
 5. Update weights: $W^l := W^l - \alpha \frac{\partial \mathcal{L}}{\partial W^l}$, $b^l := b^l - \alpha \frac{\partial \mathcal{L}}{\partial b^l}$

Back Propagation



Let's get more generic



Yann LeCun

January 5, 2018 · 🌐

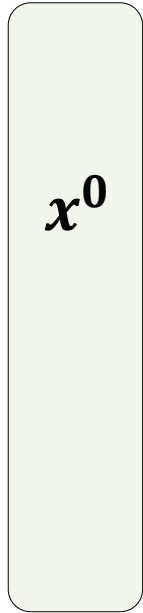


OK, Deep Learning has outlived its usefulness as a buzz-phrase. Deep Learning est mort. Vive **Differentiable Programming!**

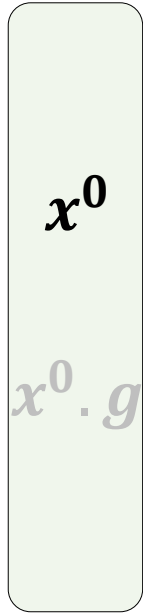
Yeah, Differentiable Programming is little more than a rebranding of the modern collection Deep Learning techniques, the same way Deep Learning was a rebranding of the modern incarnations of neural nets with more than two layers.

But the important point is that people are now building a new kind of software by assembling networks of **parameterized functional blocks** and by training them from examples using some form of gradient-based optimization.

Let's get more generic

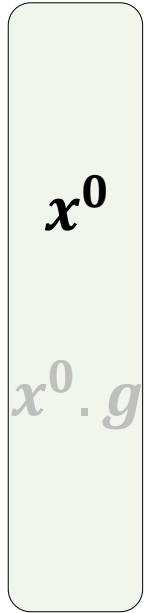


Let's get more generic

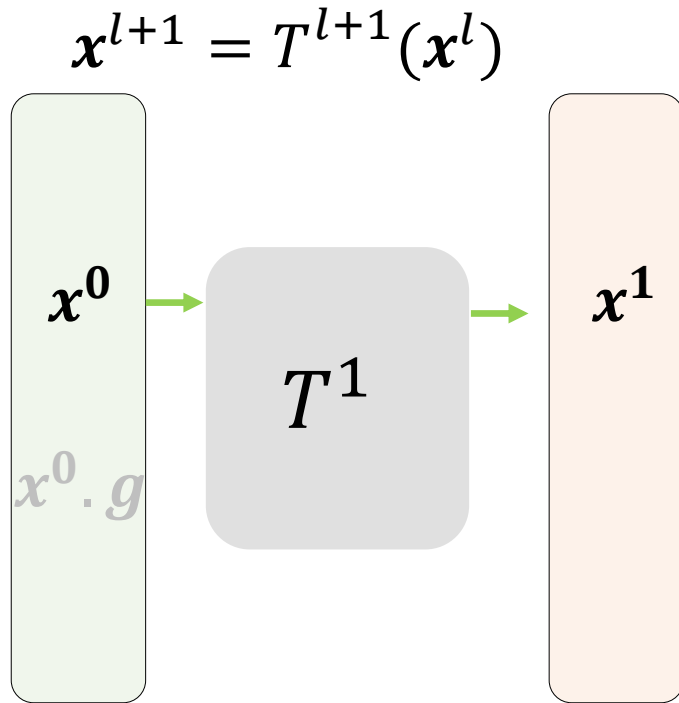


Let's get more generic

$$\mathbf{x}^{l+1} = T^{l+1}(\mathbf{x}^l)$$

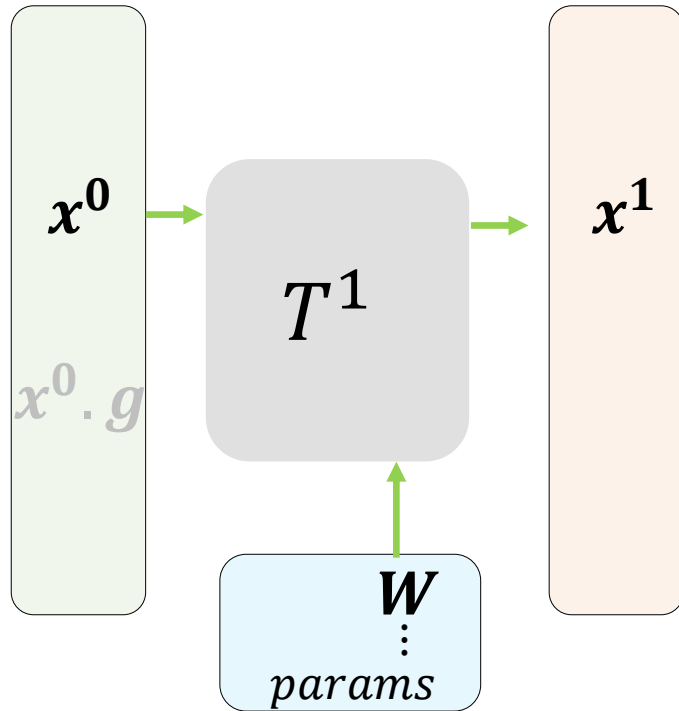


Let's get more generic



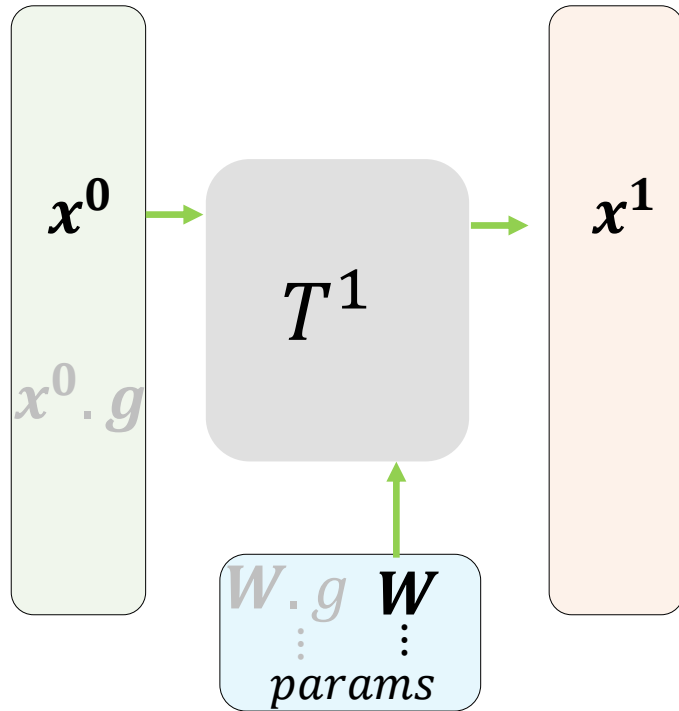
Let's get more generic

$$\mathbf{x}^{l+1} = T^{l+1}(\mathbf{x}^l; T^{l+1}.params)$$



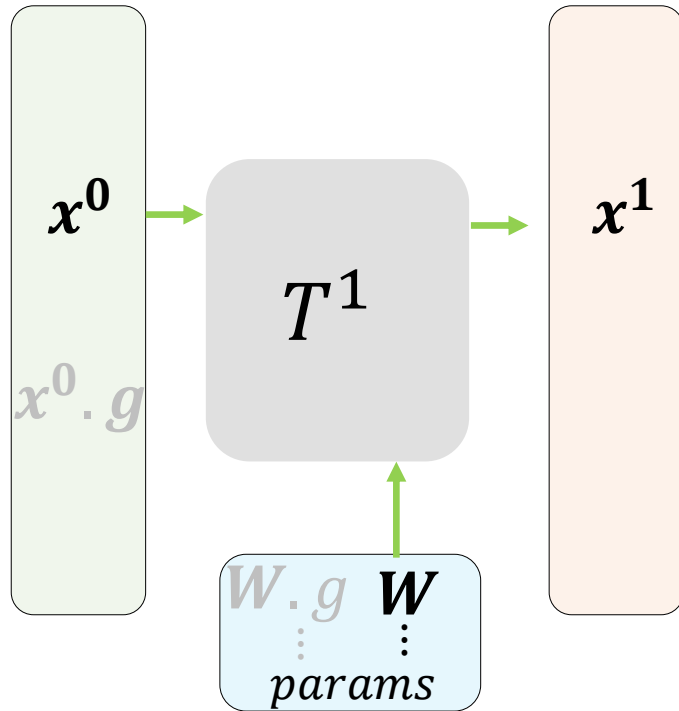
Let's get more generic

$$\mathbf{x}^{l+1} = T^{l+1}(\mathbf{x}^l; T^{l+1}.params)$$



Let's get more generic

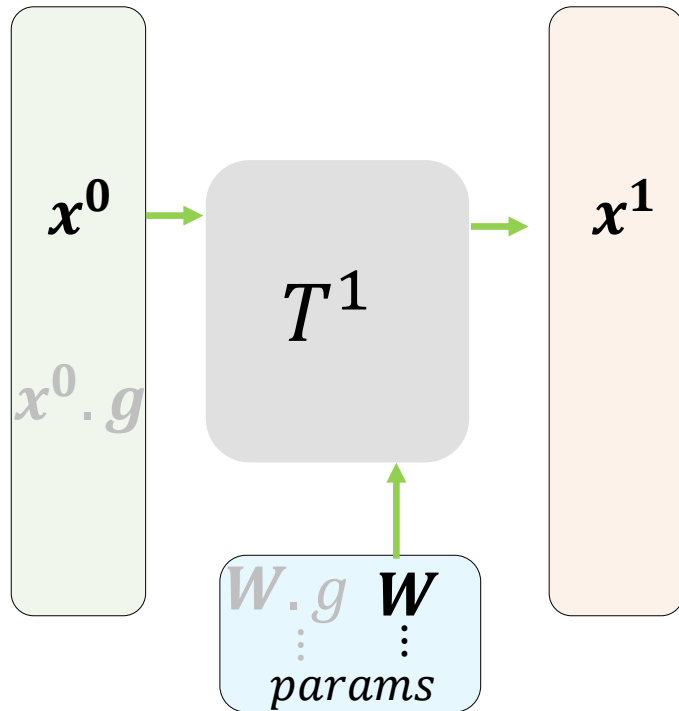
$$\mathbf{x}^{l+1} = T^{l+1}(\mathbf{x}^l; T^{l+1}.params)$$



$$\mathbf{x}^l.g_{grad} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^l}$$

Let's get more generic

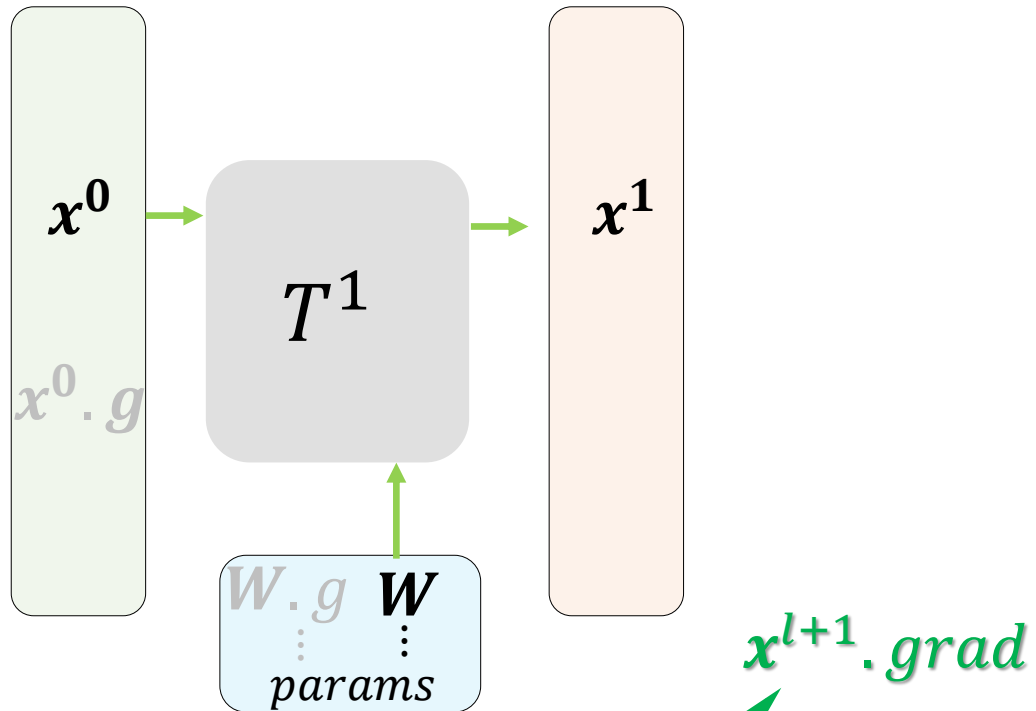
$$\mathbf{x}^{l+1} = T^{l+1}(\mathbf{x}^l; T^{l+1}.params)$$



$$\mathbf{x}^l.grad = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}} \cdot \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l}$$

Let's get more generic

$$\mathbf{x}^{l+1} = T^{l+1}(\mathbf{x}^l; T^{l+1}.params)$$

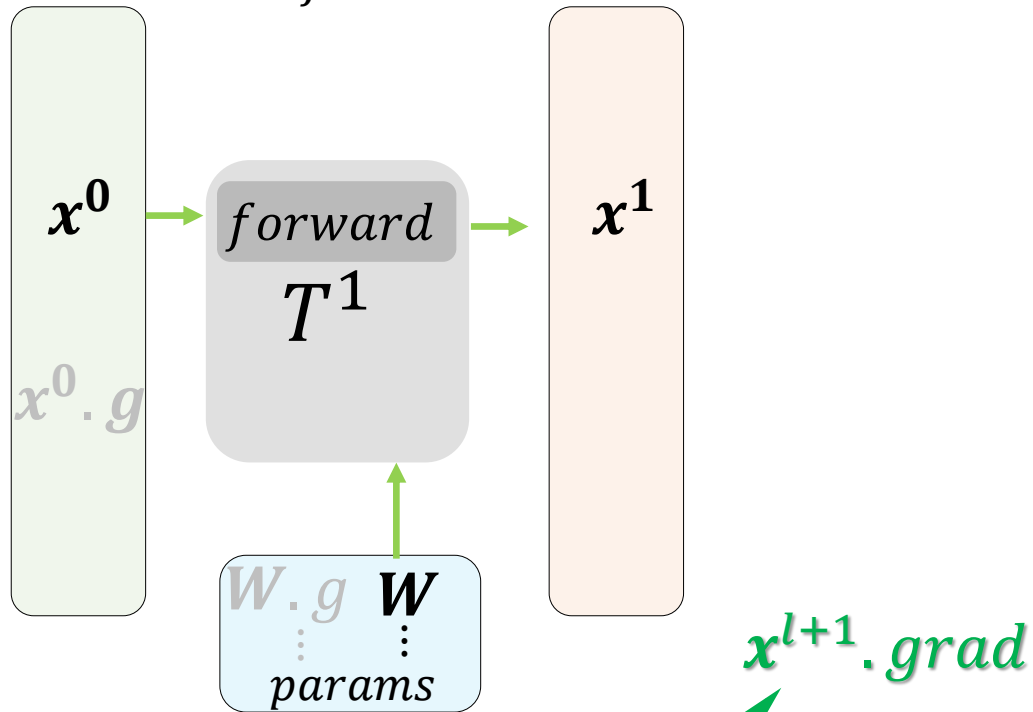


$$\mathbf{x}^l.grad = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^l} = \boxed{\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}}} \cdot \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l}$$

$\mathbf{x}^{l+1}.grad$

Let's get more generic

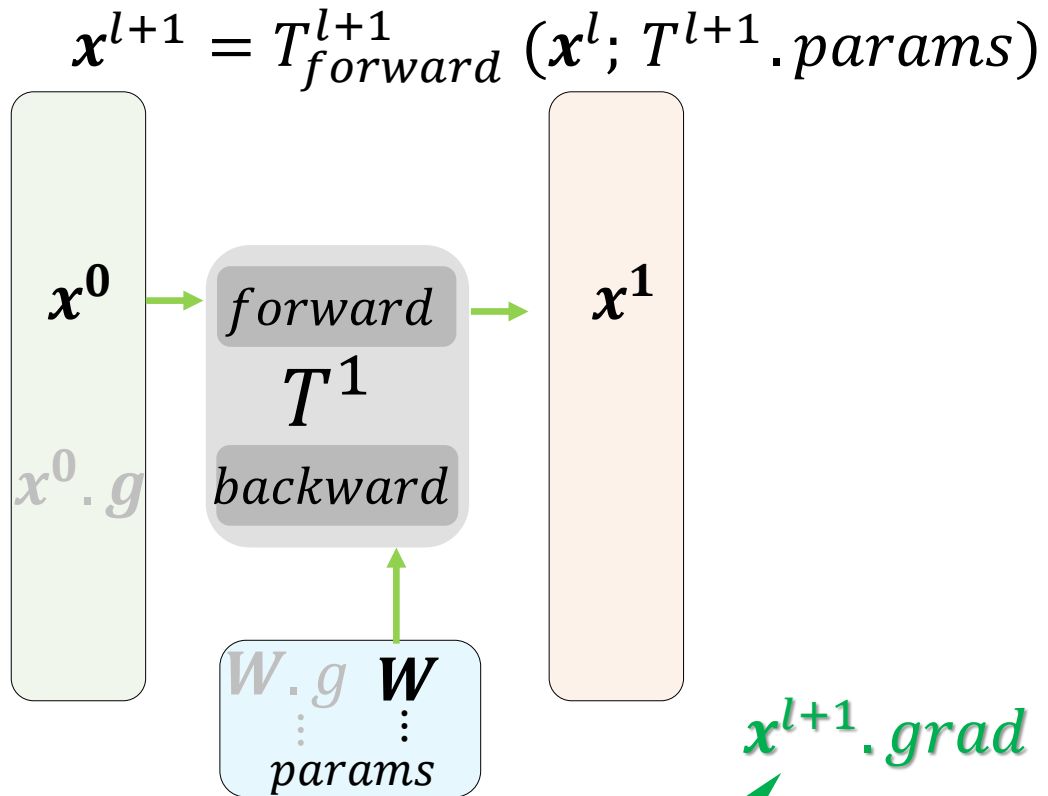
$$\mathbf{x}^{l+1} = T_{\text{forward}}^{l+1}(\mathbf{x}^l; T^{l+1}.\text{params})$$



$$\mathbf{x}^l.g\text{rad} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}} \cdot \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l}$$

The term $\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}}$ is highlighted with a green box and labeled $\mathbf{x}^{l+1}.grad$ with a green arrow.

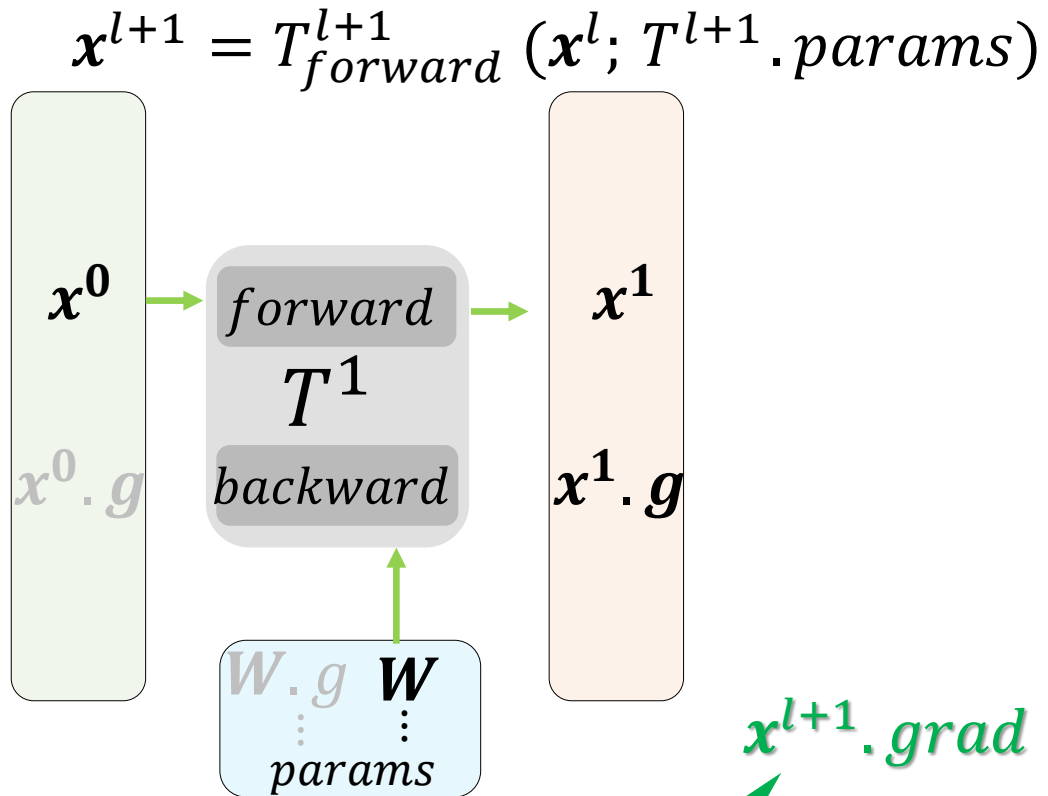
Let's get more generic



$$\mathbf{x}^l.grad = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}} \cdot \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l} = T_{backward}^{l+1}(\mathbf{x}^{l+1}.grad; T^{l+1}.params, \mathbf{x}^l)$$

$\mathbf{x}^{l+1}.grad$

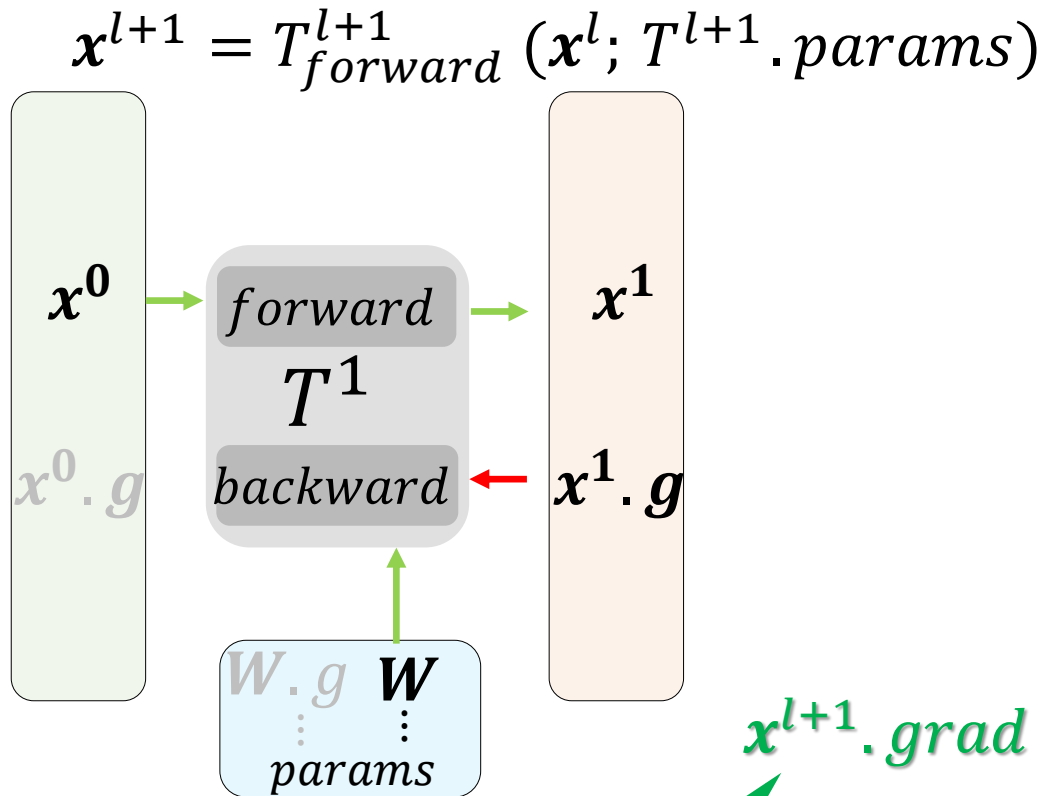
Let's get more generic



$$\mathbf{x}^l.grad = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}} \cdot \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l} = T_{backward}^{l+1}(\mathbf{x}^{l+1}.grad; T^{l+1}.params, \mathbf{x}^l)$$

The term $\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}}$ is highlighted with a green box and labeled $\mathbf{x}^{l+1}.grad$ with a green arrow.

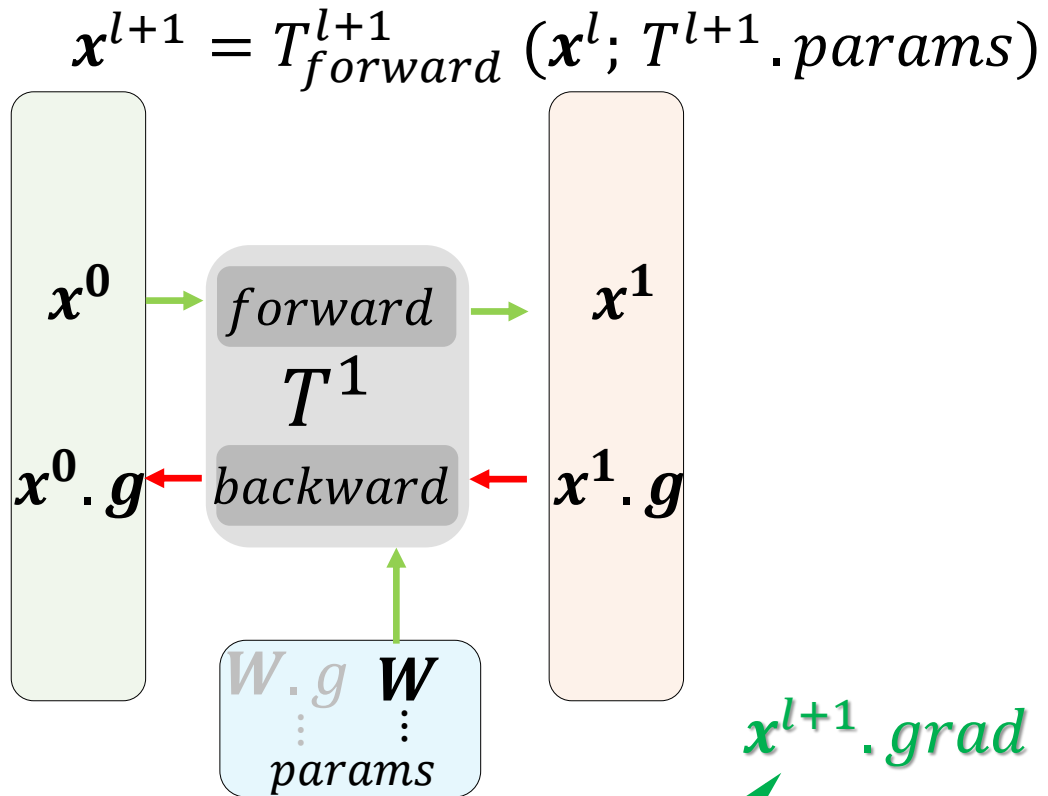
Let's get more generic



$$\mathbf{x}^l.grad = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}} \cdot \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l} = T_{backward}^{l+1}(\mathbf{x}^{l+1}.grad; T^{l+1}.params, \mathbf{x}^l)$$

$\mathbf{x}^{l+1}.grad$

Let's get more generic

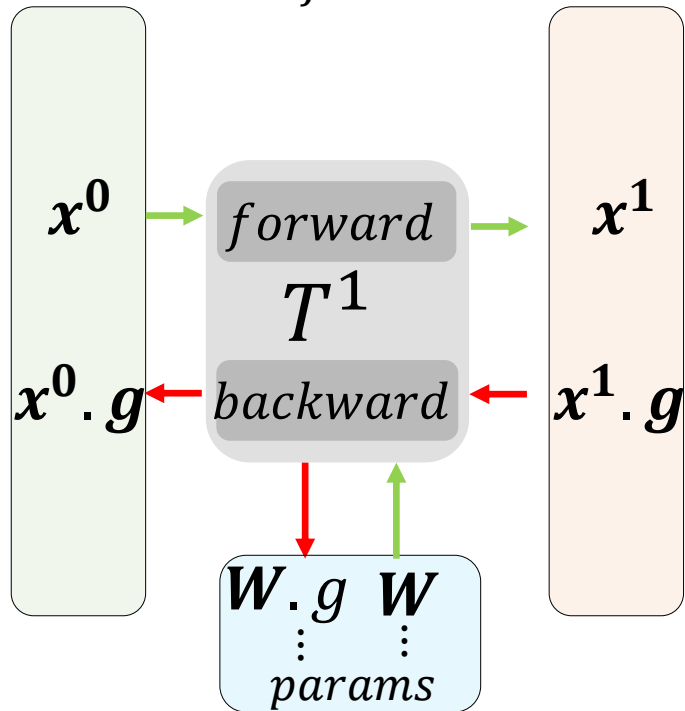


$$\mathbf{x}^l.grad = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^l} = \frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}} \cdot \frac{\partial \mathbf{x}^{l+1}}{\partial \mathbf{x}^l} = T_{backward}^{l+1}(\mathbf{x}^{l+1}.grad; T^{l+1}.params, \mathbf{x}^l)$$

The term $\frac{\partial \mathcal{L}}{\partial \mathbf{x}^{l+1}}$ is highlighted with a green box and labeled $\mathbf{x}^{l+1}.grad$ with a green arrow.

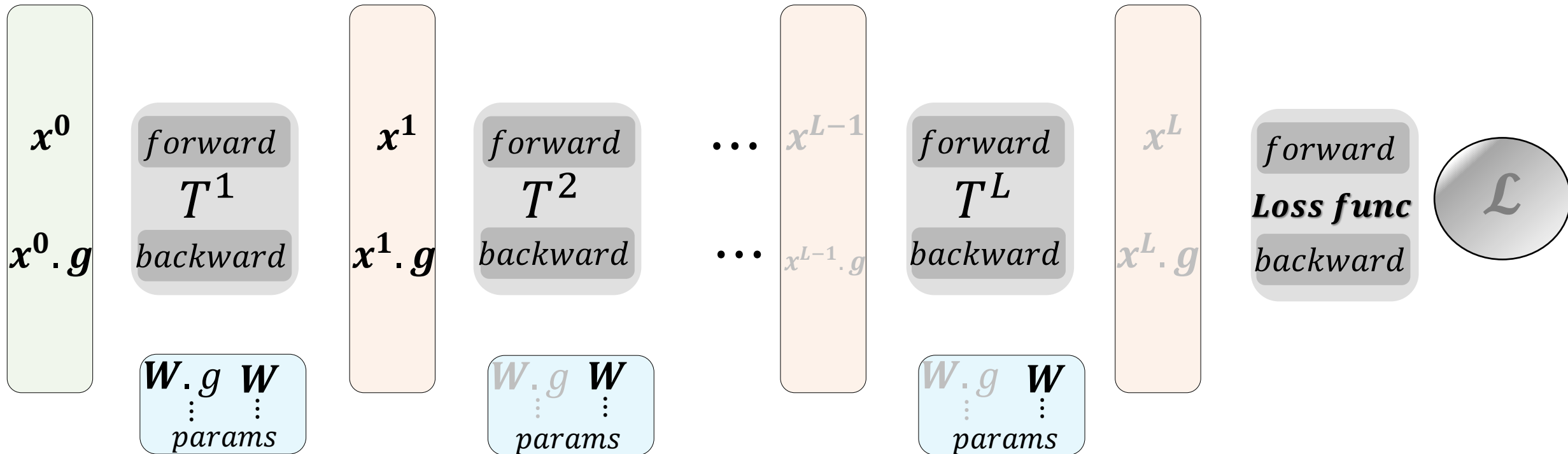
Let's get more generic

$$\mathbf{x}^{l+1} = T_{\text{forward}}^{l+1}(\mathbf{x}^l; T^{l+1}.\text{params})$$

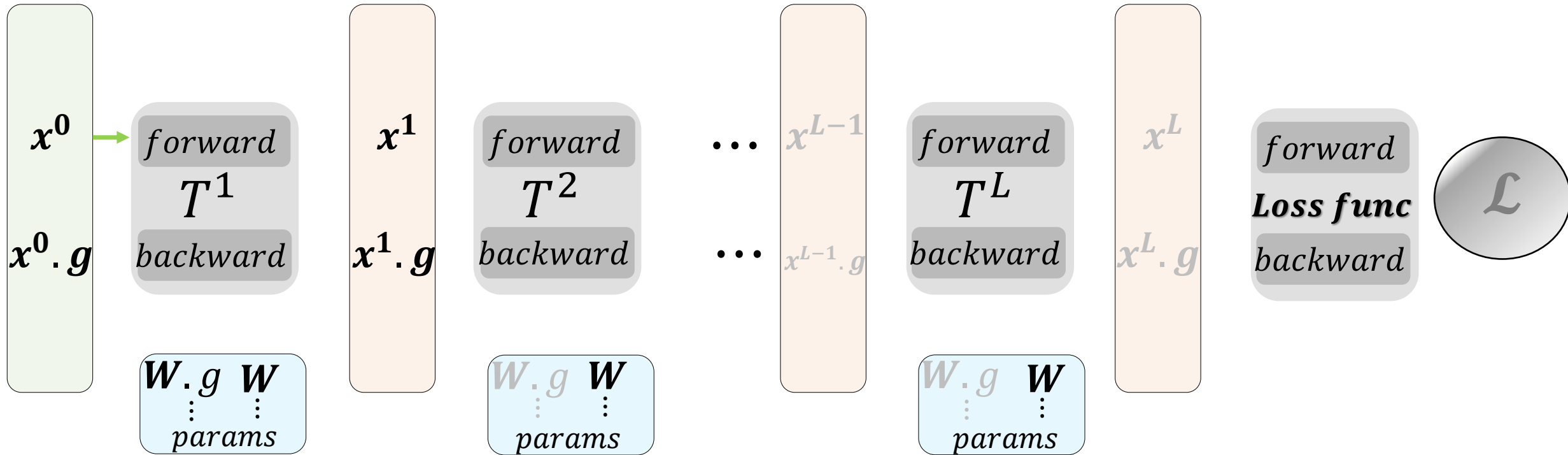


$$\begin{pmatrix} \mathbf{x}^l.\text{grad} \\ T^{l+1}.\text{params}.\text{grad} \end{pmatrix} = T_{\text{backward}}^{l+1}(\mathbf{x}^{l+1}.\text{grad}; T^{l+1}.\text{params}, \mathbf{x}^l)$$

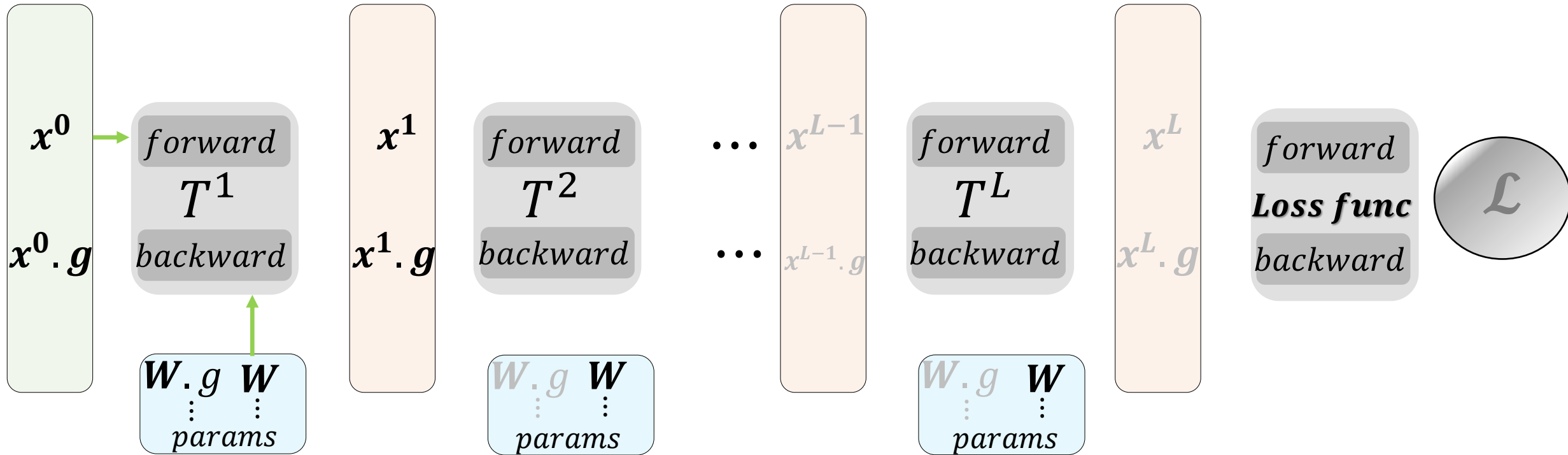
Let's get more generic



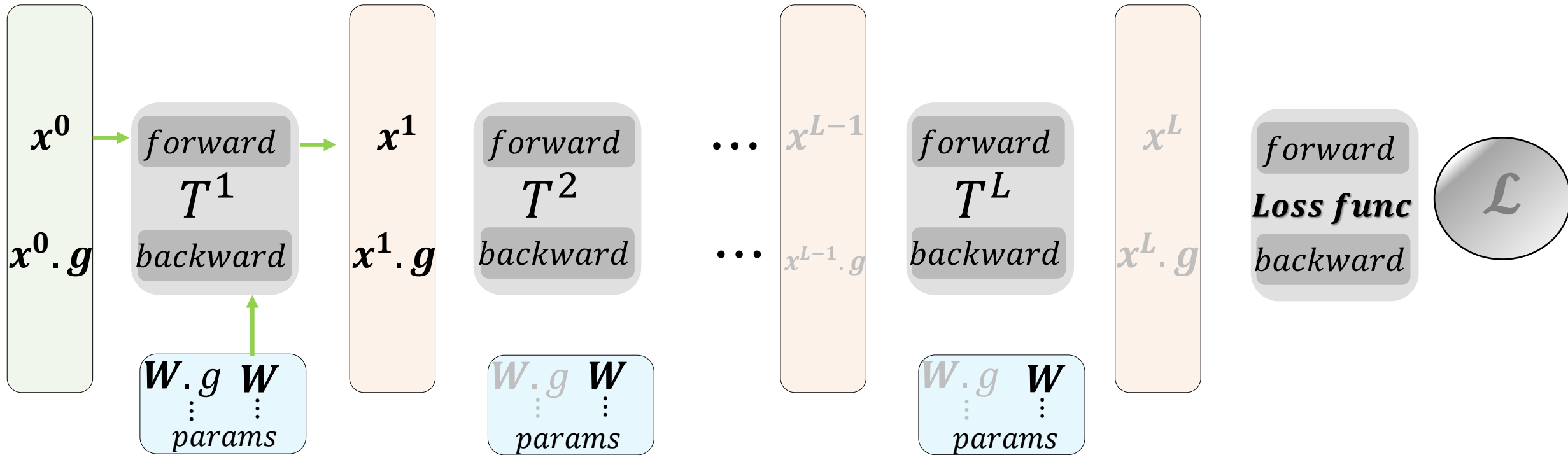
Let's get more generic



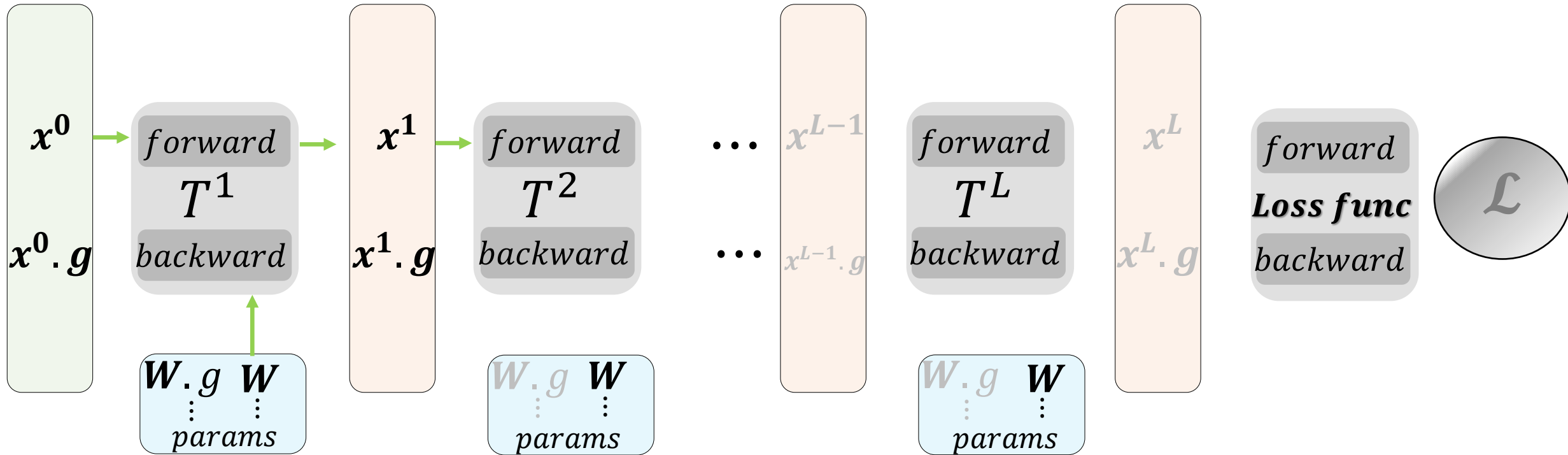
Let's get more generic



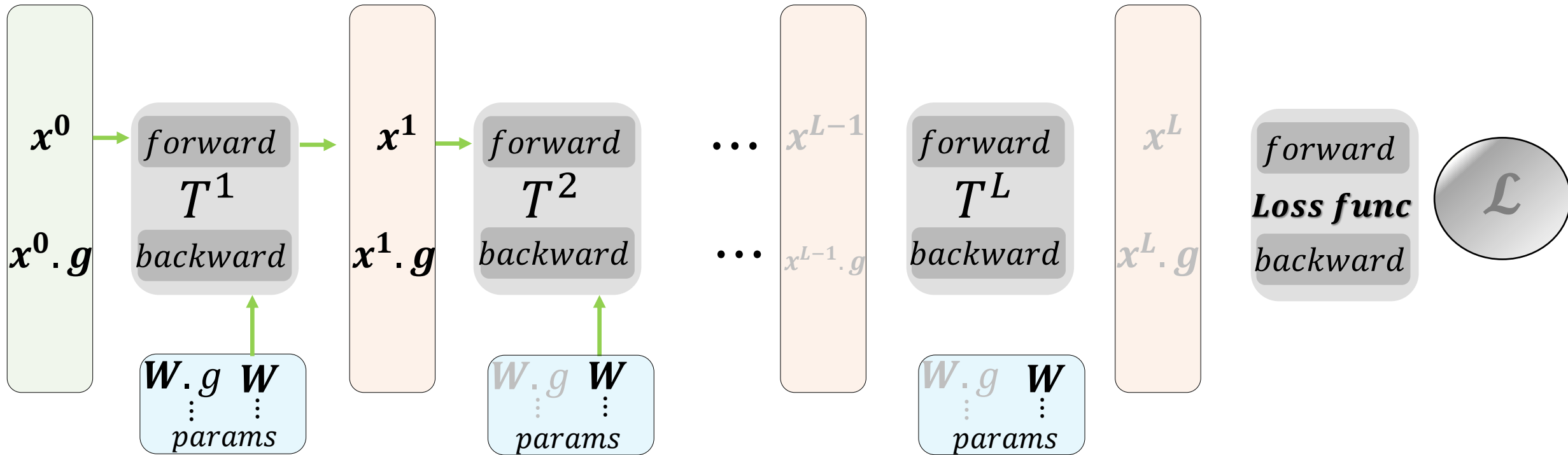
Let's get more generic



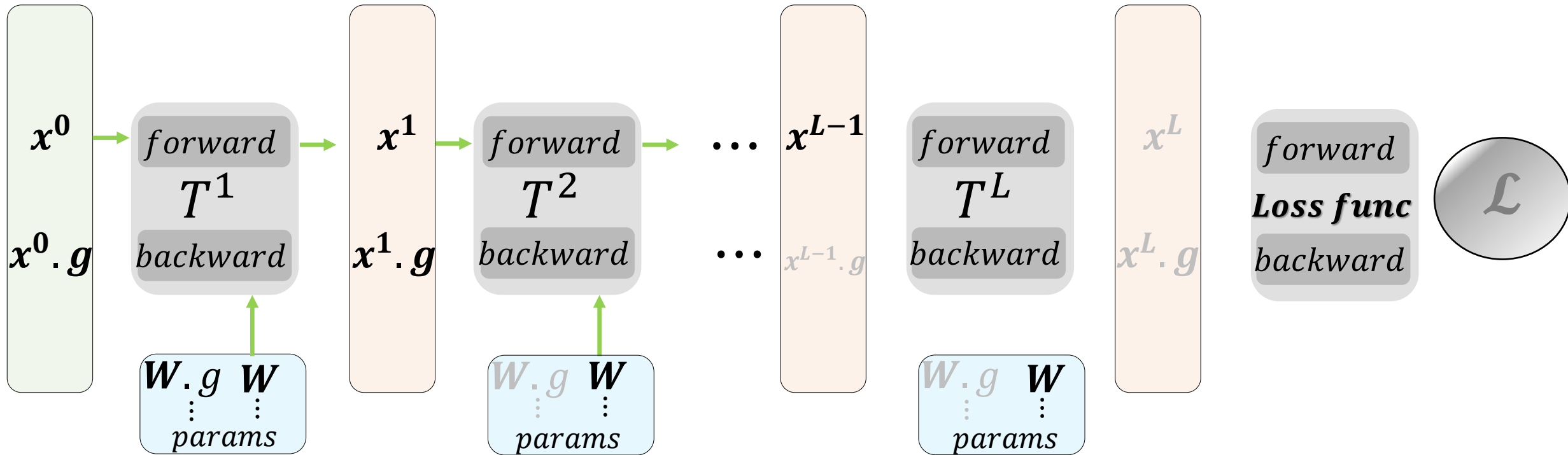
Let's get more generic



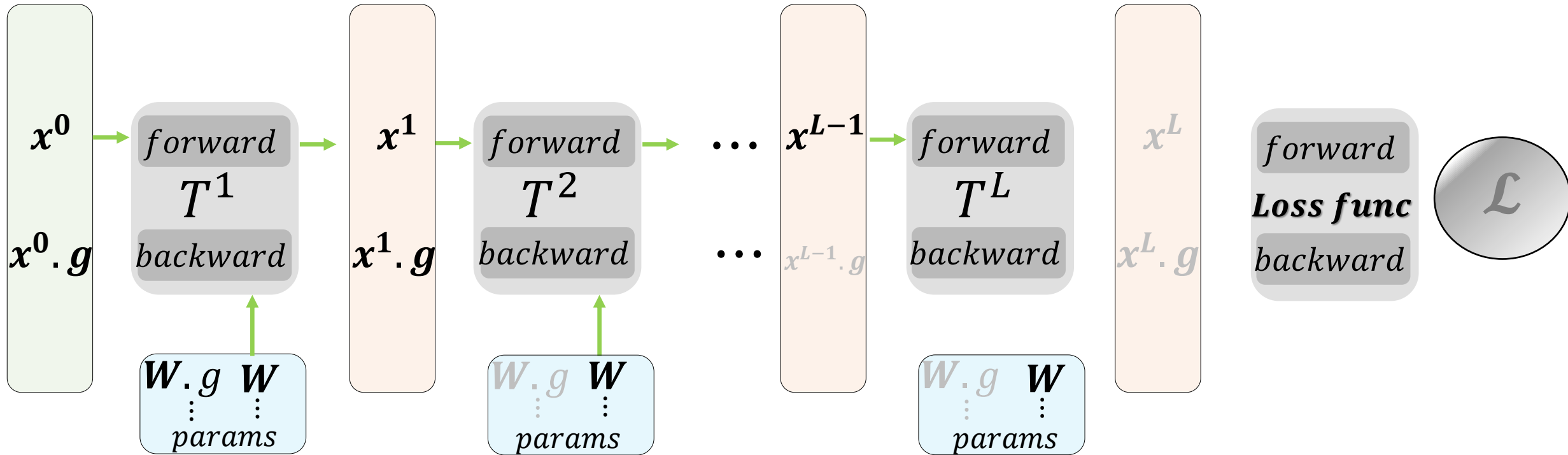
Let's get more generic



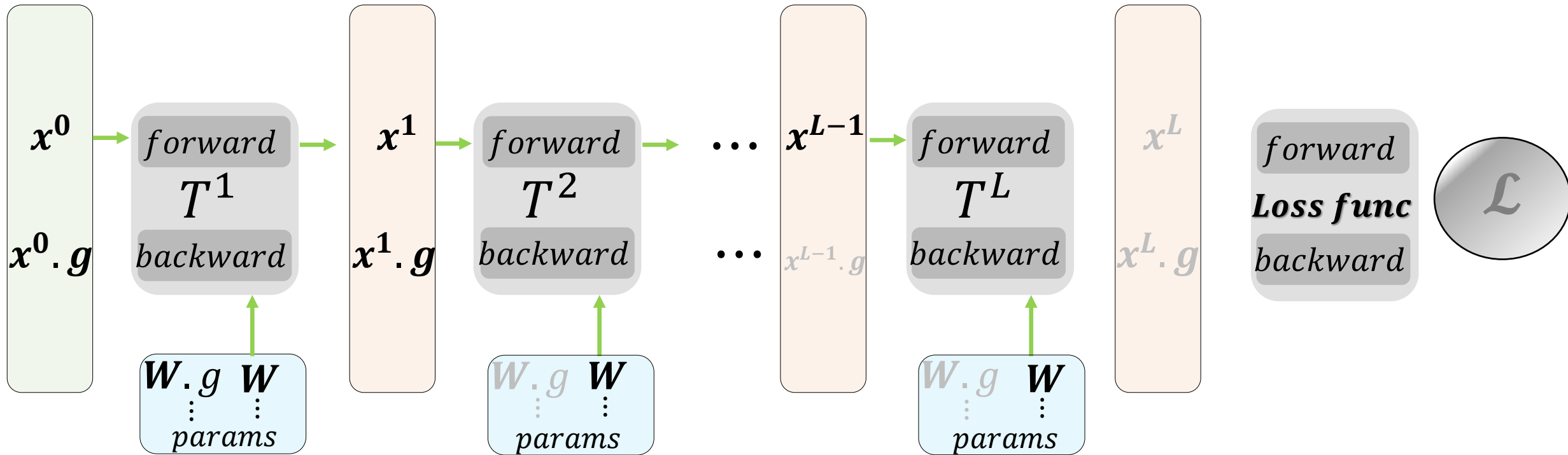
Let's get more generic



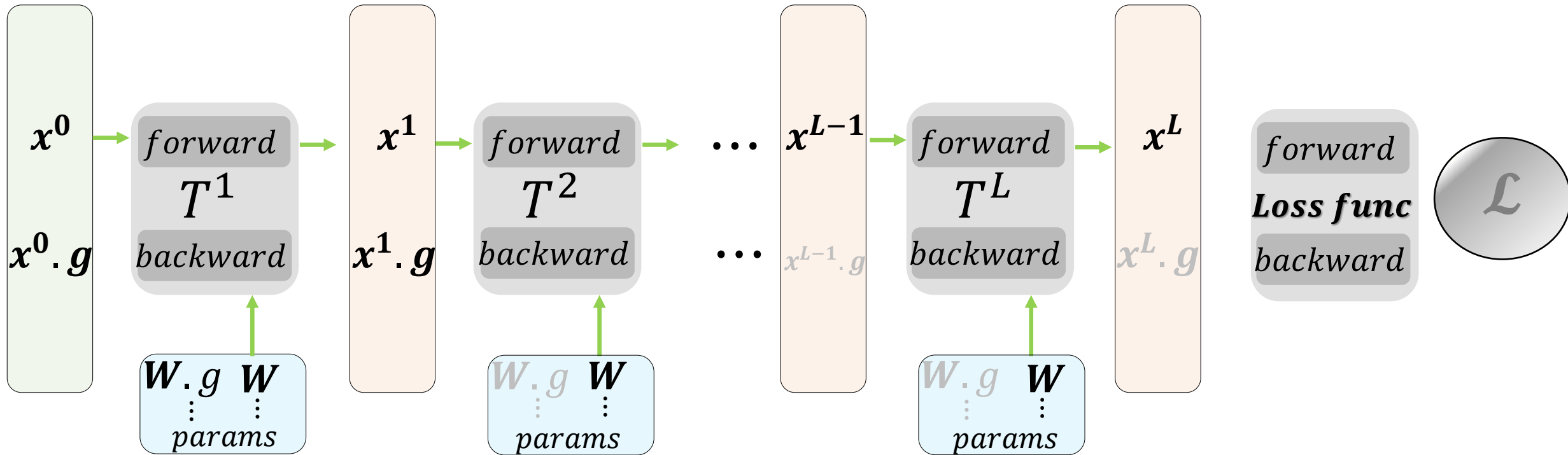
Let's get more generic



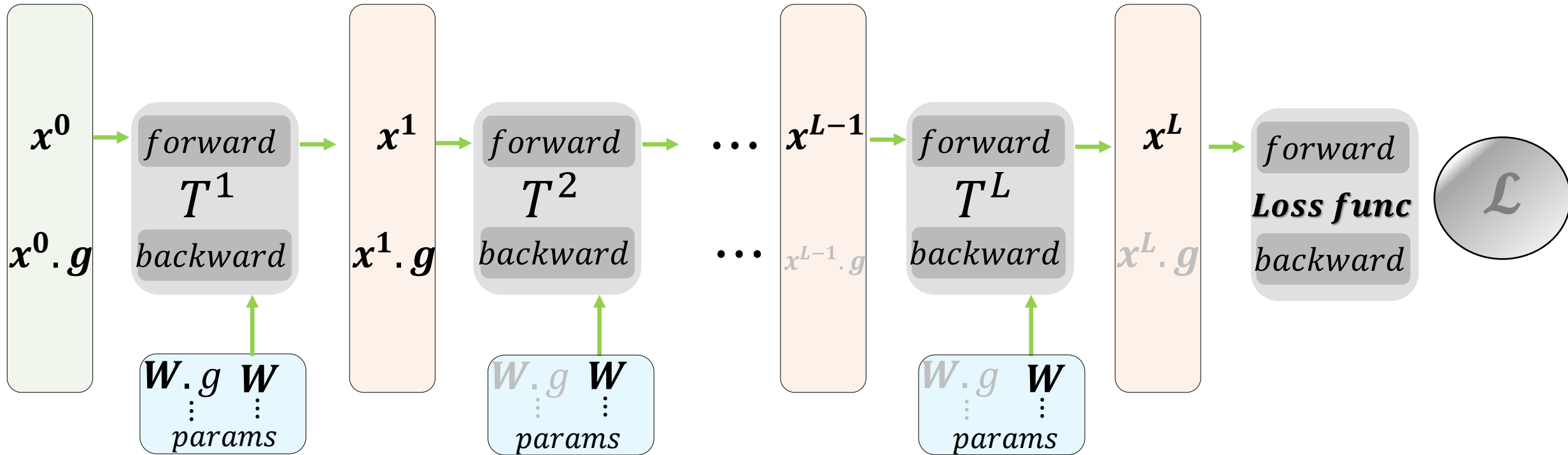
Let's get more generic



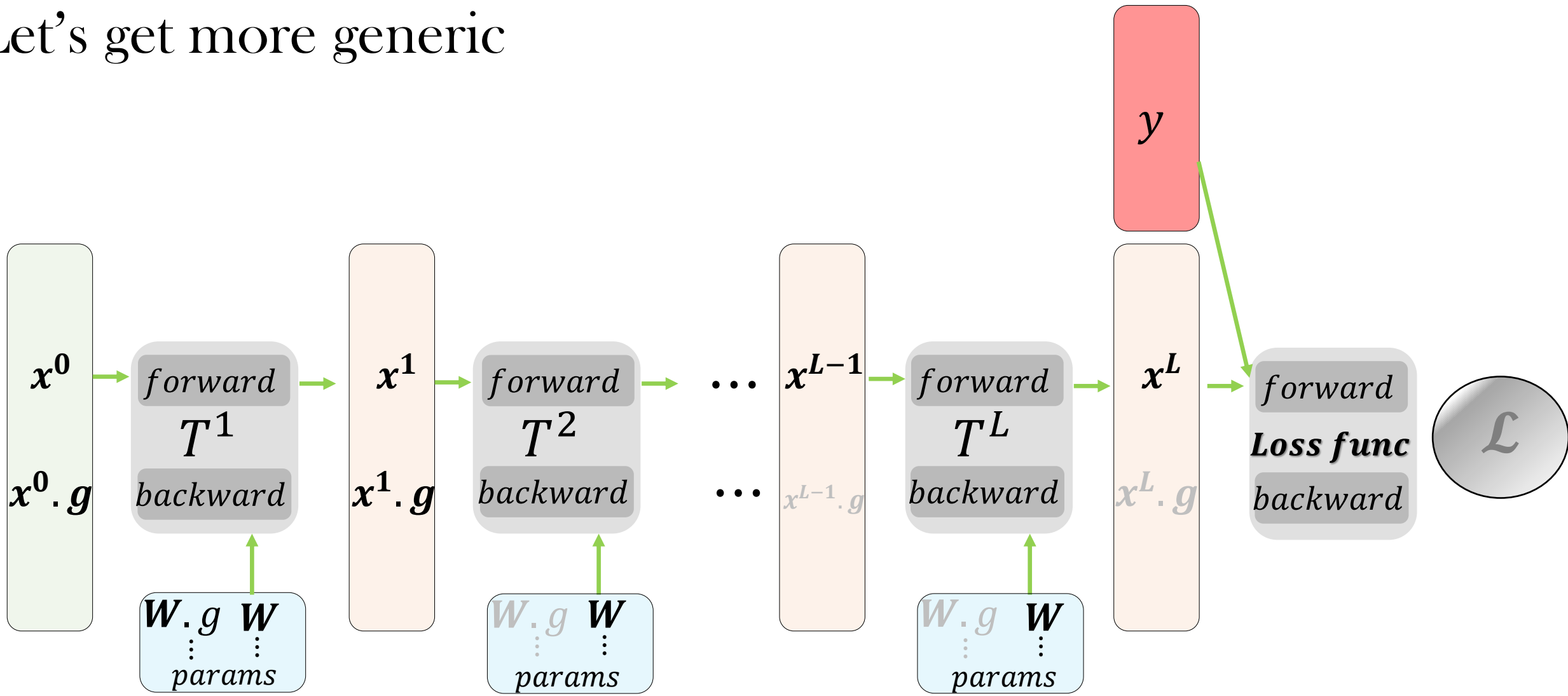
Let's get more generic



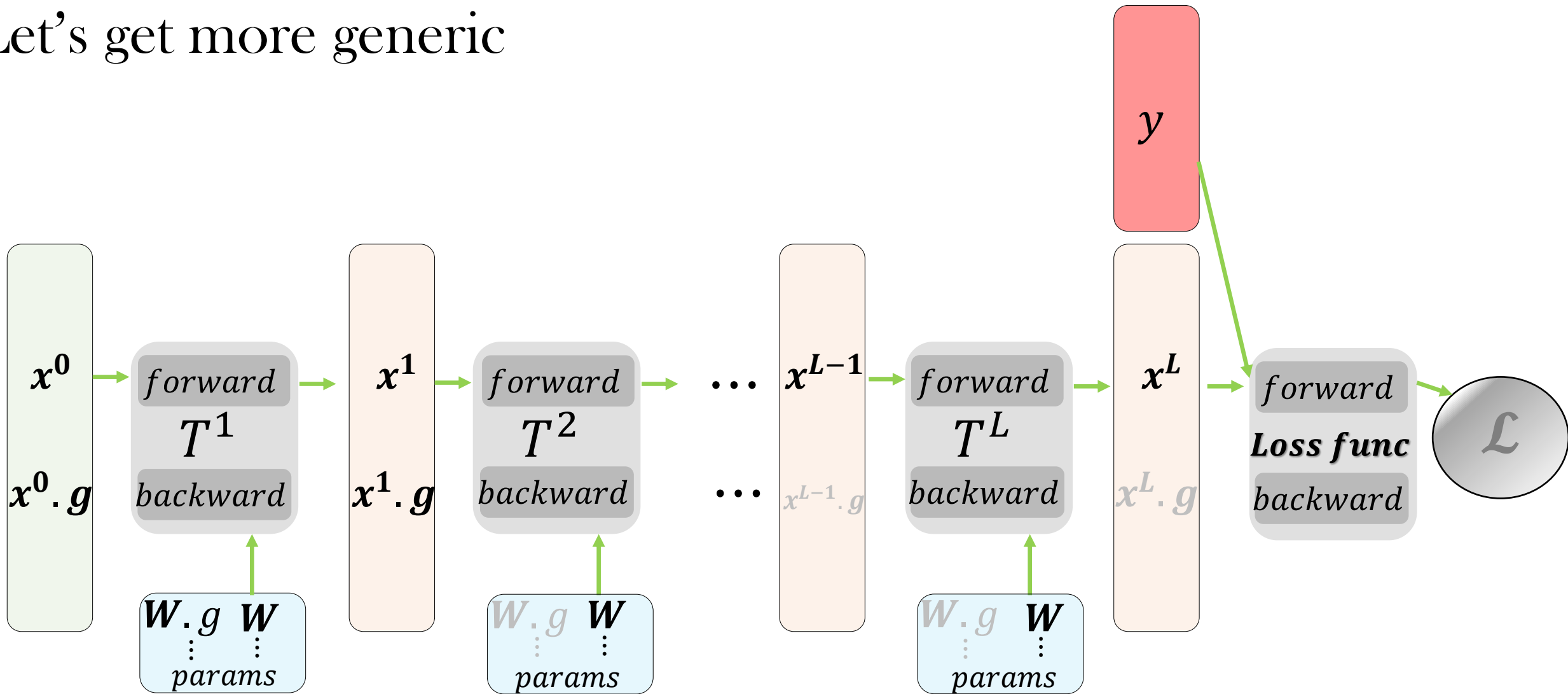
Let's get more generic



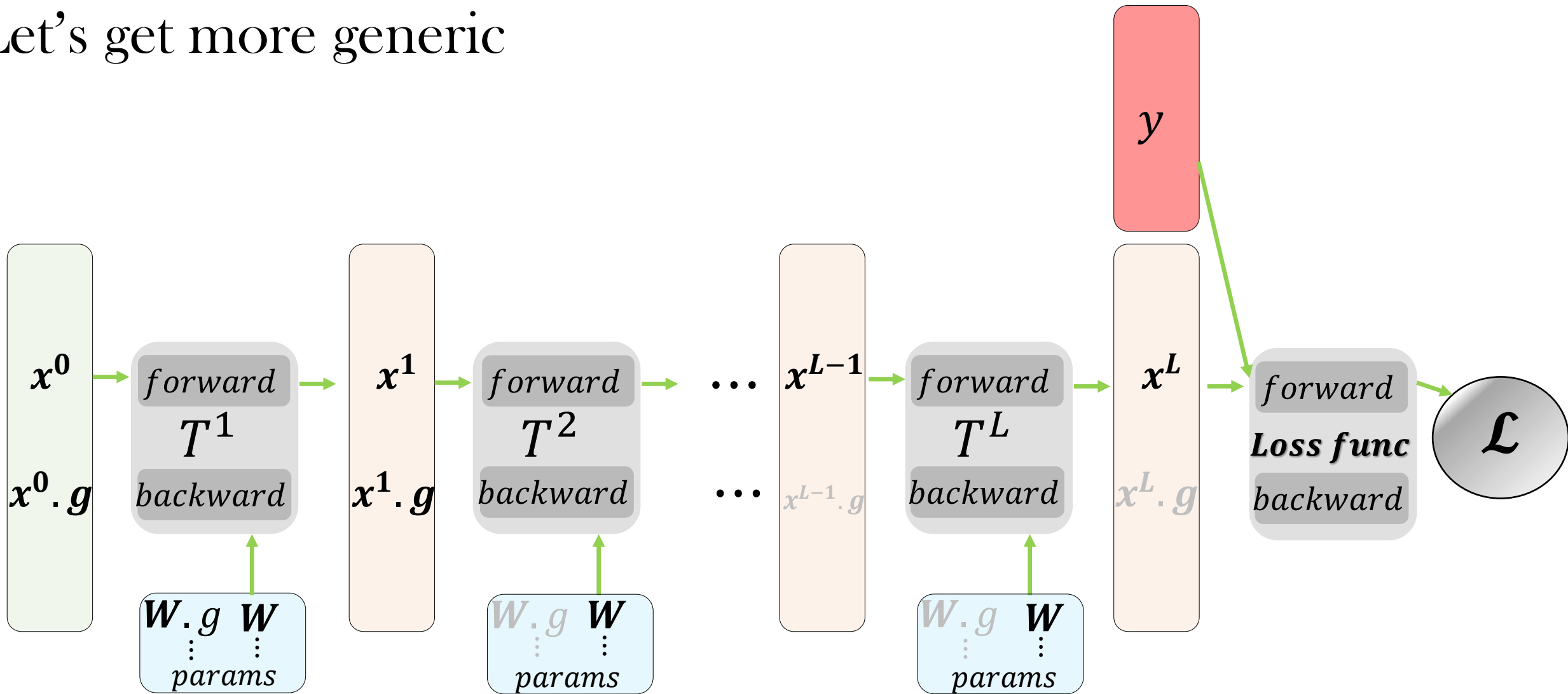
Let's get more generic



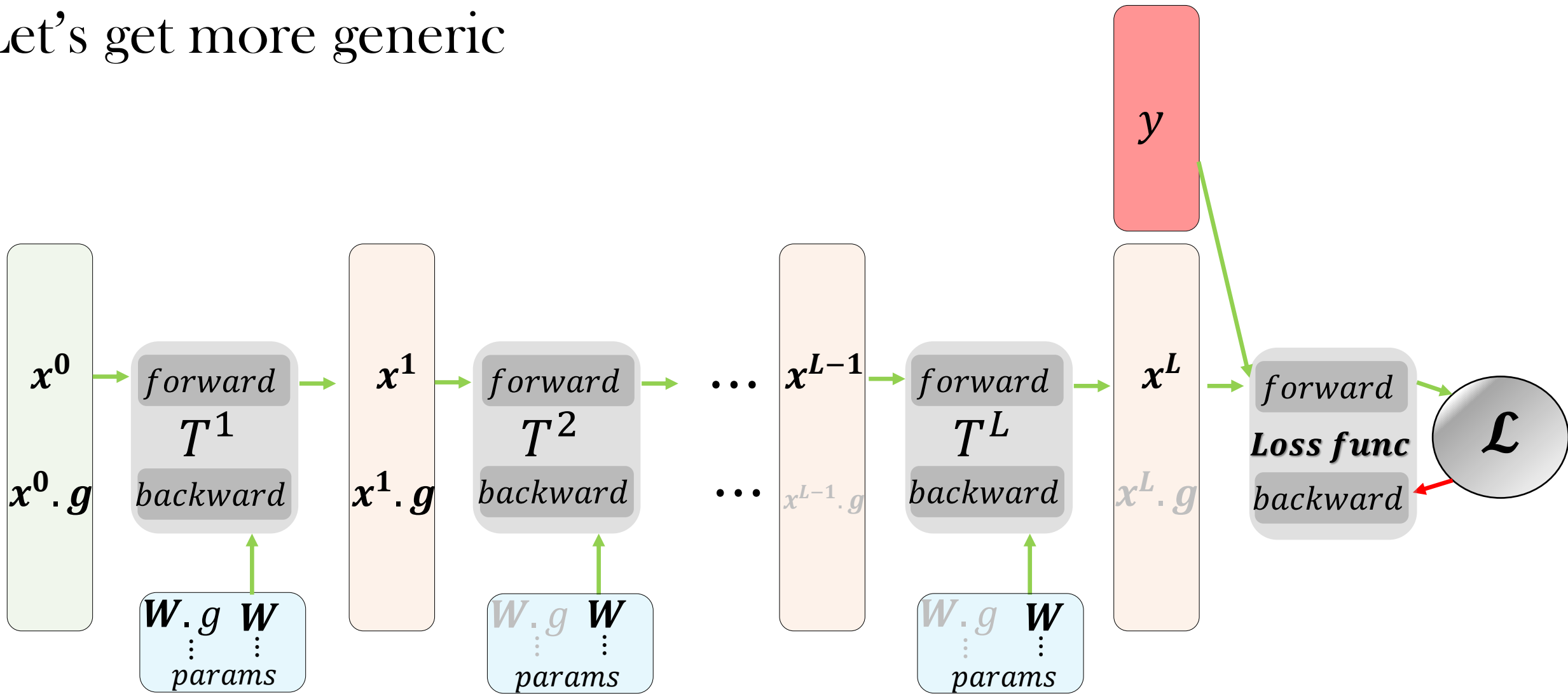
Let's get more generic



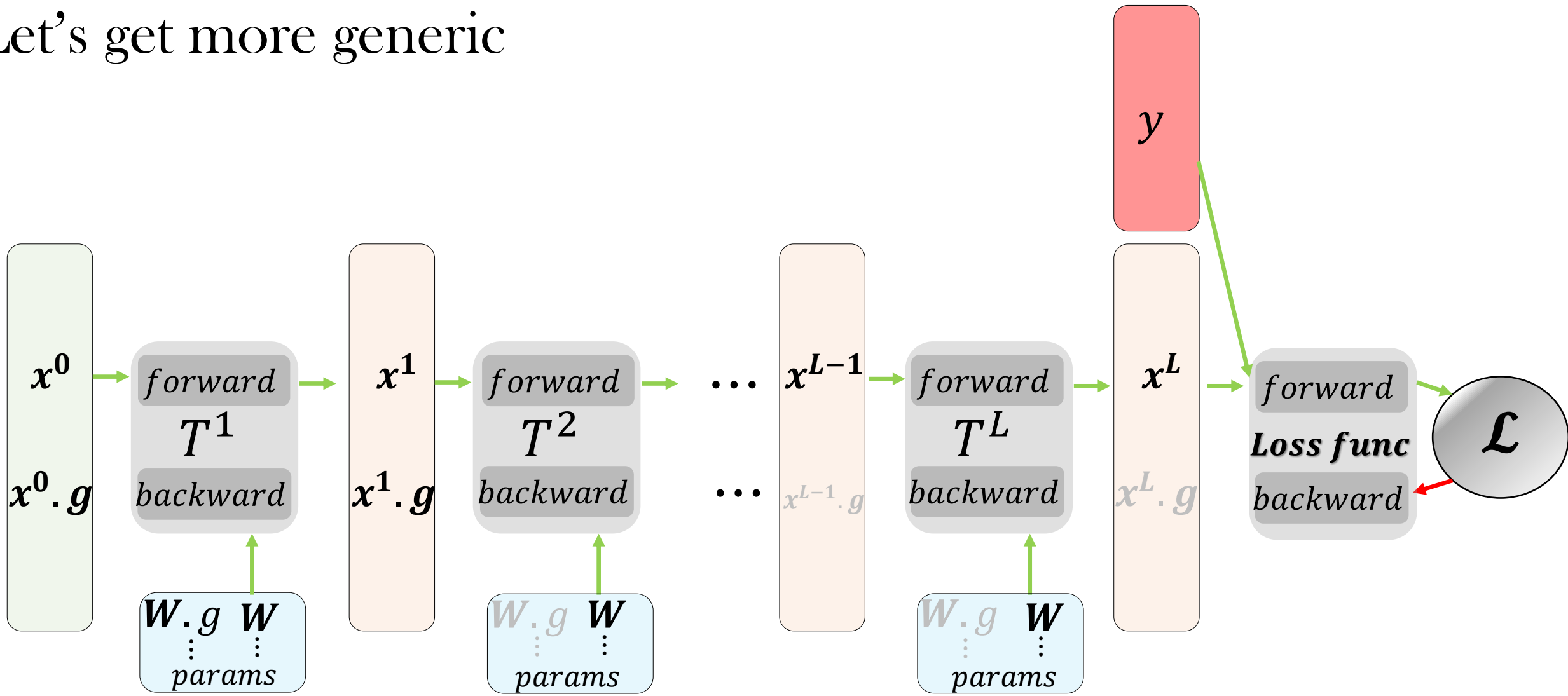
Let's get more generic



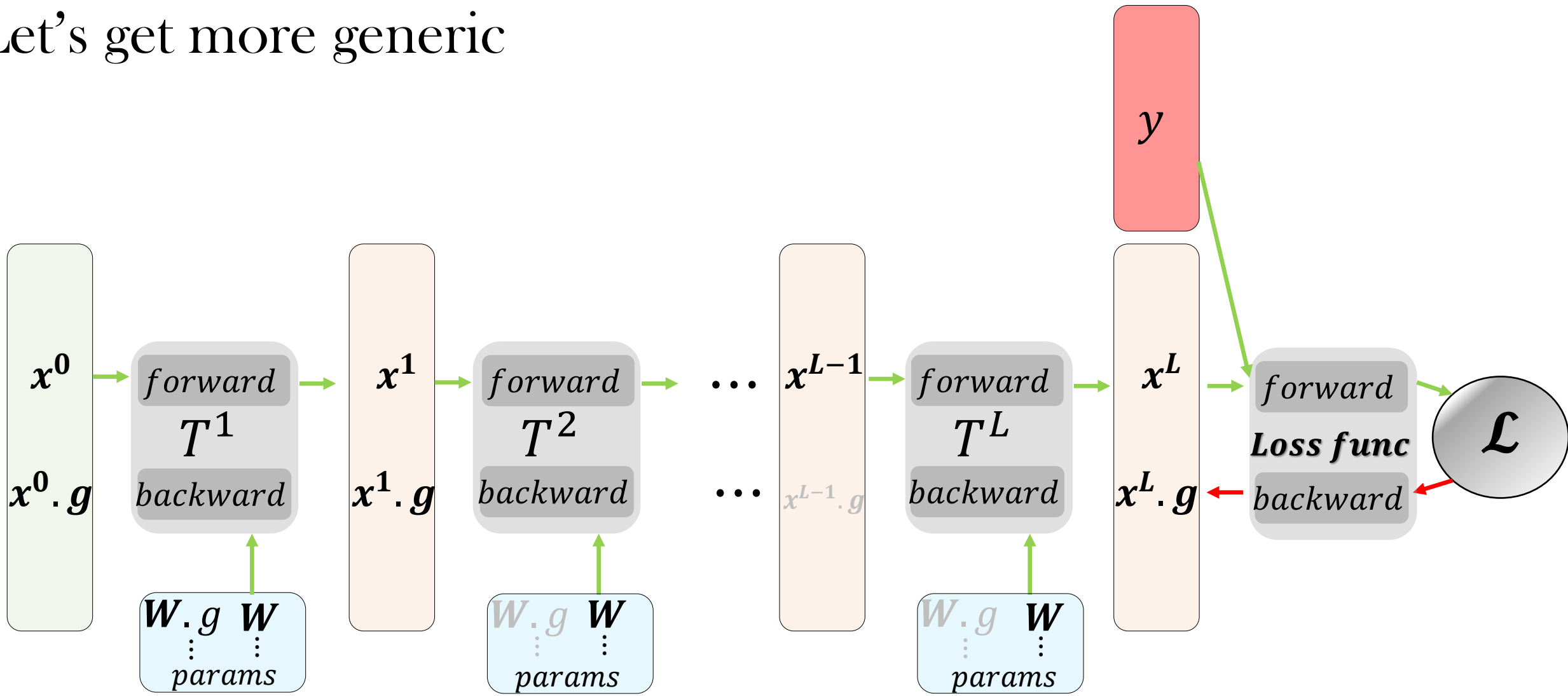
Let's get more generic



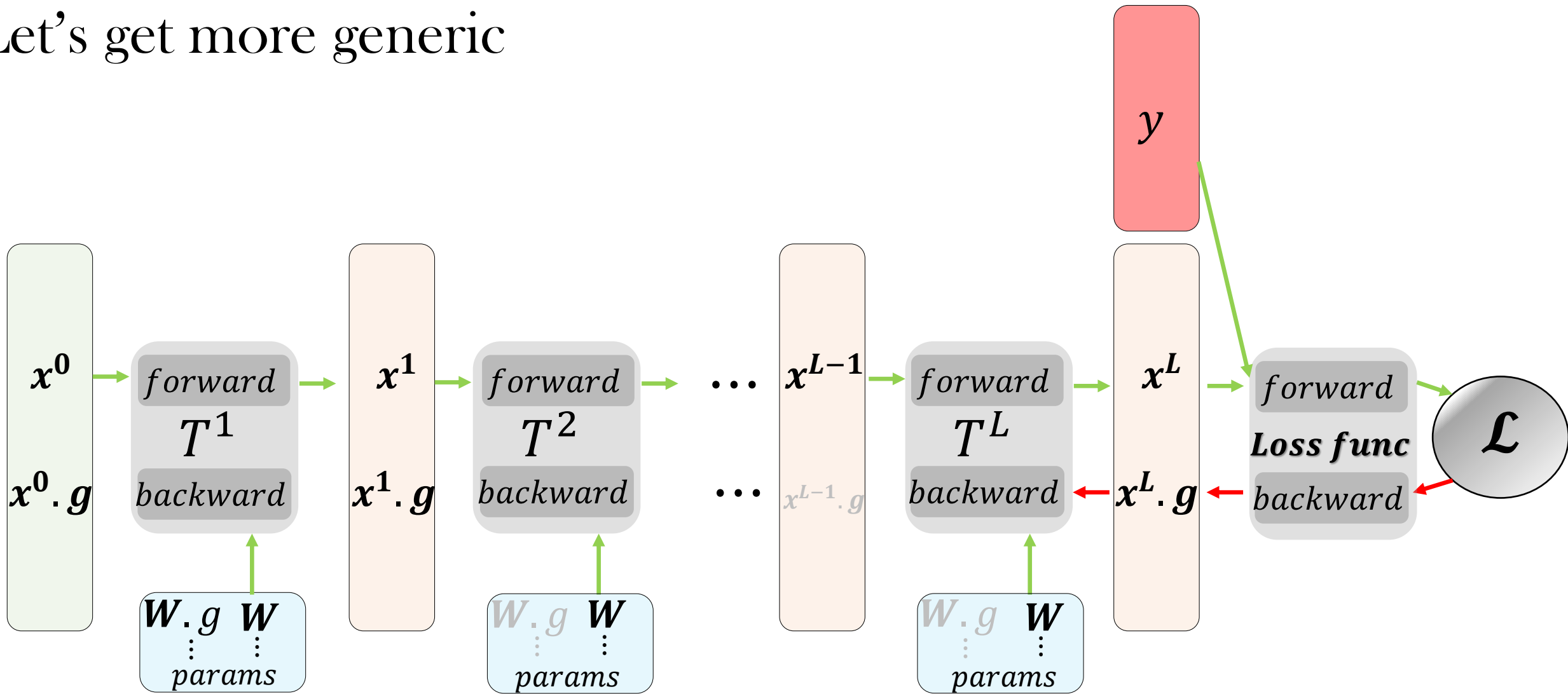
Let's get more generic



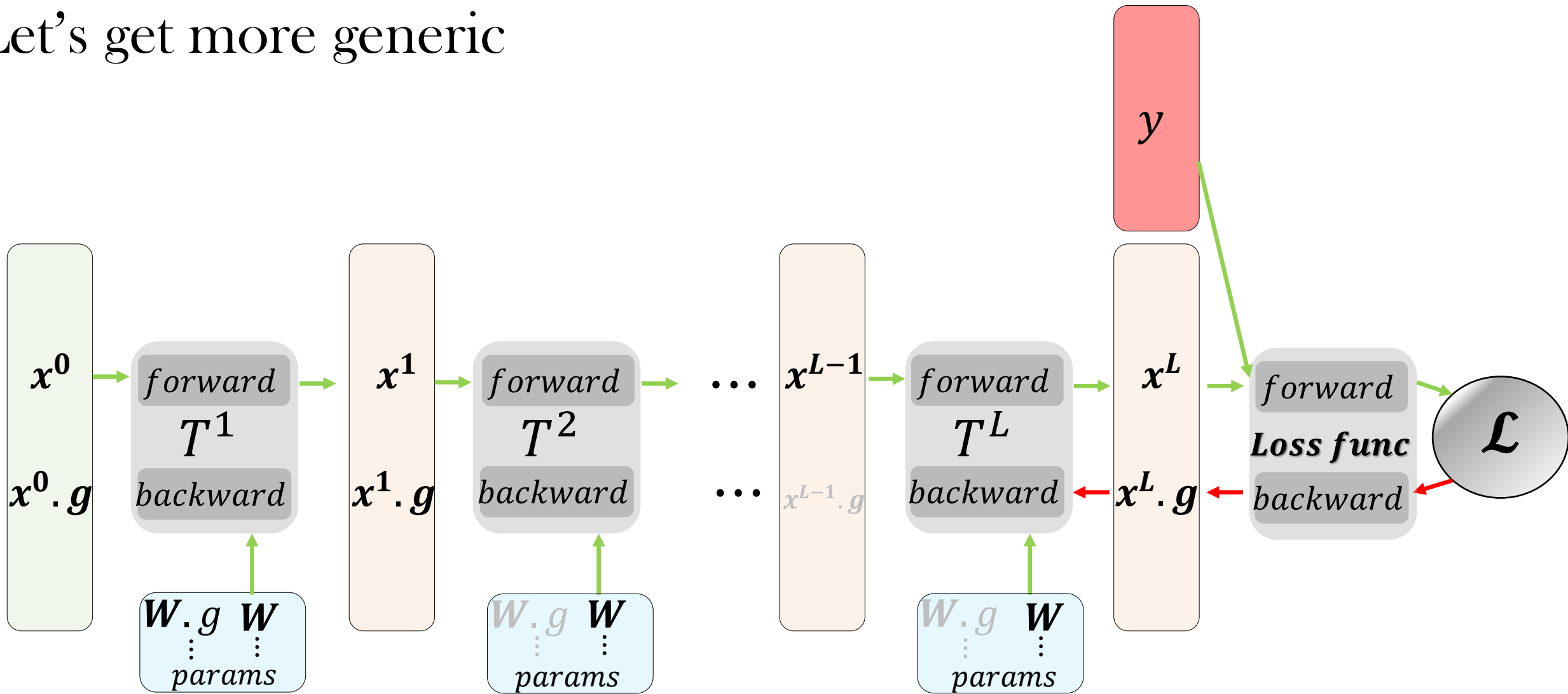
Let's get more generic



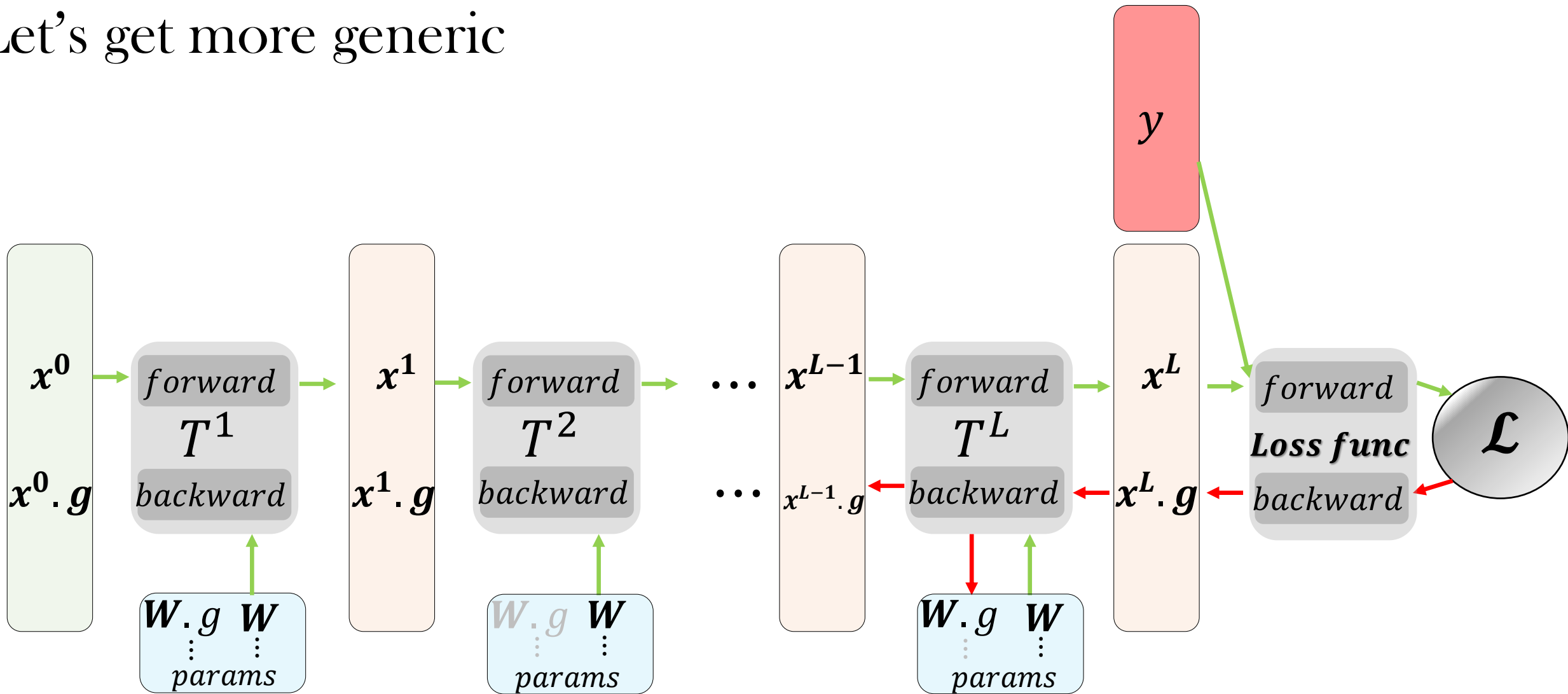
Let's get more generic



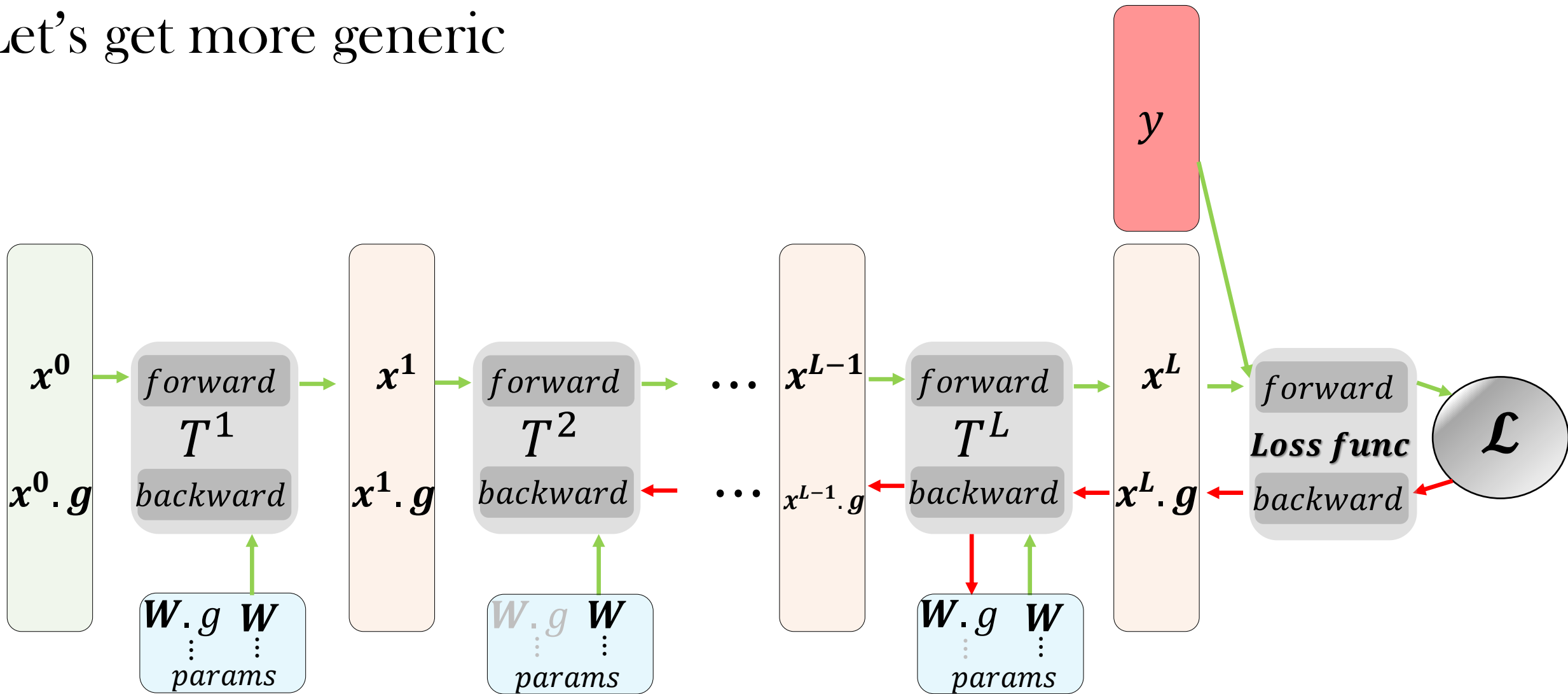
Let's get more generic



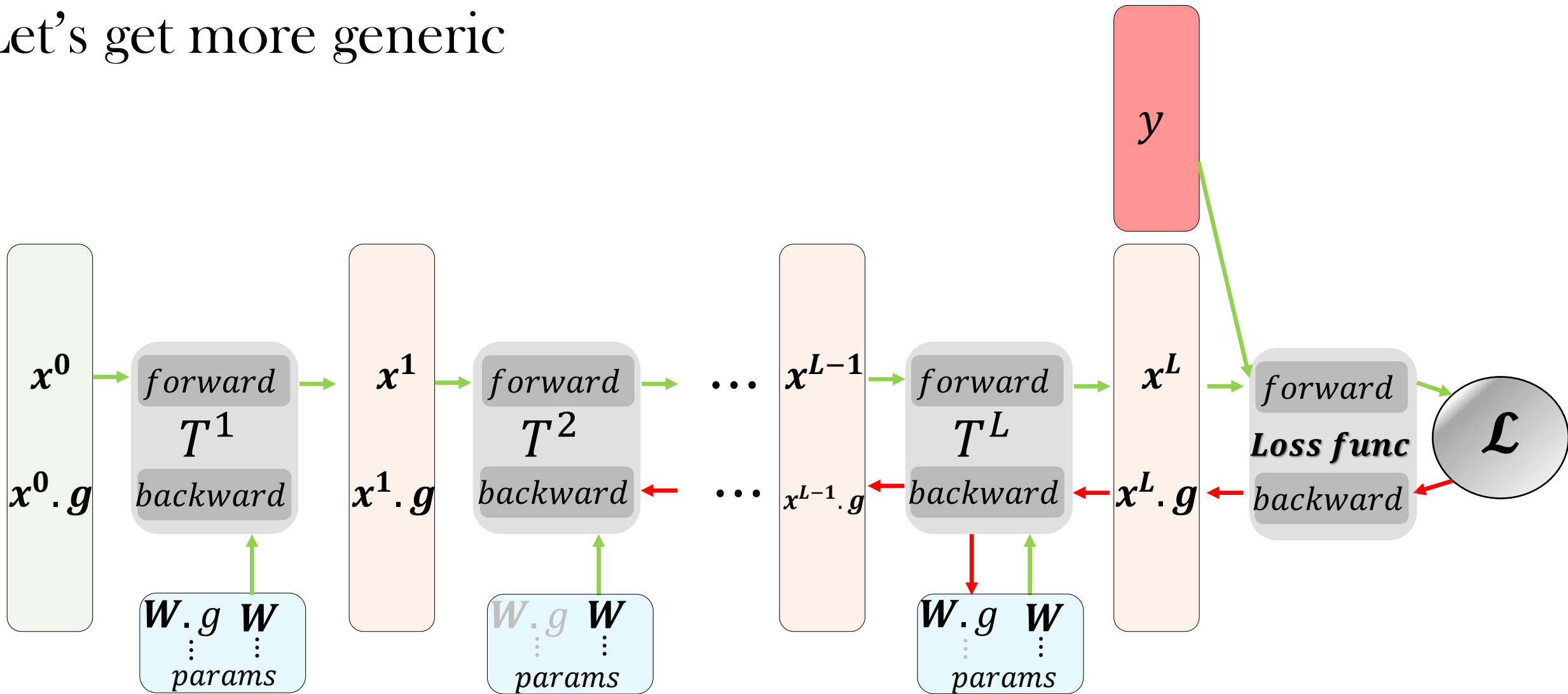
Let's get more generic



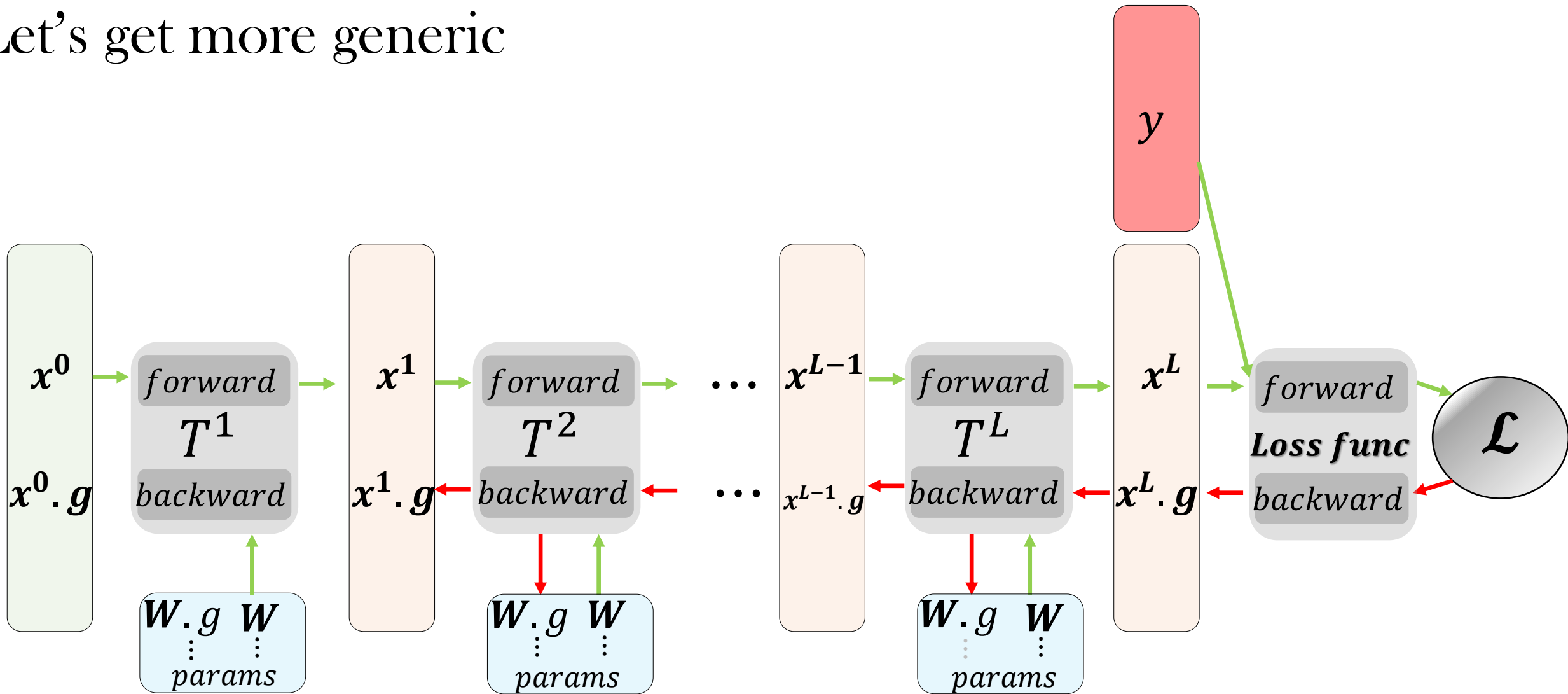
Let's get more generic



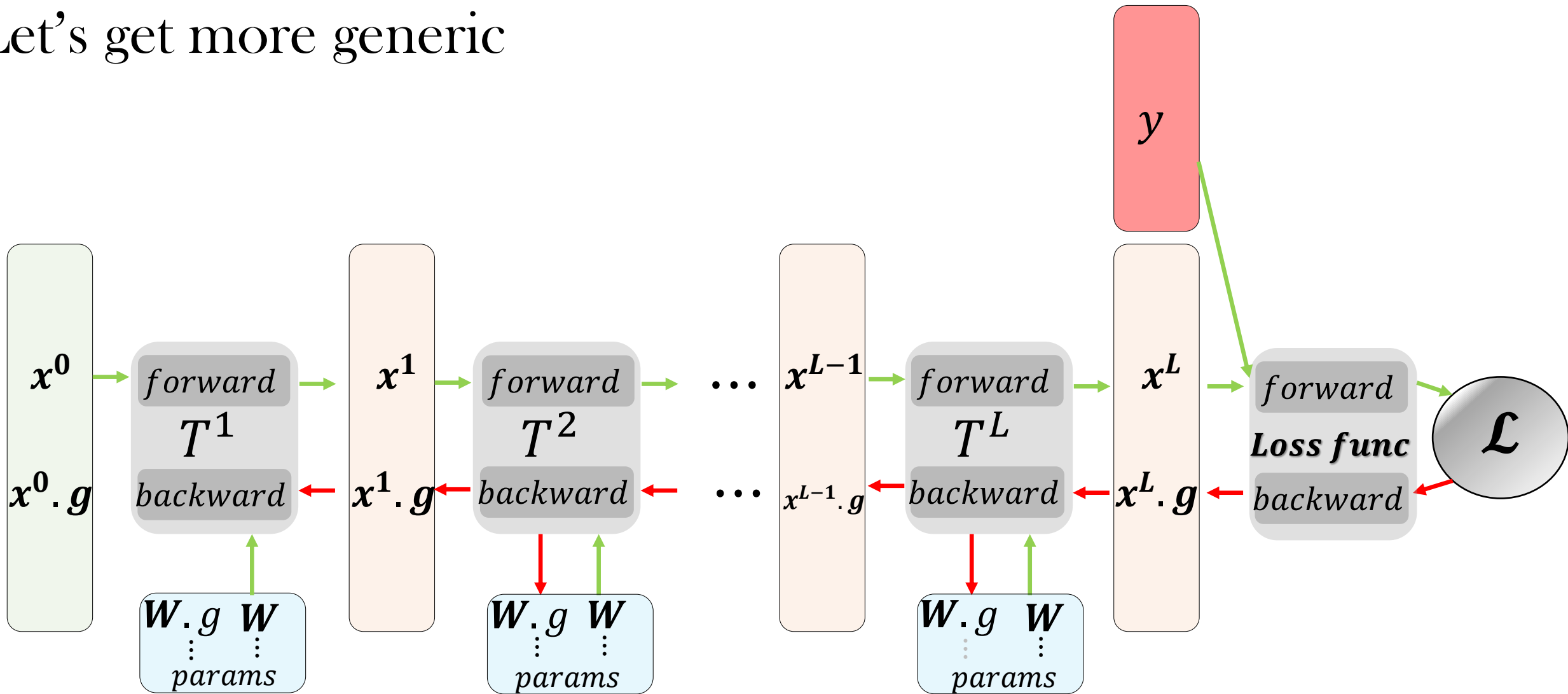
Let's get more generic



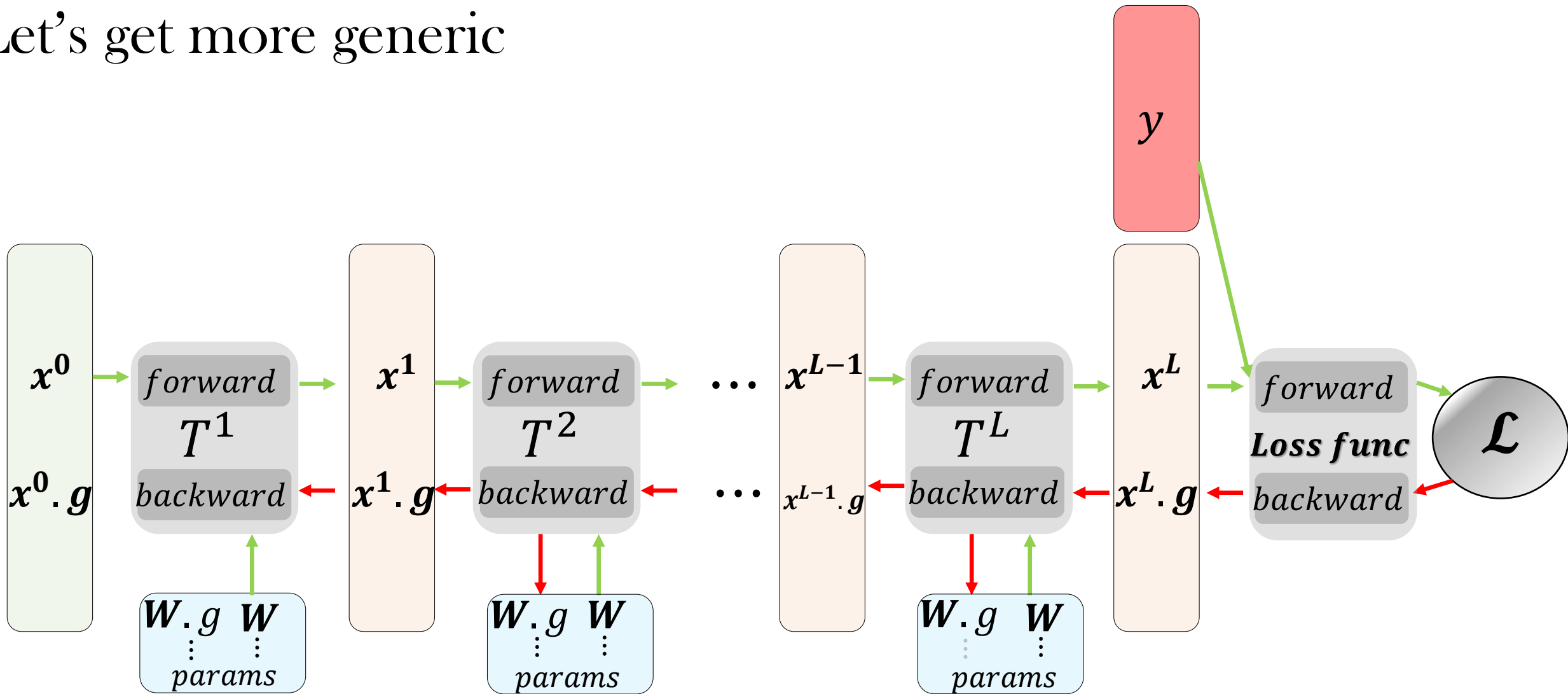
Let's get more generic



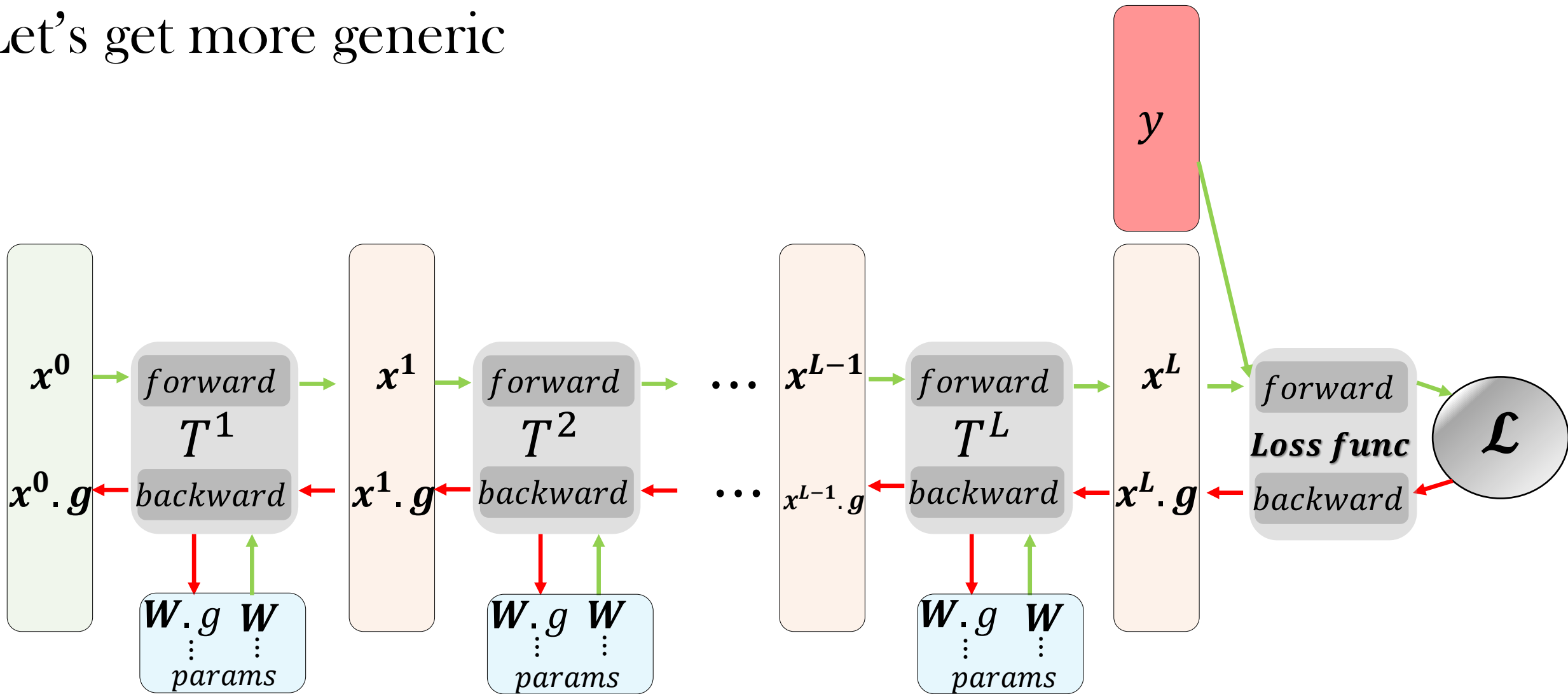
Let's get more generic



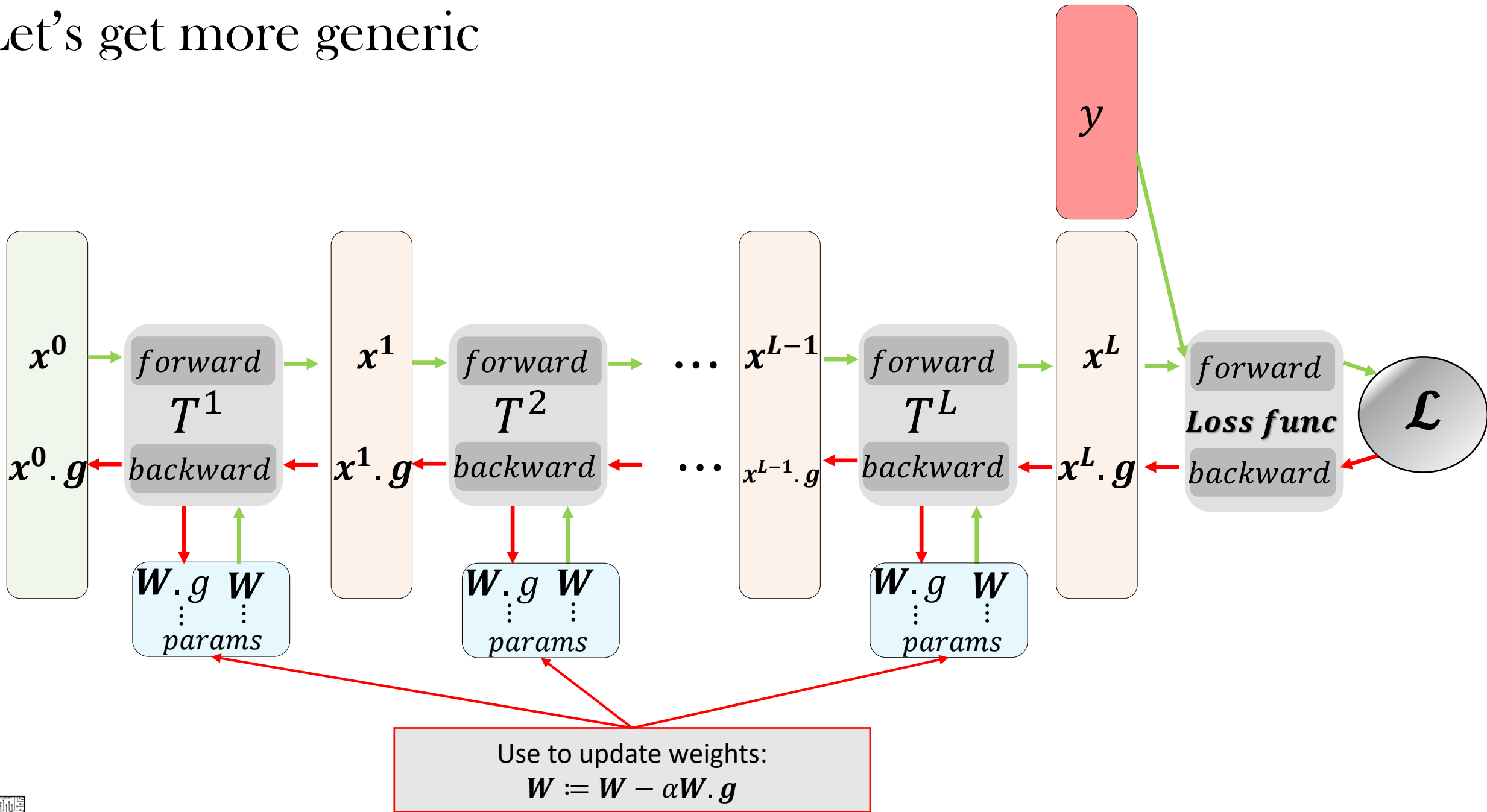
Let's get more generic



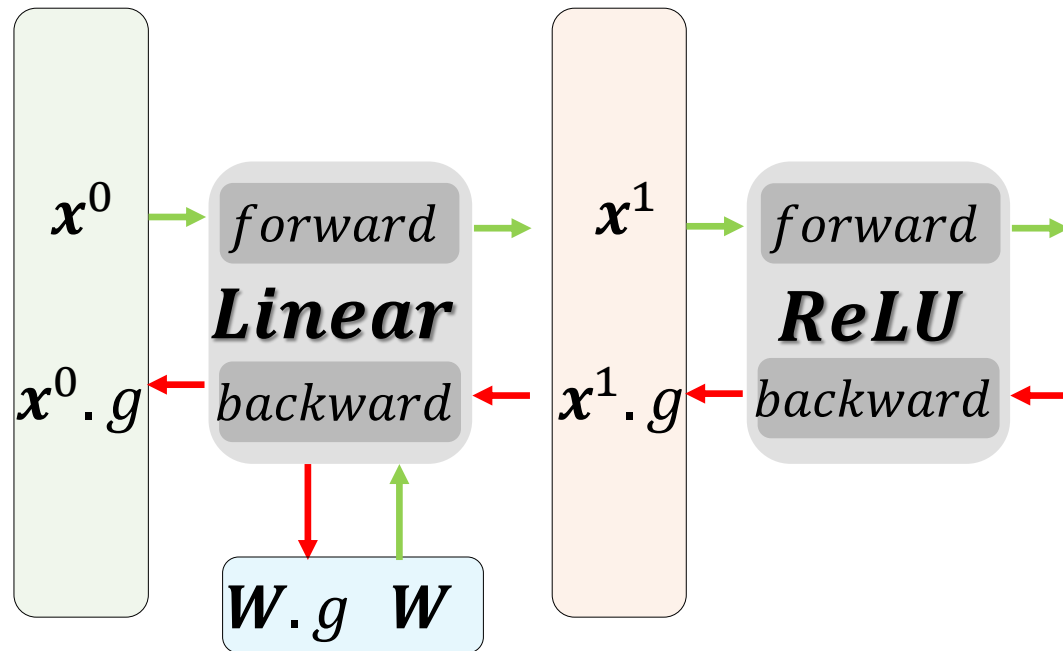
Let's get more generic



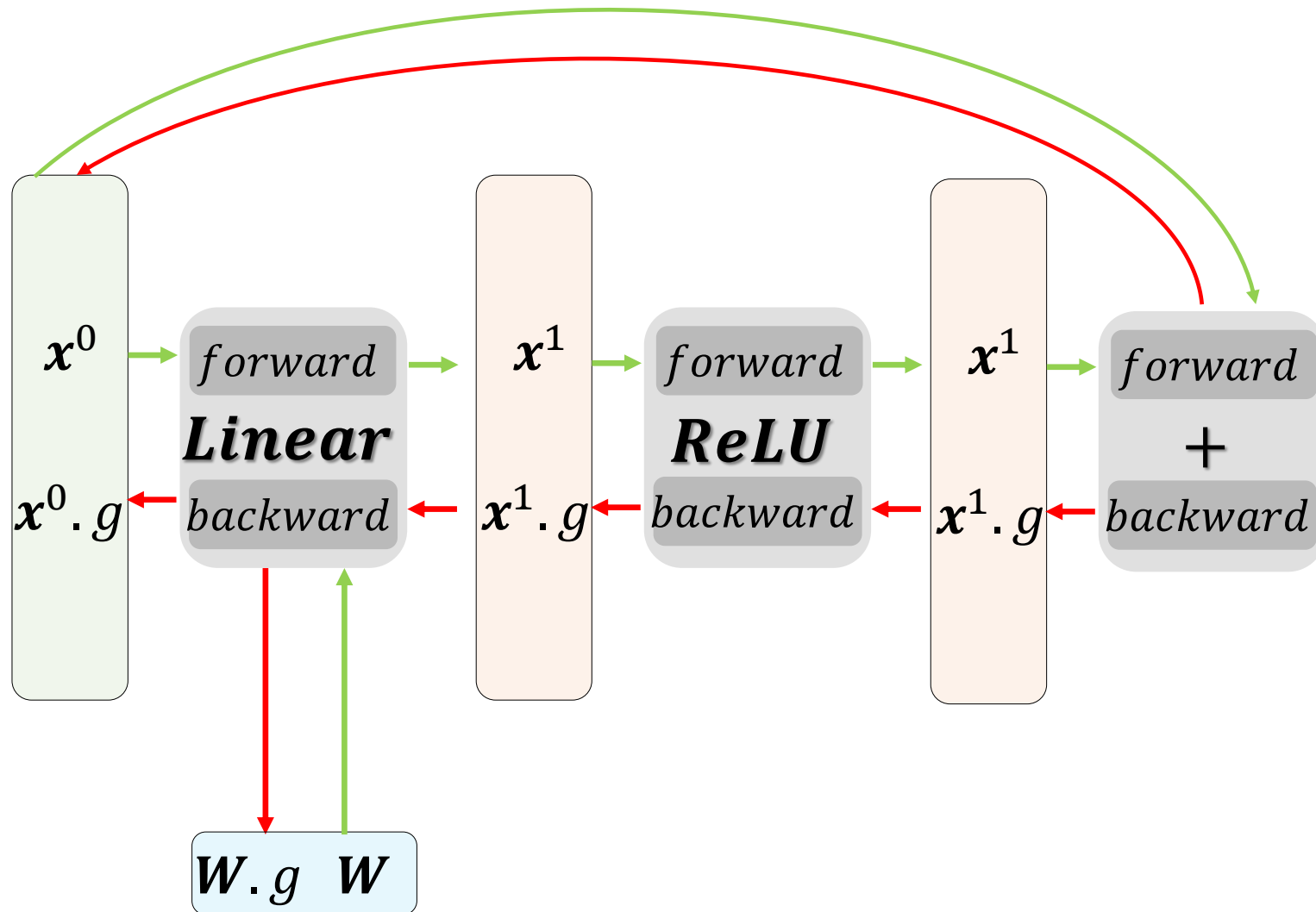
Let's get more generic



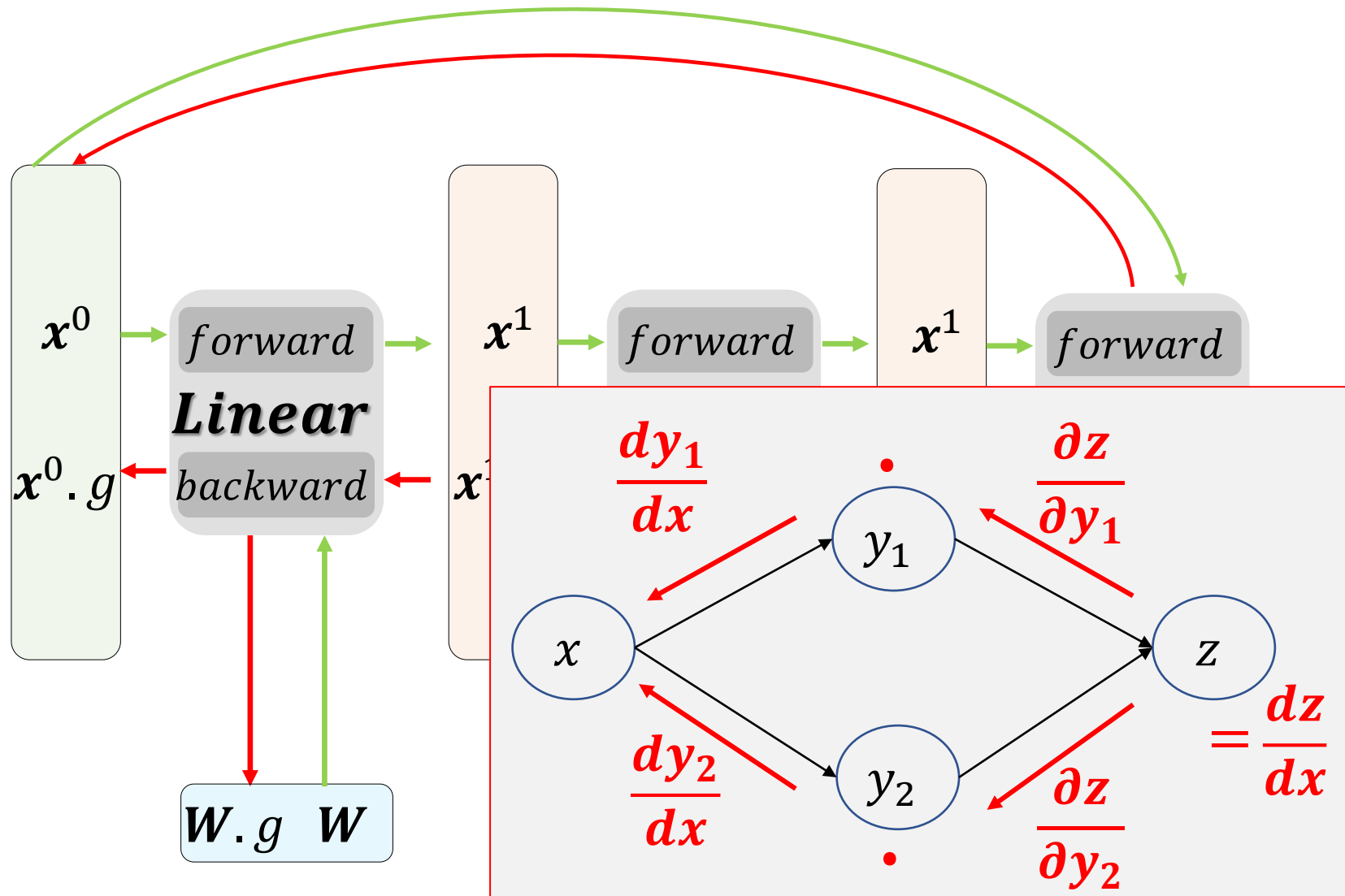
Example: Standard layer



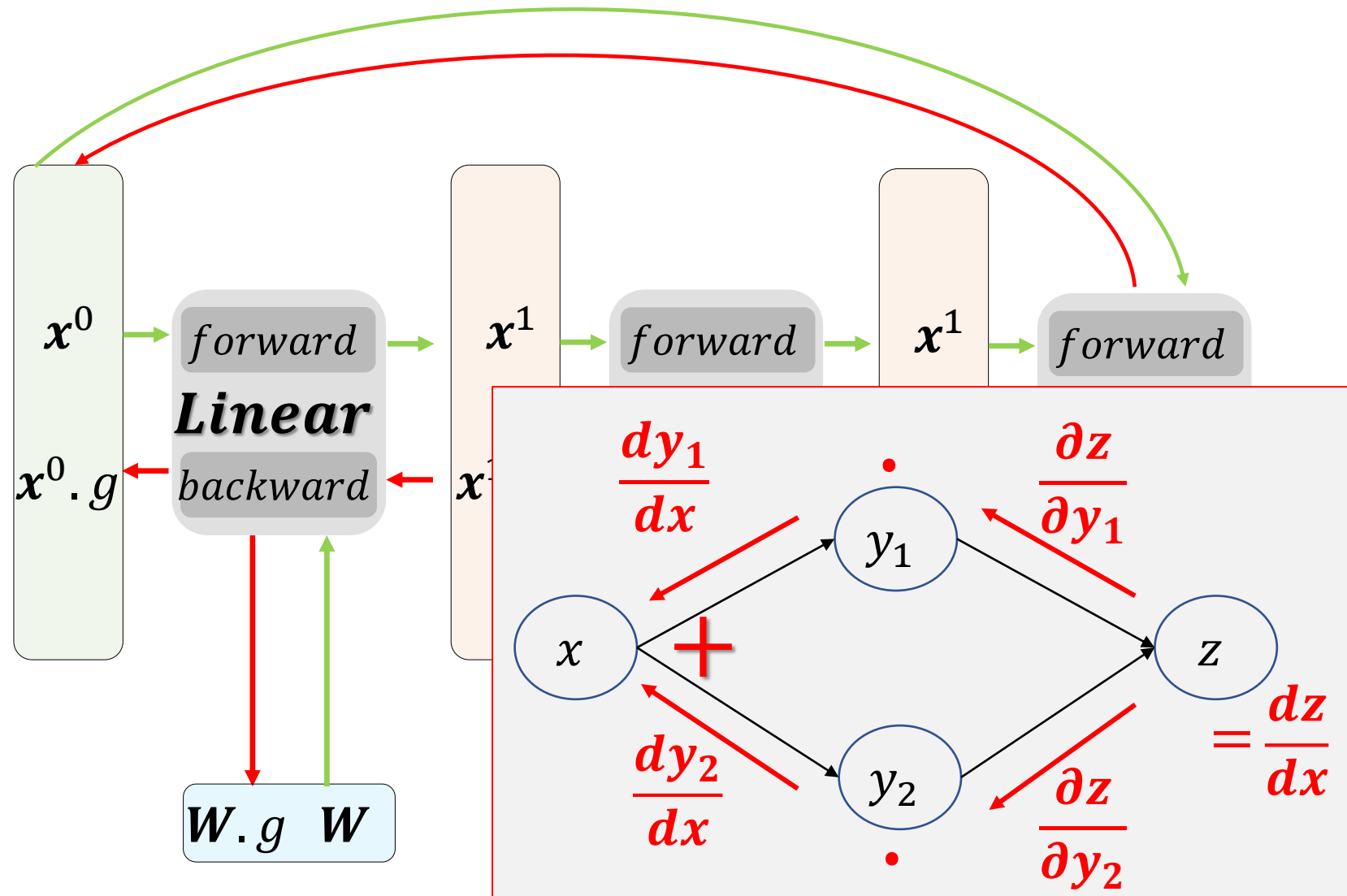
BTW : You can backprop any DAG!



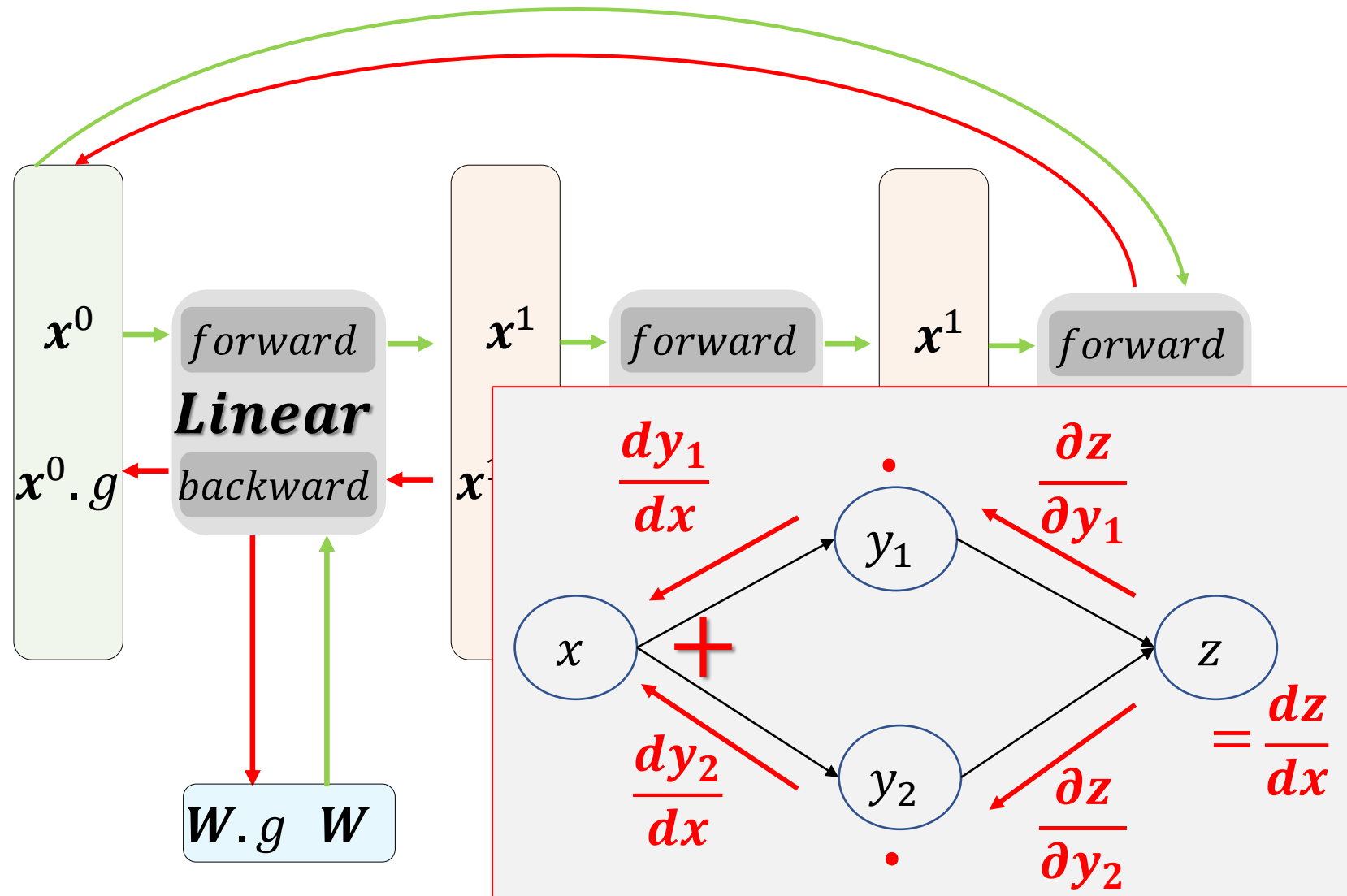
BTW : You can backprop any DAG!



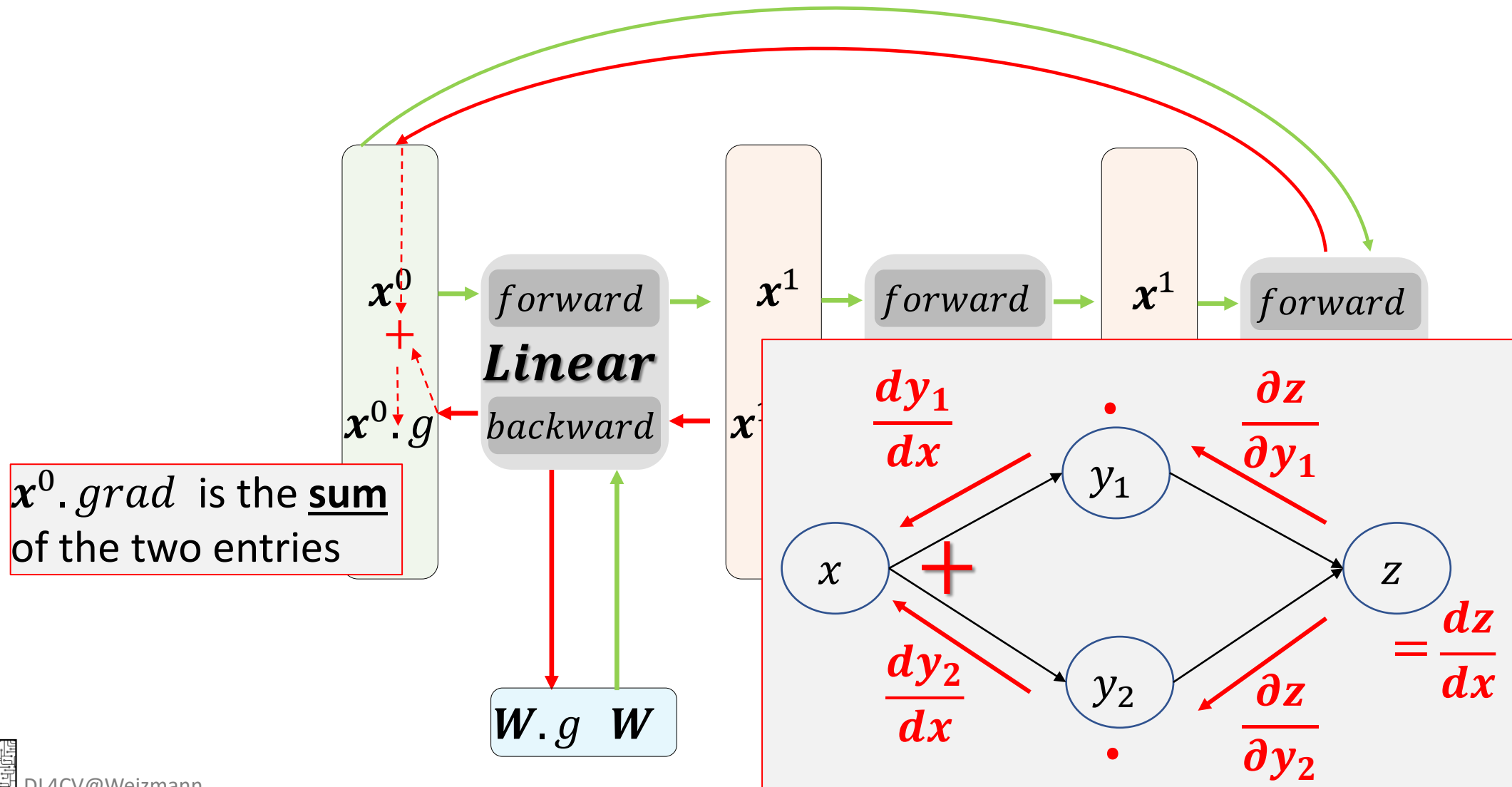
BTW : You can backprop any DAG!



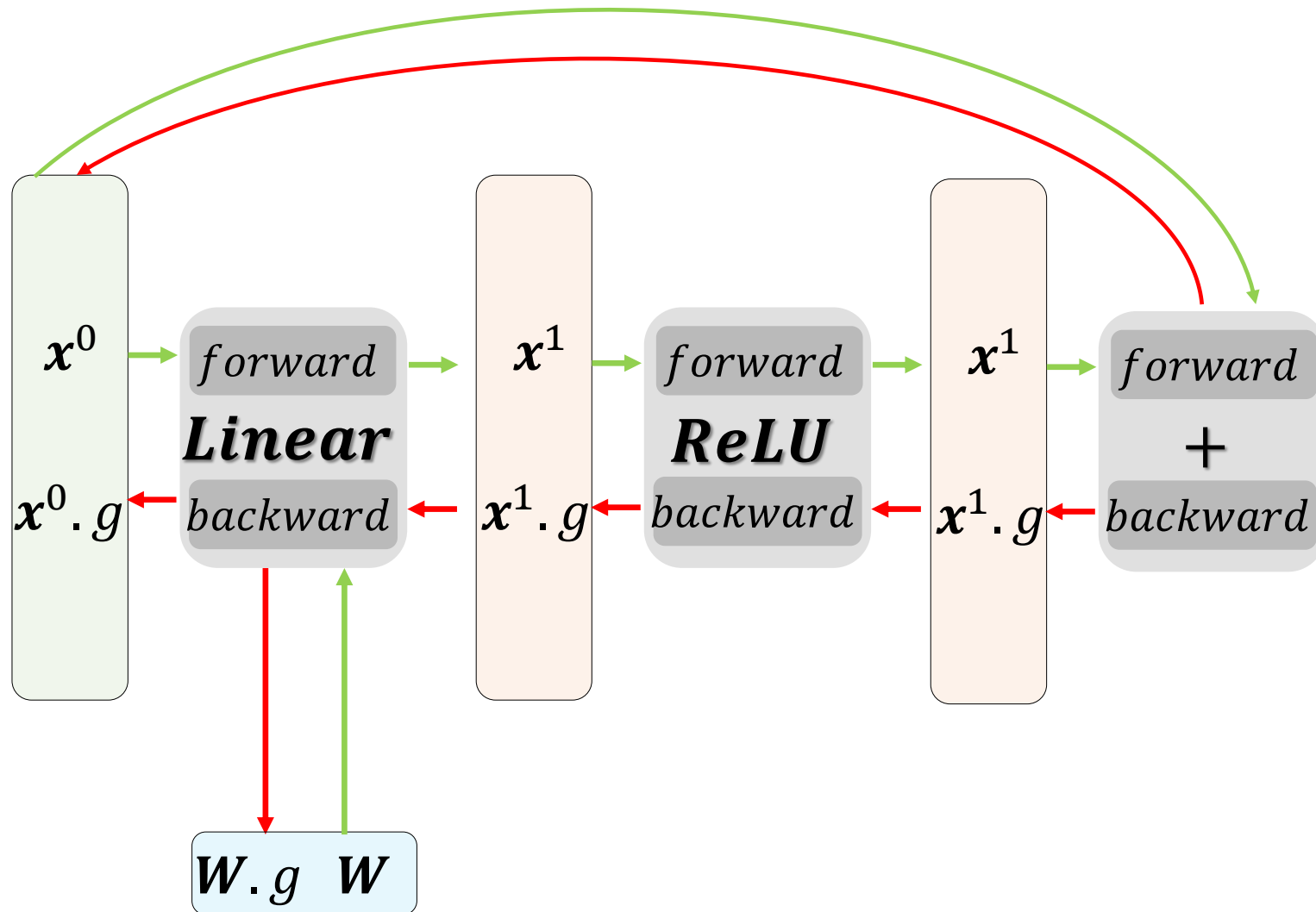
BTW : You can backprop any DAG!



BTW : You can backprop any DAG!

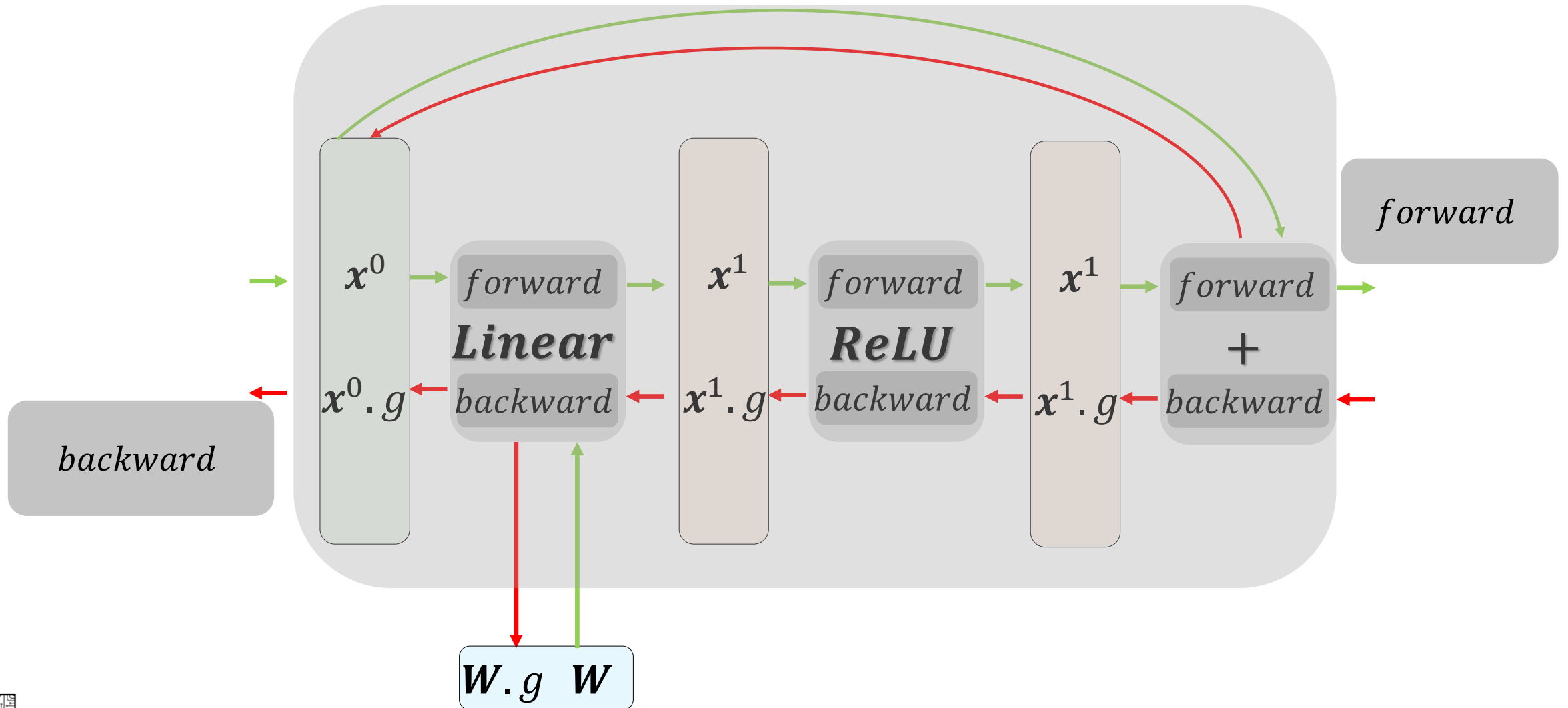


BTW : You can backprop any DAG!



BTW : You can backprop any DAG!

BTW2: Layers (NN modules) can be nested!



Q: Can we use any function inside a network?



Q: Can we use any function inside a network?

- Has to be differentiable with respect to input and params

Q: Can we use any function inside a network?

- Has to be differentiable with respect to input and params
- Implications to optimization can be deadly:

Yes you should understand backprop



Andrej Karpathy Dec 19, 2016 · 7 min read

Q: Can we use any function inside a network?

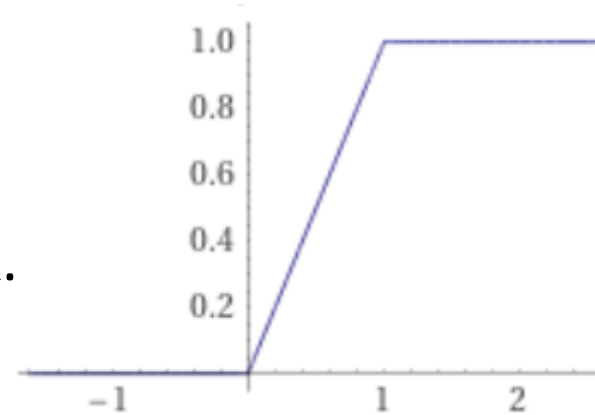
- Has to be differentiable with respect to input and params
- Implications to optimization can be deadly:

Q: Can we use any function inside a network?

- Has to be differentiable with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.

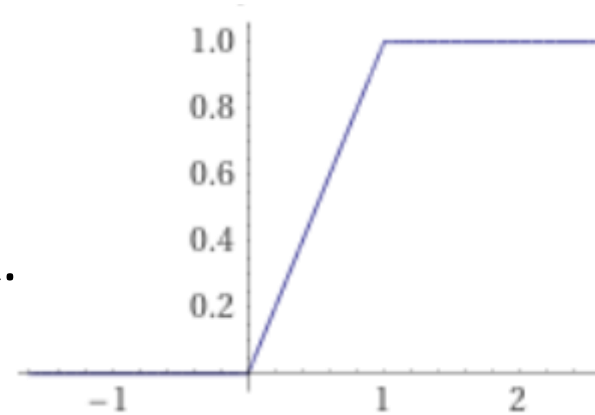
Q: Can we use any function inside a network?

- Has to be differentiable with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.



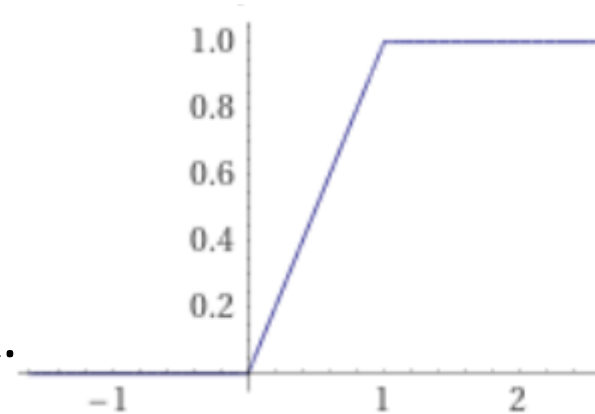
Q: Can we use any function inside a network?

- Has to be differentiable with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.
 - Ideas how to solve this?



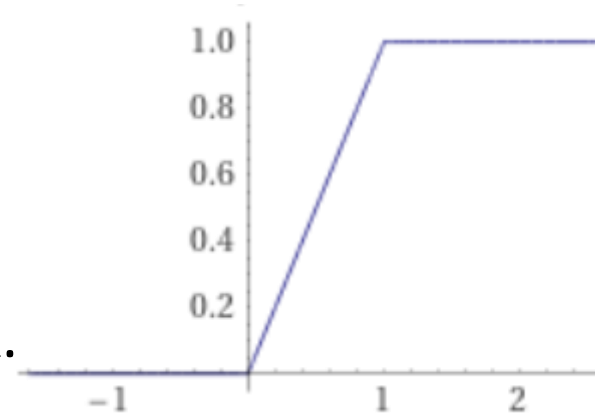
Q: Can we use any function inside a network?

- Has to be differentiable (L_2 norm)
- L_2 with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.
 - Ideas how to solve this?
 - Square root loss (L_1 norm)

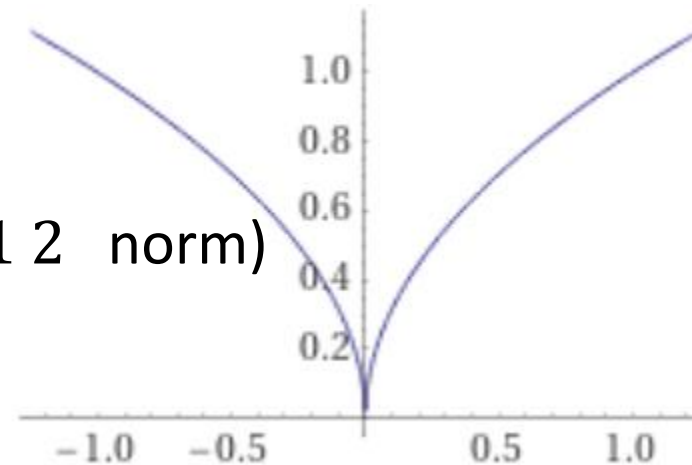


Q: Can we use any function inside a network?

- Has to be differentiable (L_2 norm)
- L_2 with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.
 - Ideas how to solve this?

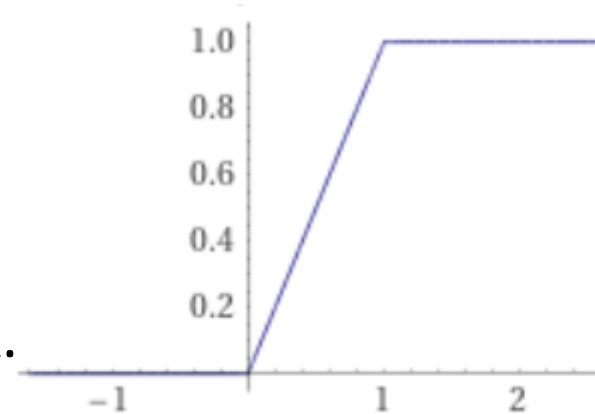


- Square root loss ($L_{1/2}$ norm)

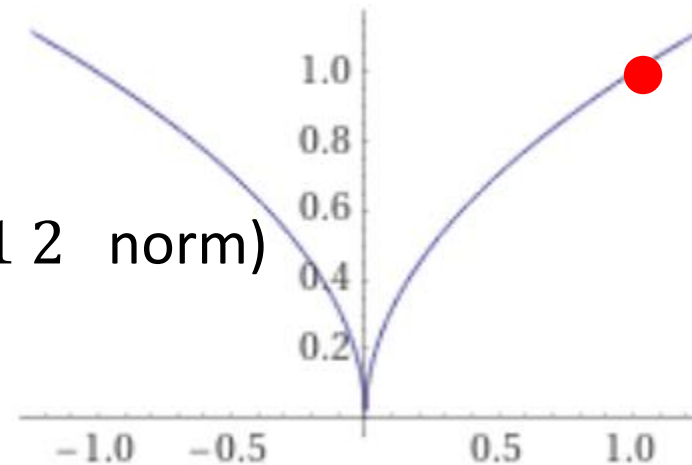


Q: Can we use any function inside a network?

- Has to be differentiable (L_2 norm)
- L_2 with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.
 - Ideas how to solve this?

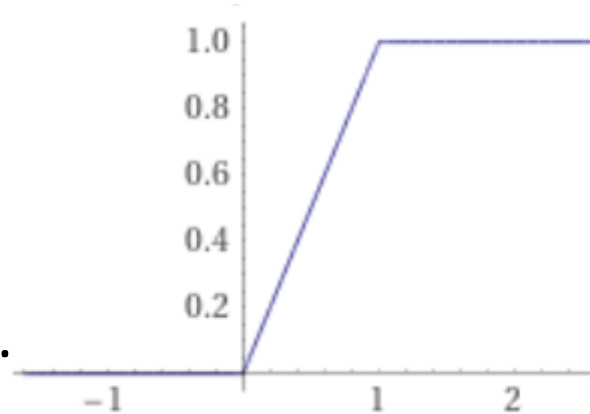


- Square root loss ($L_{1/2}$ norm)

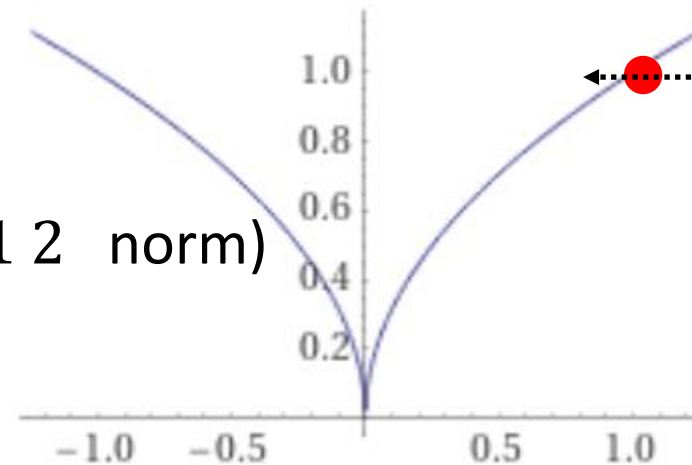


Q: Can we use any function inside a network?

- Has to be differentiable (L_2 norm)
- L_2 with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.
 - Ideas how to solve this?

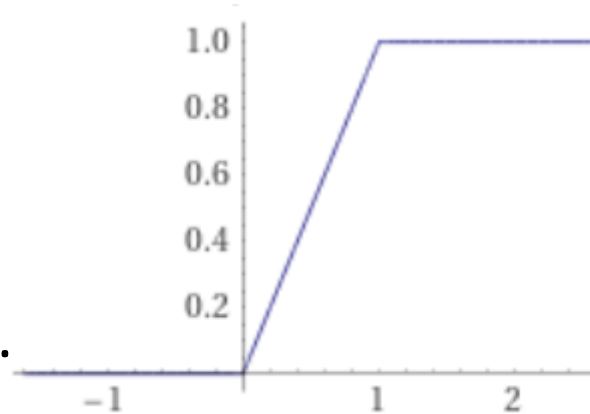


- Square root loss ($L_{1/2}$ norm)

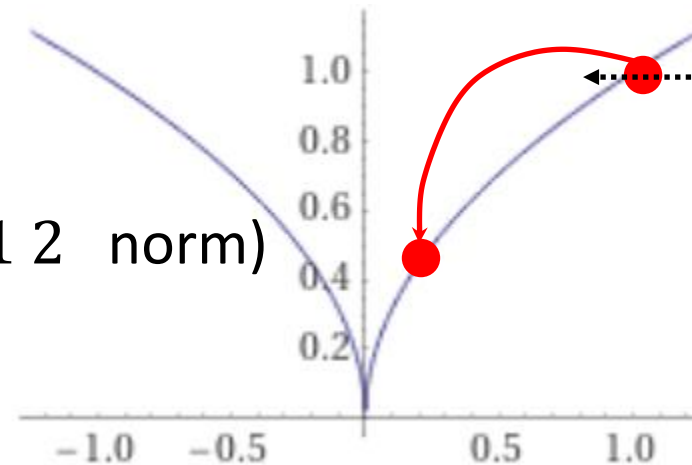


Q: Can we use any function inside a network?

- Has to be differentiable (L_2 norm)
- L_2 with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.
 - Ideas how to solve this?

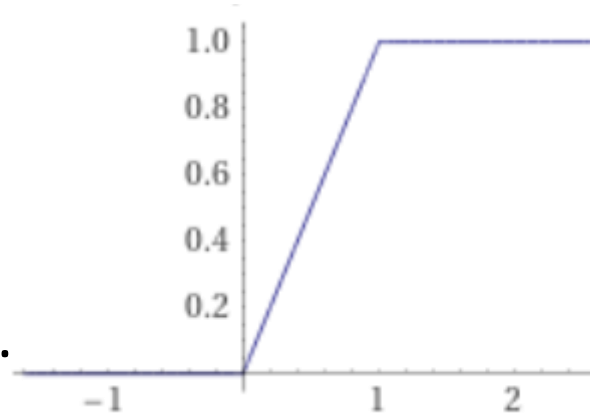


- Square root loss ($L_{1/2}$ norm)

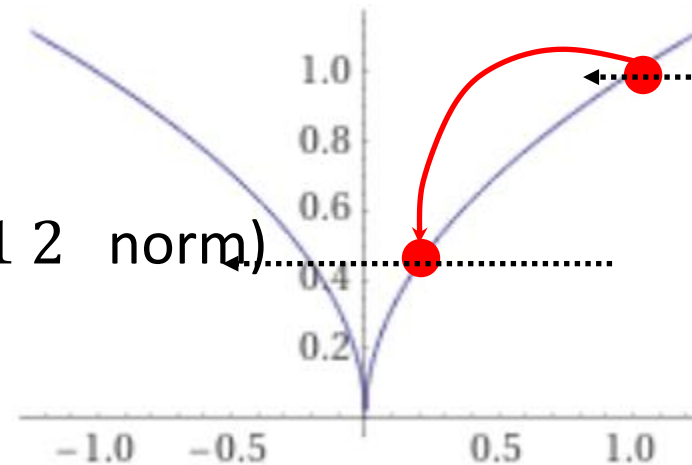


Q: Can we use any function inside a network?

- Has to be differentiable (L_2 norm)
- L_2 with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.
 - Ideas how to solve this?

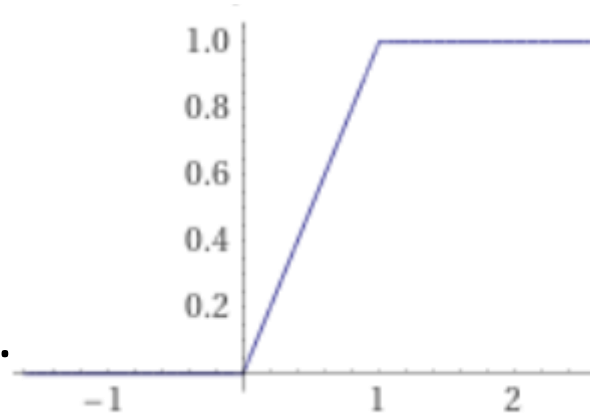


- Square root loss ($L_{1/2}$ norm)

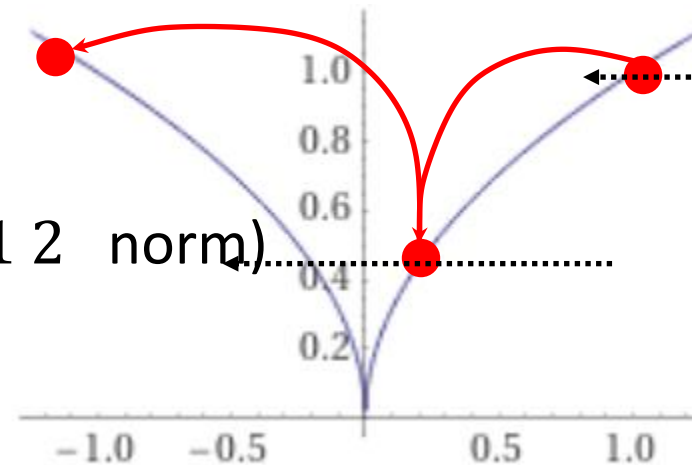


Q: Can we use any function inside a network?

- Has to be differentiable (L_2 norm)
- L_2 with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.
 - Ideas how to solve this?

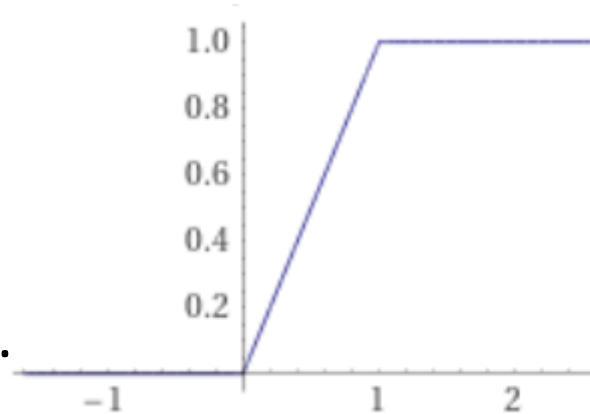


- Square root loss ($L_{1/2}$ norm)

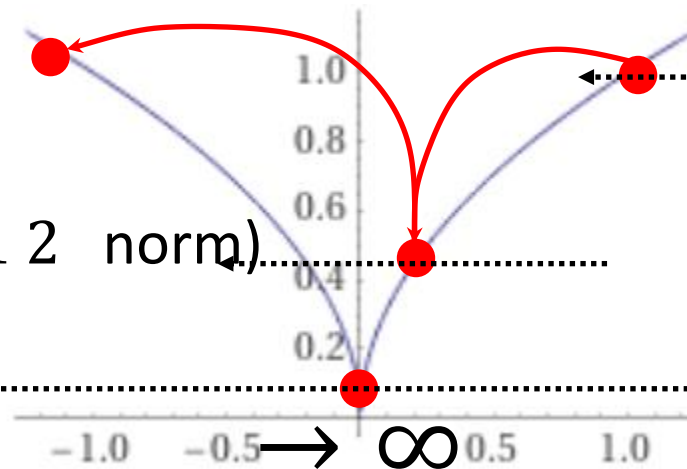


Q: Can we use any function inside a network?

- Has to be differentiable (L_2 norm)
- L_2 with respect to input and params
- Implications to optimization can be deadly:
 - 0-1 clamp layer in the end of the network.
 - Ideas how to solve this?



- Square root loss ($L_{1/2}$ norm)



Q: Can we use any function inside a network?

**Be creative,
but always watch your back(prop)!**



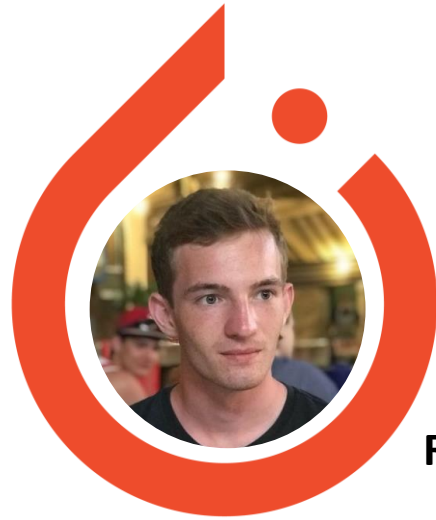
<http://playground.tensorflow.org>





DL4CV@Weizmann

This week's tutorial:



Intro to PyTorch

Rafail Fridman

This week's tutorial:



Rafail Fridman

Intro to PyTorch

Next week's lecture:

(Me
Again ☹)

Convolutional Neural Networks

