# Differential Programming

**Gabriel Peyré**

# Mathematical Coffees

Huawei-FSMP joint seminars

**https://mathematical-coffees.github.io**

**Organized by**: Mérouane Debbah & Gabriel Peyré

Optimization

Optimal Transport

Compressed Sensing

Quantum computing

Mean field games

Artificial intelligence

Deep Learning

Topos

Yves Achdou, Paris 6
Daniel Bennequin, Paris 7
Marco Cuturi, ENSAE
Jalal Fadili, ENSICaen

Alexandre Gramfort, INRIA
Olivier Grisel (INRIA)
Olivier Guéant, Paris 1
Iordanis Kerenidis, CNRS and Paris 7
Guillaume Lecué, CNRS and ENSAE

Frédéric Magniez, CNRS and Paris 7
Edouard Oyallon, CentraleSupelec
Gabriel Peyré, CNRS and ENS
Joris Van den Bossche (INRIA)

# Model Fitting in Data Sciences

$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} L(f(x, \theta), y)$$

Loss    Model    Input    Parameter    Output

# Model Fitting in Data Sciences

$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} L(f(x, \theta), y)$$

Loss    Model    Input    Parameter    Output

Deep-learning: $f(\cdot, \theta)$



$y$

class probabilities

$x$    $\theta_1$    $\theta_2$    $\theta_3$    $\theta_4$

# Model Fitting in Data Sciences

$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} L(f(x, \theta), y)$$

Loss  Model  Input  Parameter  Output

Deep-learning:  $f(\cdot, \theta)$



$y$

class probabilities

$x$  $\theta_1$  $\theta_2$  $\theta_3$  $\theta_4$

Super-resolution:



$f(x, \cdot)$

degradation

$\theta$ unknown image  $y$ observation

# Model Fitting in Data Sciences

$$\min_{\theta} \mathcal{E}(\theta) \overset{\text{def.}}{=} L(f(x, \theta), y)$$

Loss    Model    Input    Parameter    Output

Deep-learning:    $f(\cdot, \theta)$



$y$

class probabilities

$x$   $\theta_1$    $\theta_2$    $\theta_3$    $\theta_4$

Medical imaging registration:



$x$     $\xrightarrow[\text{diffeomorphism}]{f(\cdot, \theta)}$     $y$

Super-resolution:



$\xrightarrow[\text{degradation}]{f(x, \cdot)}$

$\theta$ unknown image      $y$ observation

# Gradient-based Methods

$$\min_{\theta} \mathcal{E}(\theta) \overset{\text{def.}}{=} L(f(x, \theta), y)$$

Gradient descent: $\quad \theta_{\ell+1} = \theta_\ell - \tau_\ell \nabla \mathcal{E}(\theta_\ell)$



Small $\tau_\ell$ $\qquad\qquad$ Large $\tau_\ell$ $\qquad\qquad$ Optimal $\tau_\ell = \tau_\ell^\star$

# Gradient-based Methods

$$\min_{\theta} \mathcal{E}(\theta) \stackrel{\text{def.}}{=} L(f(x, \theta), y)$$

Gradient descent: $\quad \theta_{\ell+1} = \theta_\ell - \tau_\ell \nabla \mathcal{E}(\theta_\ell)$



Small $\tau_\ell$         Large $\tau_\ell$         Optimal $\tau_\ell = \tau_\ell^\star$

**Many** generalization:

→ Nesterov / heavy-ball

→ (quasi)-Newton

→ Stochastic / incremental methods

→ Proximal splitting (non-smooth $\mathcal{E}$)

→ $\ldots$

# The Complexity of Gradient Computation

**Setup:** $\mathcal{E} : \mathbb{R}^n \rightarrow \mathbb{R}$ computable in $K$ operations.

```python
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \dots)$

and their derivatives cost $O(1)$.

# The Complexity of Gradient Computation

**Setup:** $\mathcal{E} : \mathbb{R}^n \to \mathbb{R}$ computable in $K$ operations.

```python
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \ldots)$

and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^n \to \mathbb{R}^n$?

# The Complexity of Gradient Computation

**Setup:** $\mathcal{E} : \mathbb{R}^n \to \mathbb{R}$ computable in $K$ operations.

```python
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \ldots)$

and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^n \to \mathbb{R}^n$?

Finite differences:
$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon}(\mathcal{E}(\theta + \varepsilon\delta_1) - \mathcal{E}(\theta), \ldots \mathcal{E}(\theta + \varepsilon\delta_n) - \mathcal{E}(\theta))$$
$K(n+1)$ operations, intractable for large $n$.

# The Complexity of Gradient Computation

**Setup:** $\mathcal{E} : \mathbb{R}^n \to \mathbb{R}$ computable in $K$ operations.

```python
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \dots)$
and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^n \to \mathbb{R}^n$?

Finite differences:
$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon}(\mathcal{E}(\theta + \varepsilon\delta_1) - \mathcal{E}(\theta), \dots \mathcal{E}(\theta + \varepsilon\delta_n) - \mathcal{E}(\theta))$$
$K(n+1)$ operations, intractable for large $n$.

*Theorem:* there is an algorithm to compute $\nabla \mathcal{E}$ in $O(K)$ operations.
[Seppo Linnainmaa, 1970]

# The Complexity of Gradient Computation

**Setup:** $\mathcal{E} : \mathbb{R}^n \to \mathbb{R}$ computable in $K$ operations.

```python
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

*Hypothesis:* elementary operations $(a \times b, \log(a), \sqrt{a} \ldots)$
and their derivatives cost $O(1)$.

**Question:** What is the complexity of computing $\nabla \mathcal{E} : \mathbb{R}^n \to \mathbb{R}^n$?

Finite differences: 
$$\nabla \mathcal{E}(\theta) \approx \frac{1}{\varepsilon}(\mathcal{E}(\theta + \varepsilon \delta_1) - \mathcal{E}(\theta), \ldots \mathcal{E}(\theta + \varepsilon \delta_n) - \mathcal{E}(\theta))$$
$K(n+1)$ operations, intractable for large $n$.

*Theorem:* there is an algorithm to compute $\nabla \mathcal{E}$ in $O(K)$ operations.
[Seppo Linnainmaa, 1970]

This algorithm is reverse mode automatic differentiation

```python
def BackwardNN(A,b,X):
    gx = lossG(X[R],Y) # initialize the gradient
    for r in arange(R-1,-1,-1):
        M = rhoG( A[r].dot(X[r]) + tile(b[r],[1,n]) ) * gx
        gx = A[r].transpose().dot(M)
        gA[r] =  M.dot(X[r].transpose())
        gb[r] =  MakeCol(M.sum(axis=1))
    return [gA,gb]
```

Seppo Linnainmaa

# Differentiating Composition of Functions

$x = x_0 \xrightarrow{g_0} x_1 \xrightarrow{g_1} x_2 \to \cdots \to x_R \xrightarrow{g_R} x_{R+1} \in \mathbb{R}$

$\partial g_r(x_r) \in \mathbb{R}^{n_{r+1} \times n_r}$

$x_{r+1} = g_r(x_r) \qquad g_r : \mathbb{R}^{n_r} \to \mathbb{R}^{n_{r+1}}$

$\nabla g_R(x_r) = [\partial g_r(x_r)]^\top \in \mathbb{R}^{n_{r+1} \times 1}$

# Differentiating Composition of Functions



$$x = x_0 \xrightarrow{g_0} x_1 \xrightarrow{g_1} x_2 \longrightarrow \cdots \longrightarrow x_R \xrightarrow{g_R} x_{R+1} \in \mathbb{R}$$

$$x_{r+1} = g_r(x_r) \qquad g_r : \mathbb{R}^{n_r} \to \mathbb{R}^{n_{r+1}}$$

$$\partial g_r(x_r) \in \mathbb{R}^{n_{r+1} \times n_r}$$

$$\nabla g_R(x_r) = [\partial g_r(x_r)]^\top \in \mathbb{R}^{n_{r+1} \times 1}$$

Chain rule:

$$\partial g(x) = \partial g_R(x_R) \times \partial g_{R-1}(x_{R-1}) \times \ldots \times \partial g_1(x_1) \times \partial g_0(x_0)$$

$$\underset{n_R}{\overset{A_R}{1\updownarrow \boxed{\phantom{A_R}}}} \times \underset{n_{R-1}}{\boxed{A_{R-1}}} \times \cdots \times \underset{n_2}{\boxed{A_1}} \times \underset{n_1}{\boxed{A_0}}\updownarrow n_0$$

# Differentiating Composition of Functions



$$x = x_0 \xrightarrow{g_0} x_1 \xrightarrow{g_1} x_2 \longrightarrow \cdots \longrightarrow x_R \xrightarrow{g_R} x_{R+1} \in \mathbb{R}$$
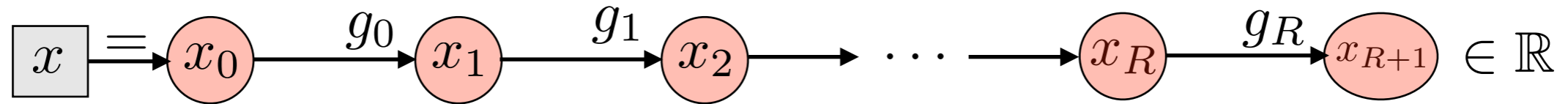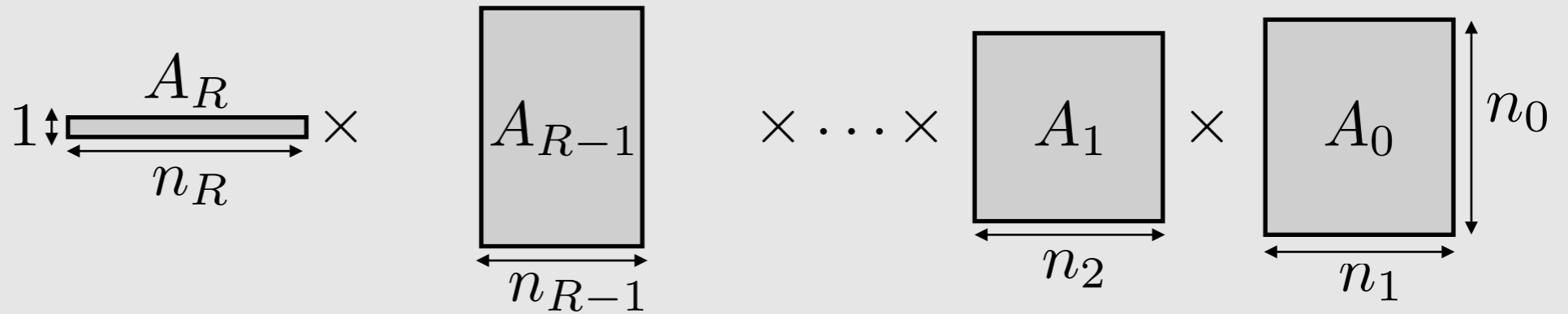
$$x_{r+1} = g_r(x_r) \qquad g_r : \mathbb{R}^{n_r} \to \mathbb{R}^{n_{r+1}}$$
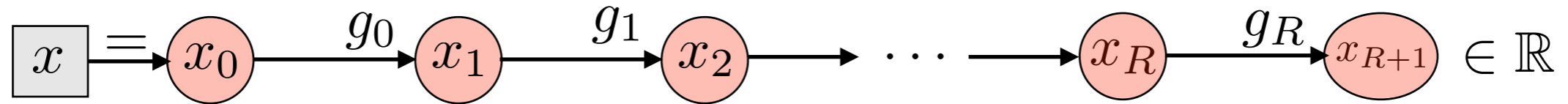
$$\partial g_r(x_r) \in \mathbb{R}^{n_{r+1} \times n_r}$$

$$\nabla g_R(x_r) = [\partial g_r(x_r)]^\top \in \mathbb{R}^{n_{r+1} \times 1}$$
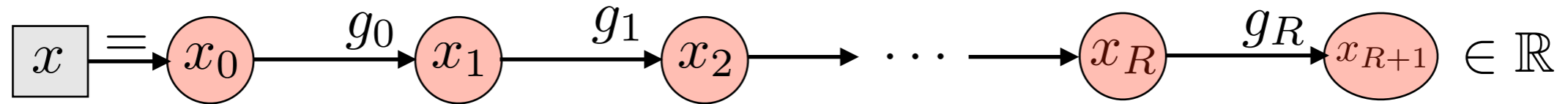
**Chain rule:**

$$\partial g(x) = \partial g_R(x_R) \times \partial g_{R-1}(x_{R-1}) \times \ldots \times \partial g_1(x_1) \times \partial g_0(x_0)$$



**Forward** $O(n^3)$

$$\partial g(x) = ((\ldots((\underbrace{A_0 \times A_1)}_{n_0 n_1 n_2} \times A_2)\ldots \times A_{R-2})}_{n_1 n_2 n_3} \times \underbrace{A_{R-1}) \times A_R}_{n_{R-2} n_{R-1} n_R}$$
$$\underbrace{\phantom{A_{R-1}) \times A_R}}_{n_{R-1} n_R}$$

*Complexity:* (if $n_r = 1$ for $r = 0, \ldots, R-1$) $(R-1)n^3 + n^2$

# Differentiating Composition of Functions



$$x = x_0 \xrightarrow{g_0} x_1 \xrightarrow{g_1} x_2 \rightarrow \cdots \rightarrow x_R \xrightarrow{g_R} x_{R+1} \in \mathbb{R}$$

$$x_{r+1} = g_r(x_r) \qquad g_r : \mathbb{R}^{n_r} \to \mathbb{R}^{n_{r+1}}$$

$$\partial g_r(x_r) \in \mathbb{R}^{n_{r+1} \times n_r}$$

$$\nabla g_R(x_r) = [\partial g_r(x_r)]^\top \in \mathbb{R}^{n_{r+1} \times 1}$$

**Chain rule:**

$$\partial g(x) = \partial g_R(x_R) \times \partial g_{R-1}(x_{R-1}) \times \ldots \times \partial g_1(x_1) \times \partial g_0(x_0)$$



**Forward** $O(n^3)$

$$\partial g(x) = ((\ldots((\underbrace{A_0 \times A_1}_{n_0 n_1 n_2}) \times A_2)\underbrace{}_{n_1 n_2 n_3} \ldots \times \underbrace{A_{R-2}) \times A_{R-1}}_{n_{R-2} n_{R-1} n_R}) \times A_R \underbrace{}_{n_{R-1} n_R}$$

*Complexity:* (if $n_r = 1$ for $r = 0, \ldots, R-1$) $(R-1)n^3 + n^2$

**Backward** $O(n^2)$

$$\partial g(x) = A_0 \times (\underbrace{A_1 \times (A_2}_{\underbrace{n_1 n_2}_{n_0 n_1}} \times \ldots \times (\underbrace{A_{R-2} \times (A_{R-1} \times A_R}_{\underbrace{n_{R-1} n_R}_{n_{R-2} n_{R-1}}}))\ldots))$$

*Complexity:* $Rn^2$

# Feedfordward Computational Graphs

# Feedfordward Computational Graphs

$$x_{r+1} = g_r(x_r, \theta_r) \qquad \mathcal{E}(x) = L(x_{R+1}, y)$$



*Example:* deep neural network (here fully connected)



$$x_{r+1} = \rho(A_r x_r + b_r) \qquad \theta_r = (A_r, b_r)$$

$$x_r \in \mathbb{R}^{d_r}$$

$$A_r \in \mathbb{R}^{d_{r+1} \times d_r}$$

$$b_r \in \mathbb{R}^{d_{r+1}}$$

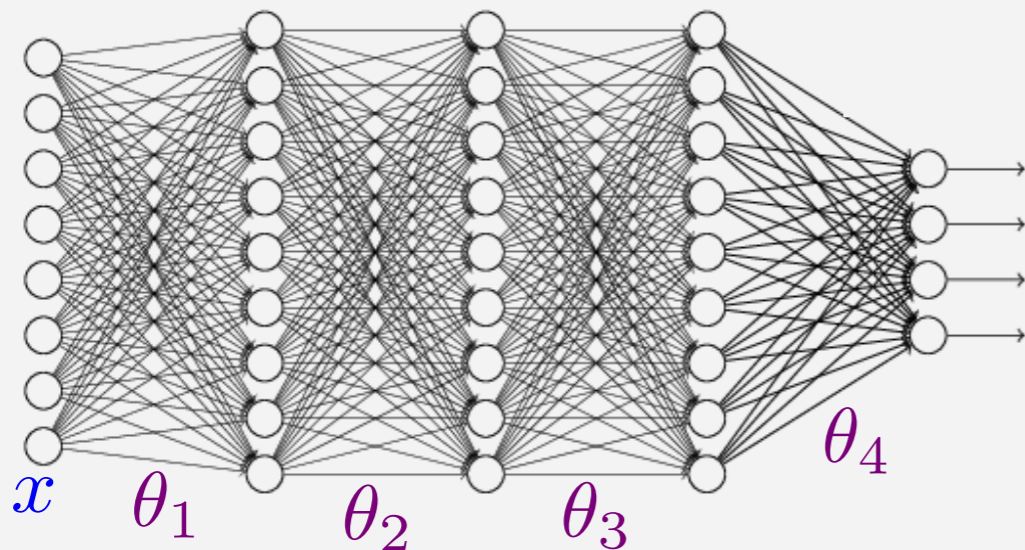# Feedfordward Computational Graphs

$$x_{r+1} = g_r(x_r, \theta_r) \qquad \mathcal{E}(x) = L(x_{R+1}, y)$$



*Example:* deep neural network (here fully connected)
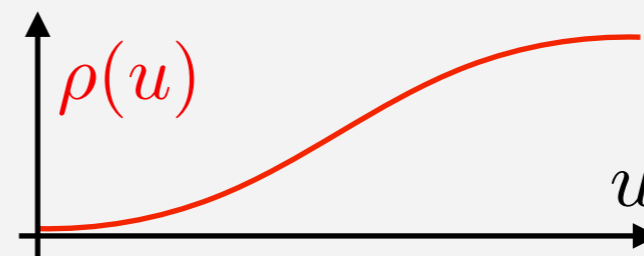


$$x_{r+1} = \rho(A_r x_r + b_r) \qquad \theta_r = (A_r, b_r)$$

$$x_r \in \mathbb{R}^{d_r}$$
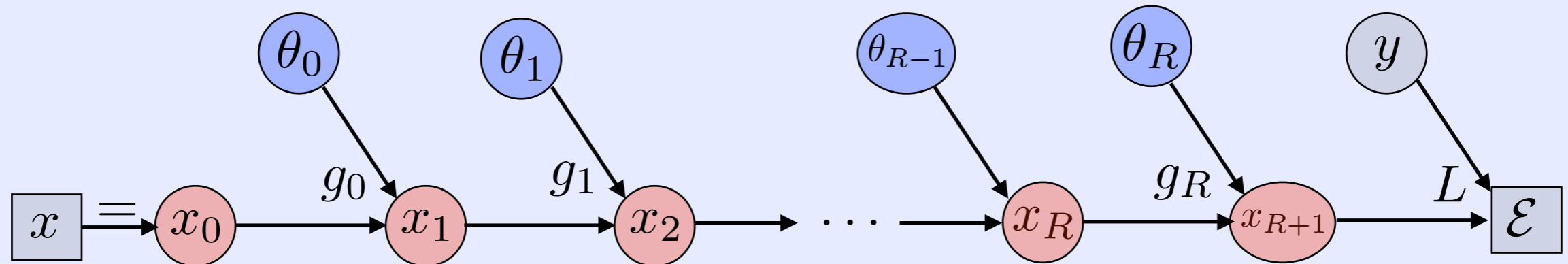
$$A_r \in \mathbb{R}^{d_{r+1} \times d_r}$$
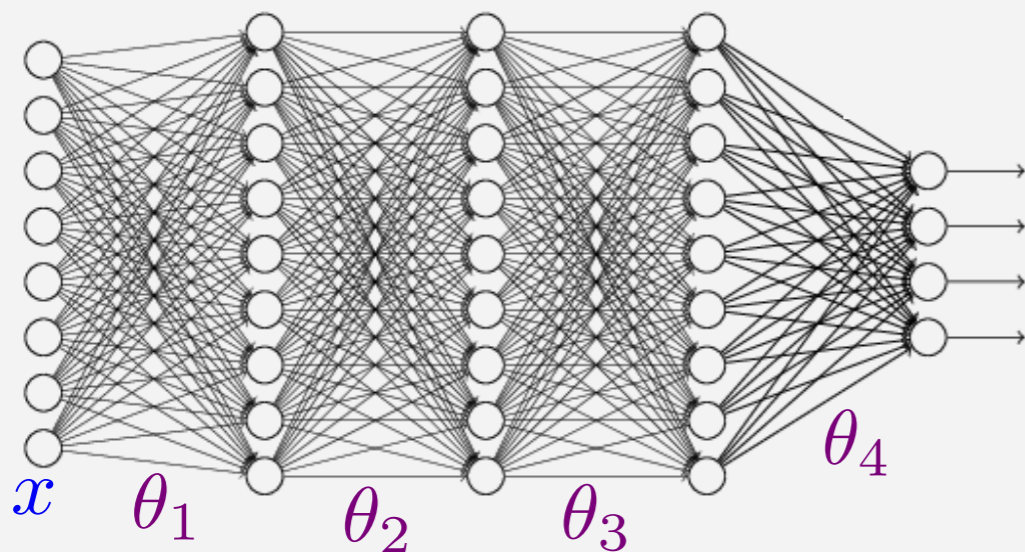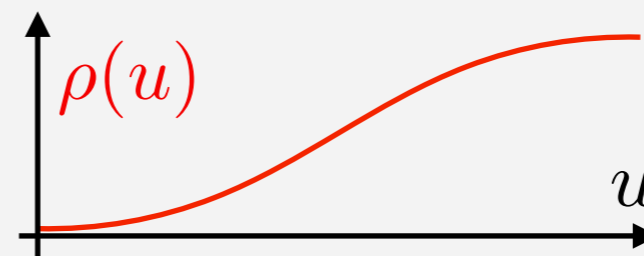
$$b_r \in \mathbb{R}^{d_{r+1}}$$

Logistic loss:
(classification)

$$L(x_{R+1}, y) \overset{\text{def.}}{=} \log \sum_i \exp(x_{R+1,i}) - x_{R+1,i} y_i$$

$$\nabla_{x_{R+1}} L(x_{R+1}, y) = \frac{e^{x_{R+1}}}{\sum_i e^{x_{R+1,i}}} - y$$

# Backpropagation Algorithm

# Backpropagation Algorithm

$$x_{r+1} = g_r(x_r, \theta_r) \qquad \mathcal{E}(x) = L(x_{R+1}, y)$$



*Proposition:* $\quad \forall r = R, \ldots, 0, \qquad \nabla_{x_r}\mathcal{E} = [\partial_{x_r} g_R(x_r, \theta_r)]^\top (\nabla_{x_{r+1}}\mathcal{E})$

$$\nabla_{\theta_r}\mathcal{E} = [\partial_{\theta_r} g_R(x_r, \theta_r)]^\top (\nabla_{x_{r+1}}\mathcal{E})$$

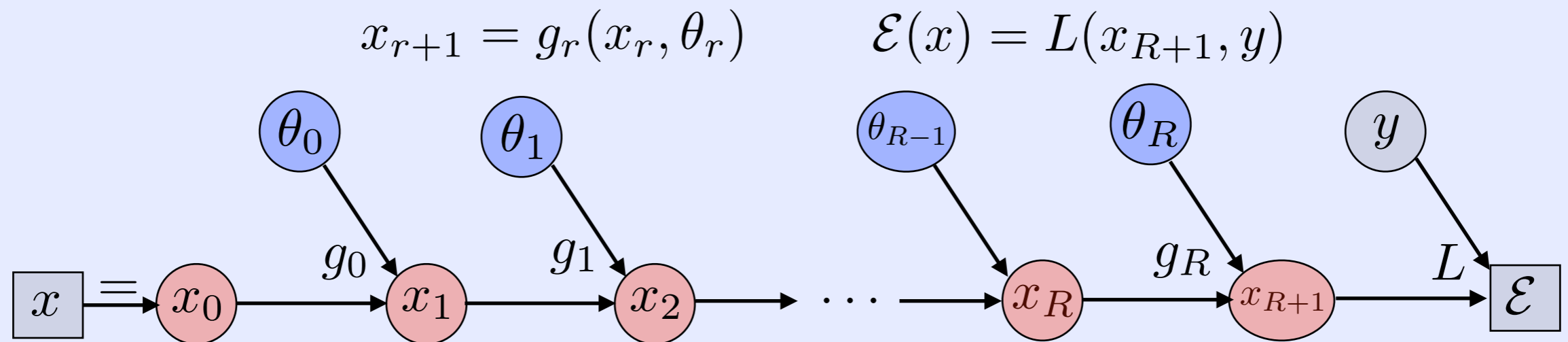# Backpropagation Algorithm

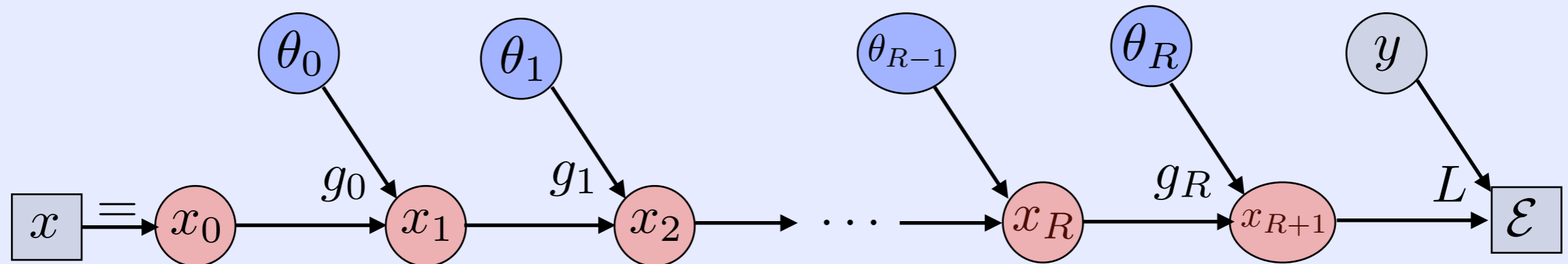$$x_{r+1} = g_r(x_r, \theta_r) \qquad \mathcal{E}(x) = L(x_{R+1}, y)$$



*Proposition:* $\forall r = R, \dots, 0,$

$$\nabla_{x_r} \mathcal{E} = [\partial_{x_r} g_R(x_r, \theta_r)]^\top (\nabla_{x_{r+1}} \mathcal{E})$$

$$\nabla_{\theta_r} \mathcal{E} = [\partial_{\theta_r} g_R(x_r, \theta_r)]^\top (\nabla_{x_{r+1}} \mathcal{E})$$

*Example:* deep neural network $\quad x_{r+1} = \rho(A_r x_r + b_r)$

$$\nabla_{x_r} \mathcal{E} = A_r^\top M_r$$

$\forall r = R, \dots, 0, \qquad \nabla_{A_r} \mathcal{E} = M_r x_r^\top \qquad M_r \overset{\text{def.}}{=} \rho'(A_r x_r + b_r) \odot \nabla_{x_{r+1}} \mathcal{E}$

$$\nabla_{b_r} \mathcal{E} = M_r \mathbb{1}$$

```python
def ForwardNN(A,b,Z):
    X = []
    X.append(Z)
    for r in arange(0,R):
        X.append( rhoF( A[r].dot(X[r]) + tile(b[r],[1,Z.shape[1]]) ) )
    return X
```

```python
def BackwardNN(A,b,X):
    gx = lossG(X[R],Y) # initialize the gradient
    for r in arange(R-1,-1,-1):
        M = rhoG( A[r].dot(X[r]) + tile(b[r],[1,n]) ) * gx
        gx = A[r].transpose().dot(M)
        gA[r] =  M.dot(X[r].transpose())
        gb[r] =  MakeCol(M.sum(axis=1))
    return [gA,gb]
```

# Recurrent Architectures

Shared parameters:  $x_{r+1} = g_r(x_r, \theta)$

$\theta$

$y$

$x$ $=$ $x_0$ $\xrightarrow{g_0}$ $x_1$ $\xrightarrow{g_1}$ $x_2$ $\longrightarrow$ $\cdots$ $\longrightarrow$ $x_R$ $\xrightarrow{g_R}$ $x_{R+1}$ $\xrightarrow{L}$ $\mathcal{E}$

# Recurrent Architectures
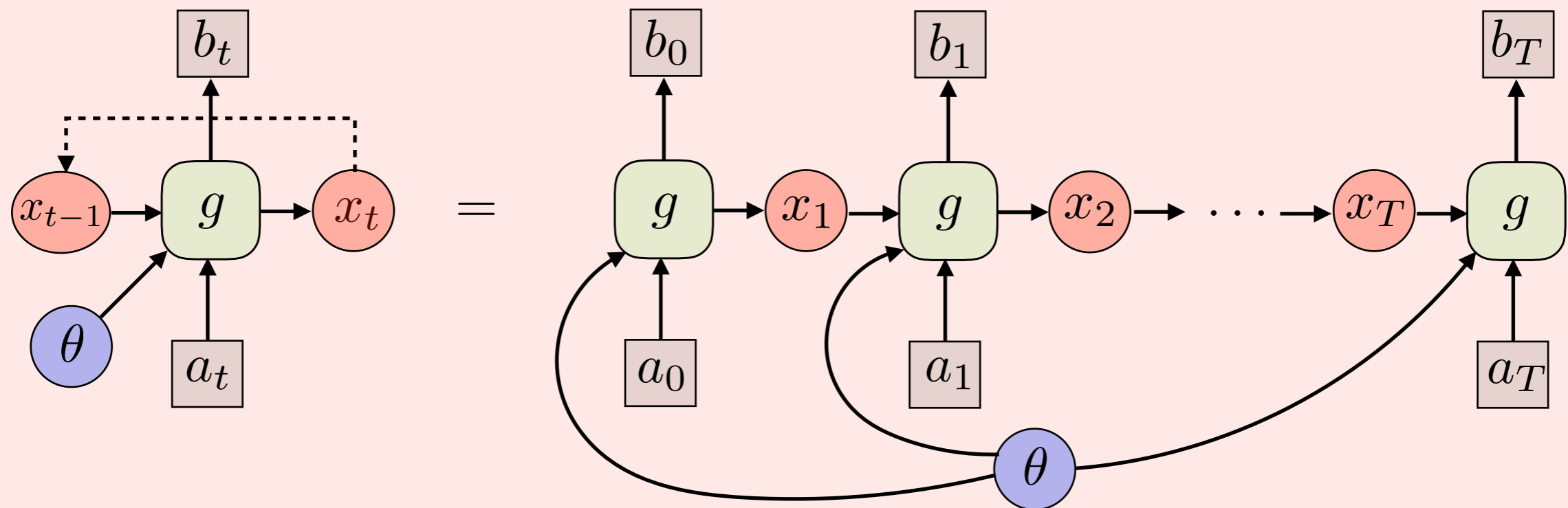
Shared parameters: $x_{r+1} = g_r(x_r, \theta)$



Recurrent networks for natural language processing:

# Recurrent Architectures

Shared parameters: $x_{r+1} = g_r(x_r, \theta)$



Recurrent networks for natural language processing:



*Take home message:* for complicated computational architectures, you do not want to do the computation/implementation by hand.

# Computational Graph

# Computational Graph

Computer program $\Leftrightarrow$ directed acyclic graph $\Leftrightarrow$ linear ordering of nodes $(\theta_r)_r$

# Example

$$\ell(\theta_1, \theta_2) \overset{\text{def.}}{=} \theta_2 e^{\theta_1} \sqrt{\theta_1 + \theta_2 e^{\theta_1}}$$

# Example

$$\ell(\theta_1, \theta_2) \stackrel{\text{def.}}{=} \theta_2 e^{\theta_1} \sqrt{\theta_1 + \theta_2 e^{\theta_1}}$$
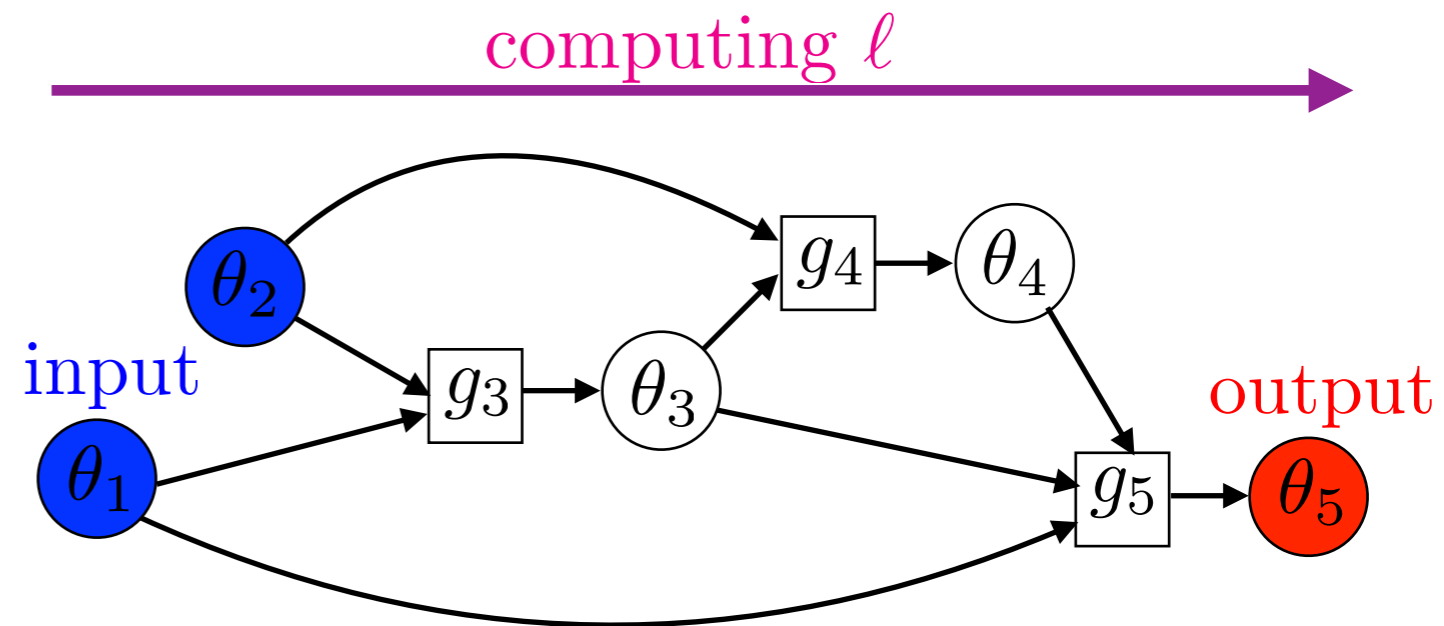


input

output $\ell$

Chain rules:

$$\text{``} \frac{\partial \theta_j}{\partial \theta_1} = \sum_{i \in \text{Parent}(j)} \frac{\partial \theta_j}{\partial \theta_i} \frac{\partial \theta_i}{\partial \theta_1} \text{''}$$

$\partial_i g_j(\theta)$

"Classical" evaluation: **forward**.
Complexity $\sim$ #inputs.

# Example

$$\ell(\theta_1, \theta_2) \stackrel{\text{def.}}{=} \theta_2 e^{\theta_1} \sqrt{\theta_1 + \theta_2 e^{\theta_1}}$$



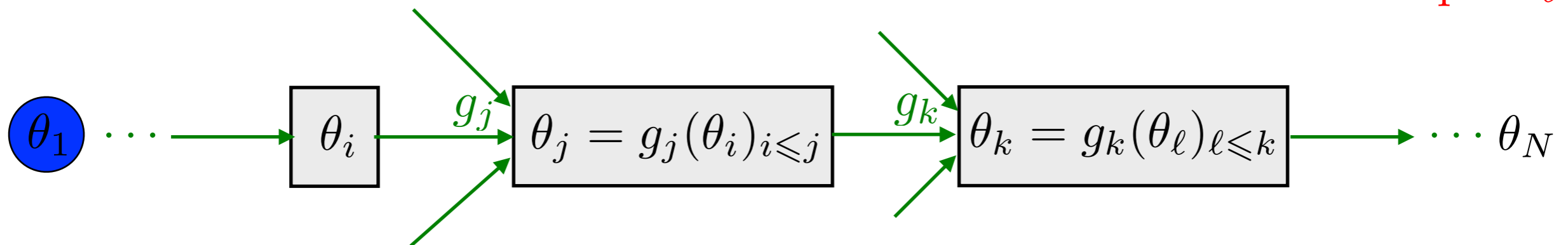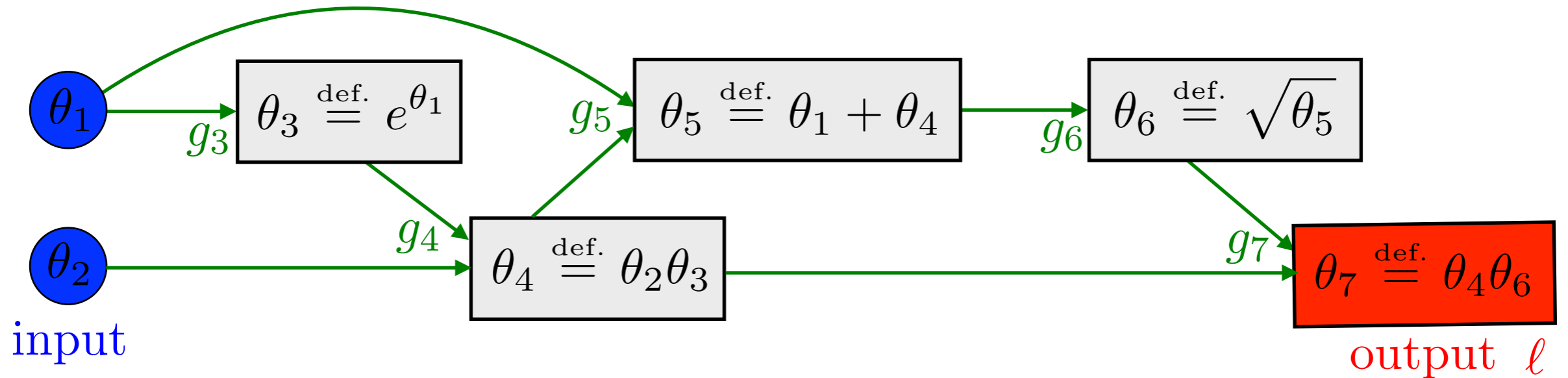$\theta_1$   $\theta_3 \stackrel{\text{def.}}{=} e^{\theta_1}$   $g_3$   $g_5$   $\theta_5 \stackrel{\text{def.}}{=} \theta_1 + \theta_4$   $g_6$   $\theta_6 \stackrel{\text{def.}}{=} \sqrt{\theta_5}$

$\theta_2$   $g_4$   $\theta_4 \stackrel{\text{def.}}{=} \theta_2 \theta_3$   $g_7$   $\theta_7 \stackrel{\text{def.}}{=} \theta_4 \theta_6$

input

output $\ell$

$\theta_1$ $\cdots$ $\theta_i$   $g_j$   $\theta_j = g_j(\theta_i)_{i \leqslant j}$   $g_k$   $\theta_k = g_k(\theta_\ell)_{\ell \leqslant k}$ $\cdots \theta_N$

Chain rules:
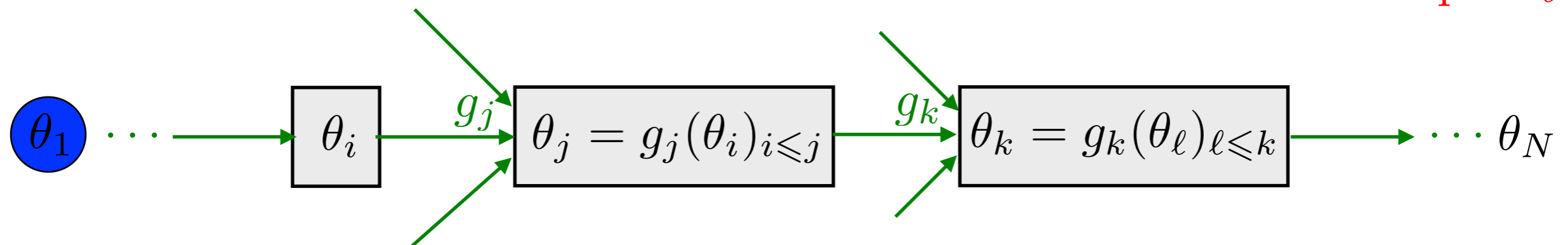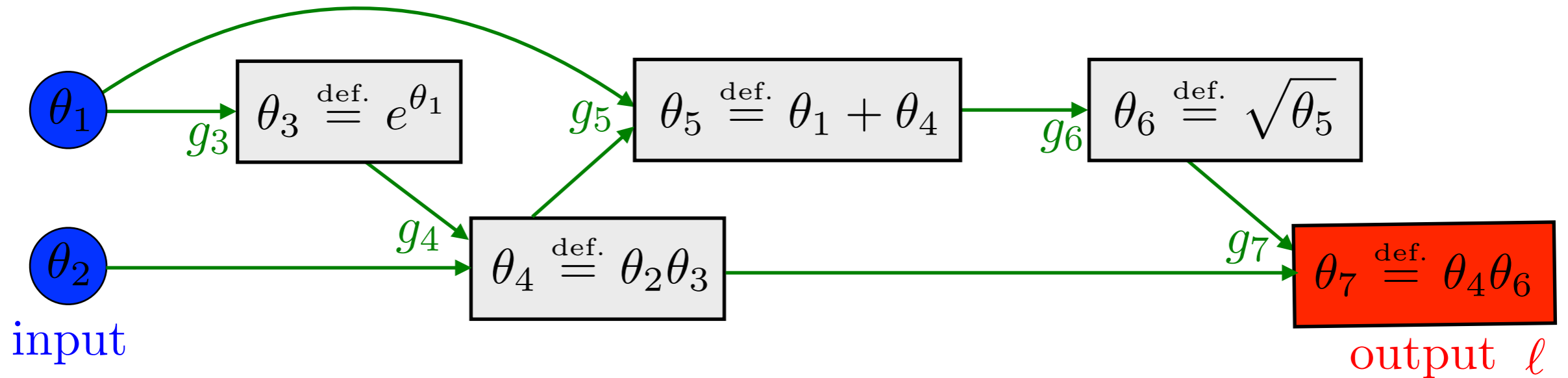
$$``\frac{\partial \theta_j}{\partial \theta_1} = \sum_{i \in \text{Parent}(j)} \frac{\partial \theta_j}{\partial \theta_i} \frac{\partial \theta_i}{\partial \theta_1}"$$

$\partial_i g_j(\theta)$

$$``\frac{\partial \theta_N}{\partial \theta_j} = \sum_{k \in \text{Child}(j)} \frac{\partial \theta_N}{\partial \theta_k} \frac{\partial \theta_k}{\partial \theta_j}"$$
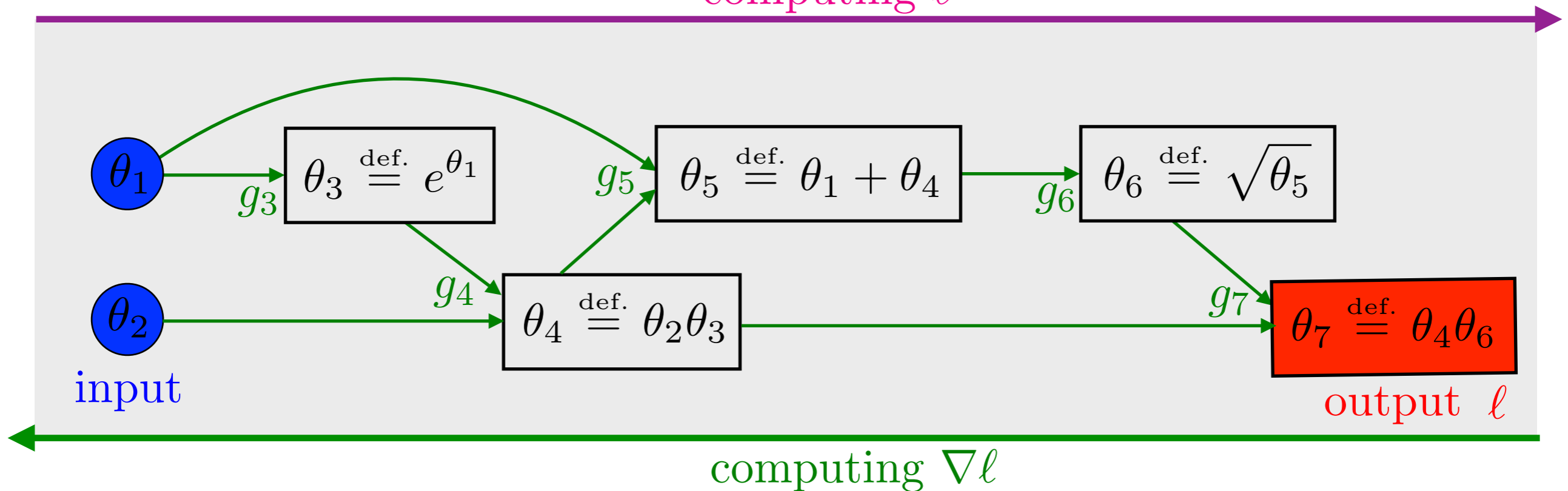
$\nabla_j \ell(\theta)$   $\nabla_k \ell(\theta)$   $\partial_j g_k(\theta)$

"Classical" evaluation: **forward**.
Complexity $\sim$ #inputs.

**Backward** evaluation.
Complexity $\sim$ #outputs (1 for grad).

# Backward Automatic Differentiation

$$\ell(\theta_1, \theta_2) \overset{\text{def.}}{=} \theta_2 e^{\theta_1} \sqrt{\theta_1 + \theta_2 e^{\theta_1}}$$

computing $\ell$



$\theta_1$

$\theta_3 \overset{\text{def.}}{=} e^{\theta_1}$

$g_3$

$g_5$

$\theta_5 \overset{\text{def.}}{=} \theta_1 + \theta_4$

$g_6$

$\theta_6 \overset{\text{def.}}{=} \sqrt{\theta_5}$

$\theta_2$

$g_4$

$\theta_4 \overset{\text{def.}}{=} \theta_2 \theta_3$

$g_7$

$\theta_7 \overset{\text{def.}}{=} \theta_4 \theta_6$

input

output $\ell$

computing $\nabla \ell$

forward

```
function ℓ(θ₁,...,θ_M)
  for r = M + 1,...,R
  |   θ_r = g_r(θ_Parents(r))
  return θ_R
```

$$\text{function } \ell(\theta_1, \ldots, \theta_M)$$
$$\text{for } r = M+1, \ldots, R$$
$$\qquad \theta_r = g_r(\theta_{\text{Parents}(r)})$$
$$\text{return } \theta_R$$

backward

$$\text{function } \nabla\ell(\theta_1, \ldots, \theta_M)$$
$$\nabla_R \ell = 1$$
$$\text{for } r = R-1, \ldots, 1$$
$$\qquad \nabla_r \ell = \sum_{s \in \text{Child}(r)} \partial_r g_s(\theta)\, \nabla_s \ell$$
$$\text{return } (\nabla_1 \ell, \ldots, \nabla_M \ell)$$

# Softwares