# Parametric Models Fitting with Automatic Differentiation
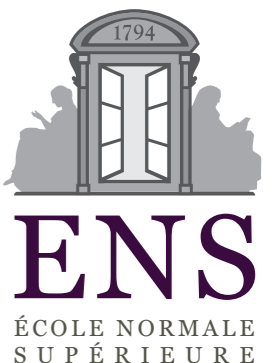
## Gabriel Peyré
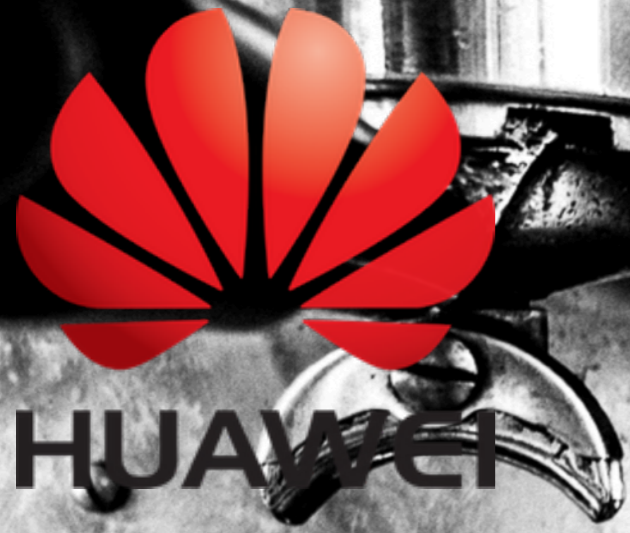
# Mathematical Coffees

Huawei-FSMP joint seminars

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

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

Optimal Transport

Geodesics

Meshes

Optimization

Deep Learning

Diagonal Line Node — Face Node — Cat Node

Sparsity

Neuro-imaging

Patches

Bayesian

Parallel/Stochastic

Alexandre Allauzen, Paris-Sud.
Pierre Alliez, INRIA.
Guillaume Charpiat, INRIA.
Emilie Chouzenoux, Paris-Est.

Nicolas Courty, IRISA.
Laurent Cohen, CNRS Dauphine.
Marco Cuturi, ENSAE.
Julie Delon, Paris 5.

Fabian Pedregosa, INRIA.
Julien Tierny, CNRS and P6.
Robin Ryder, Paris-Dauphine.
Gael Varoquaux, INRIA.

Jalal Fadili, ENSICaen.
Alexandre Gramfort, INRIA.
Matthieu Kowalski, Supelec.
Jean-Marie Mirebeau, CNRS,P-Sud.

# Parametric Models

(Noisy) observations $(x_i, y_j)$, try to infer $y = f(x)$.



Regression $\quad (x, y) \in \mathbb{R}^n \times \mathbb{R}^p$

Classification $\quad (x, y) \in \mathbb{R}^n \times \{-1, 1\}$

# Parametric Models

(Noisy) observations $(x_i, y_j)$, try to infer $y = f(x)$.



$y = f(x)$

$(x_i, y_i)$

Regression $\quad (x, y) \in \mathbb{R}^n \times \mathbb{R}^p$



$x_i$ $\quad x_j$

$f_i = -1$

$f_j = 1$

$f(x) = 0$

Classification $(x, y) \in \mathbb{R}^n \times \{-1, 1\}$

Parametric model: $y = f(x, \theta)$, find $\theta$.

Linear model: $f(x, \theta) = \langle x, \theta \rangle$.



$y = f(x)$



$f(x) = 0$

# Parametric Models

(Noisy) observations $(x_i, y_j)$, try to infer $y = f(x)$.



$y = f(x)$

$(x_i, y_i)$

Regression $\qquad (x, y) \in \mathbb{R}^n \times \mathbb{R}^p$



$x_i$ $\qquad$ $x_j$

$f_i = -1$

$f_j = 1$

$f(x) = 0$

Classification $(x, y) \in \mathbb{R}^n \times \{-1, 1\}$

Parametric model: $y = f(x, \theta)$, find $\theta$.

Linear model: $f(x, \theta) = \langle x, \theta \rangle$.



$y = f(x)$



$f(x) = 0$

Deep network:
$f(x, \theta) = \theta_K(\dots \rho(\theta_2(\rho(\theta_1(x) \dots))$



$x$ $\quad \theta_1 \quad \rho \quad \theta_2 \quad \rho \quad \theta_3 \quad \rho \quad \theta_4 \qquad f(x)$

# Empirical Loss Minimization

Regression: $\quad(y, y') \in \mathbb{R}^d \times \mathbb{R}^d, \qquad L(y, y') = \|y - y'\|^2$

Classification: $\quad(y, y') \in \mathbb{R}^d \times \{-1, 1\}, \quad L(y, y') = \log(\exp(-y'y) + 1)$

Loss minimization:

$$\min_{\theta} \sum_i L(f(x_i, \theta), y_i) \qquad\qquad \min_{\theta} \mathbb{E}_{(X,Y)}(L(f(X, \theta), Y))$$

# Empirical Loss Minimization

Regression: $\qquad (y, y') \in \mathbb{R}^d \times \mathbb{R}^d, \qquad L(y, y') = \|y - y'\|^2$

Classification: $\quad (y, y') \in \mathbb{R}^d \times \{-1, 1\}, \quad L(y, y') = \log(\exp(-y'y) + 1)$

Loss minimization:

$$\min_\theta \sum_i L(f(x_i, \theta), y_i) \qquad\qquad \min_\theta \mathbb{E}_{(X,Y)}(L(f(X, \theta), Y))$$

Stochastic gradient descent:

– Sample: $\qquad (x, y) \in \{(x_i, y_i)\}_i \qquad\qquad (x, y) \sim (X, Y)$

– Update: $\qquad \theta^{(\ell+1)} \overset{\text{def.}}{=} \theta^{(\ell)} - \tau_\ell \nabla_\theta \ell_{x,y}(\theta)$

$$\text{where} \quad \ell_{x,y}(\theta) \overset{\text{def.}}{=} L(f(x, \theta), y)$$

# Gradient Computation

How to compute $\nabla \ell_{x,y}(\theta)$? $\qquad$ $\ell_{x,y}(\theta) \overset{\text{def.}}{=} L(f(x, \theta), y)$

Chain rule: $\nabla \ell_{x,y}(\theta) = [\partial f(x, \theta)]^{\top} (\nabla L(f(x, \theta), y))$

Linear $f(x, \theta) = \theta \times x$: $\partial f(x, \theta) = \theta$.

Non-linear $f(x, \theta)$: painful ... but $\ell_{x,y}$ it is just a computer program.

# Gradient Computation

How to compute $\nabla \ell_{x,y}(\theta)$?     $\ell_{x,y}(\theta) \overset{\text{def.}}{=} L(f(x,\theta), y)$

Chain rule:  $\nabla \ell_{x,y}(\theta) = [\partial f(x,\theta)]^\top (\nabla L(f(x,\theta), y))$

Linear $f(x,\theta) = \theta \times x$: $\partial f(x,\theta) = \theta$.

Non-linear $f(x,\theta)$: painful ...   but $\ell_{x,y}$ it is just a computer program.

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



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

# Example

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

# Example

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



$\theta_1$

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

$g_3$

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

$g_5$

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

$g_6$

$\theta_2$

$g_4$
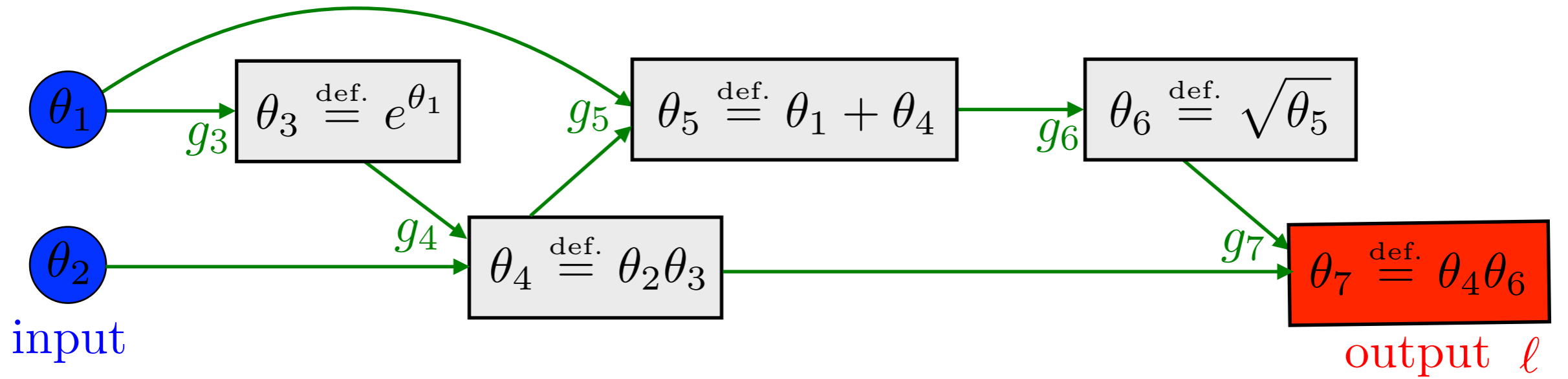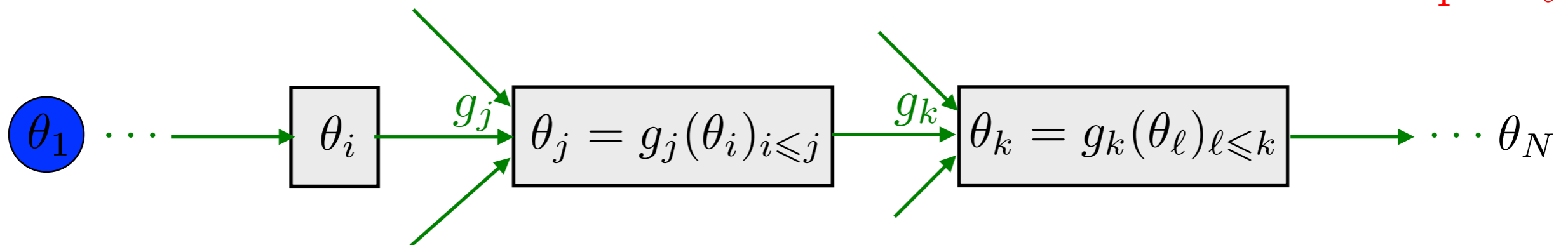
$\theta_4 \overset{\text{def.}}{=} \theta_2 \theta_3$

$g_7$

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

input

output $\ell$

$\theta_1$

$\theta_i$

$g_j$

$\theta_j = g_j(\theta_i)_{i \leqslant j}$

$g_k$

$\theta_k = g_k(\theta_\ell)_{\ell \leqslant k}$

$\theta_N$

Chain rules:

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

$\partial_i g_j(\theta)$
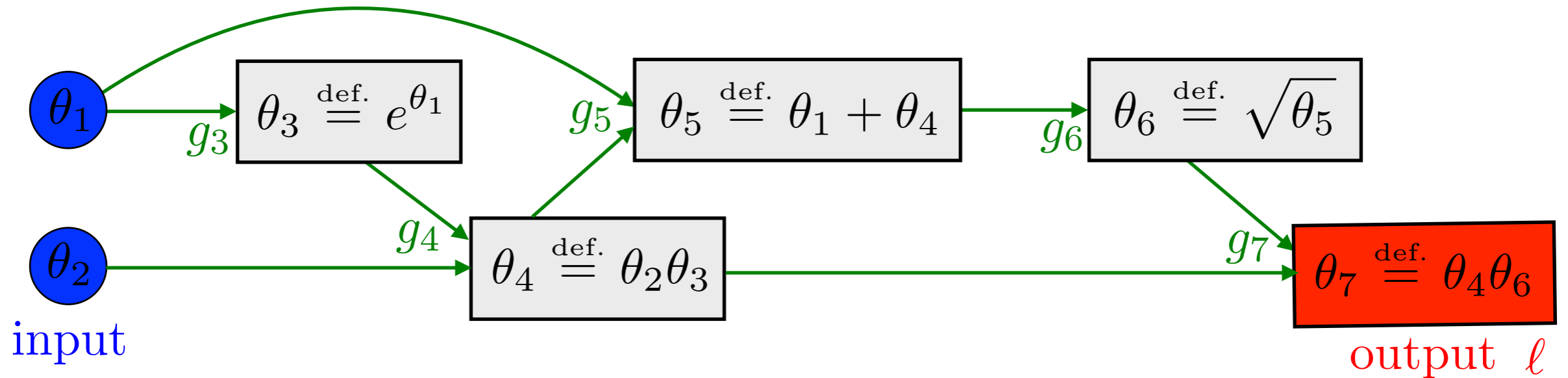
"Classical" evaluation: **forward**.
Complexity $\sim$ #inputs.

# Example

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



$\theta_1$ (input)

$g_3$   $\theta_3 \overset{\text{def.}}{=} e^{\theta_1}$

$g_5$   $\theta_5 \overset{\text{def.}}{=} \theta_1 + \theta_4$

$g_6$   $\theta_6 \overset{\text{def.}}{=} \sqrt{\theta_5}$

$g_4$

$\theta_2$ (input)

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

$g_7$   $\theta_7 \overset{\text{def.}}{=} \theta_4 \theta_6$   output $\ell$

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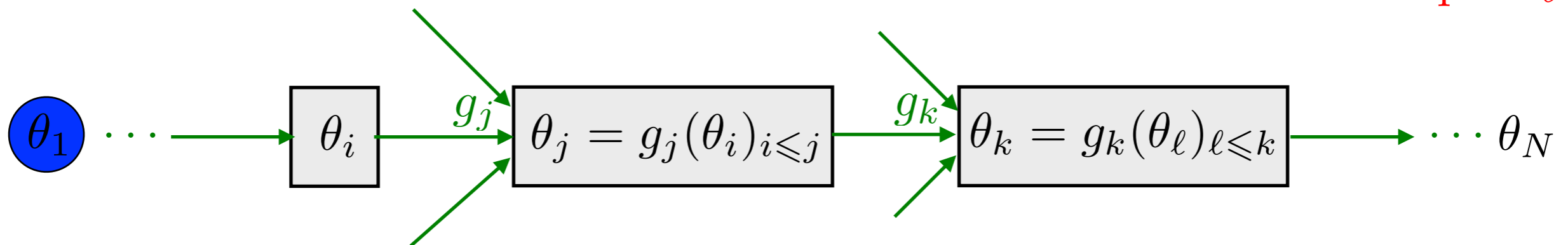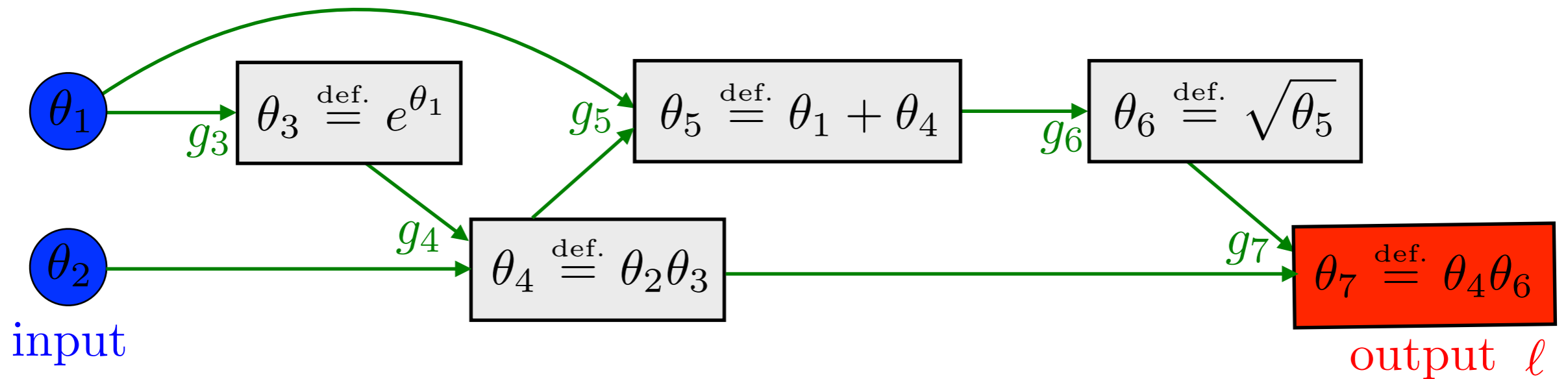
Chain rules:

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

$$\partial_i g_j(\theta)$$

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

$$``\frac{\partial \theta_N}{\partial \theta_j} = \sum_{k \in \text{Child}(j)} \frac{\partial \theta_N}{\partial \theta_k} \frac{\partial \theta_k}{\partial \theta_j}"$$

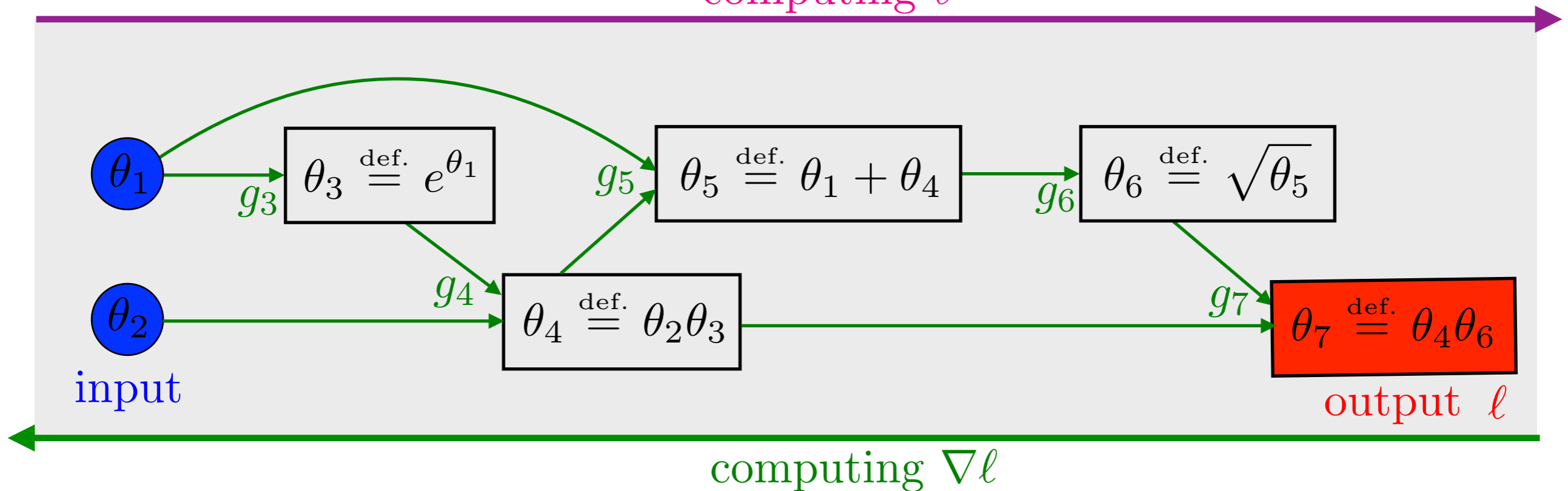$$\nabla_j \ell(\theta) \qquad \nabla_k \ell(\theta) \qquad \partial_j g_k(\theta)$$

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

# Backward Automatic Differentiation

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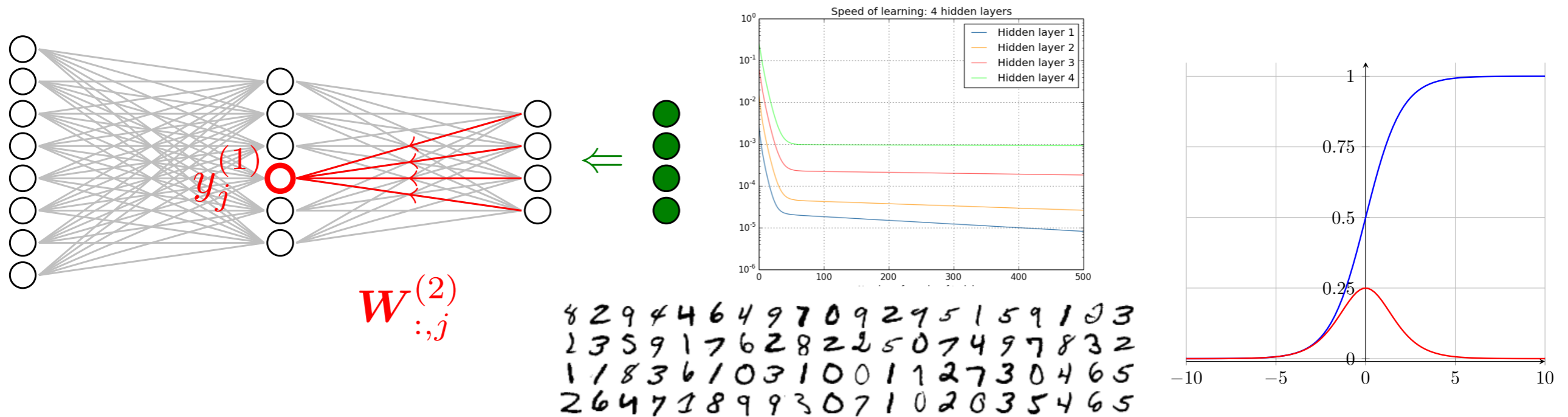
computing $\ell$



computing $\nabla\ell$

forward

```
function ℓ(θ₁,…,θ_M)
  for r = M + 1,…,R
  │   θ_r = g_r(θ_{Parents(r)})
  return θ_R
```

backward

```
function ∇ℓ(θ₁,…,θ_M)
  ∇_R ℓ = 1
  for r = R − 1,…,1
  │     ∇_r ℓ = ∑_{s∈Child(r)} ∂_r g_s(θ) ∇_s ℓ
  return (∇₁ℓ,…,∇_M ℓ)
```

# What's Next

**Alexandre Allauzen:** deep neural networks training.



**Guillaume Charpiat:** architecture of deep neural networks.