

95-865 Unstructured Data Analytics

Some slides for part of recitation:
How to train a deep net

Slides by George H. Chen

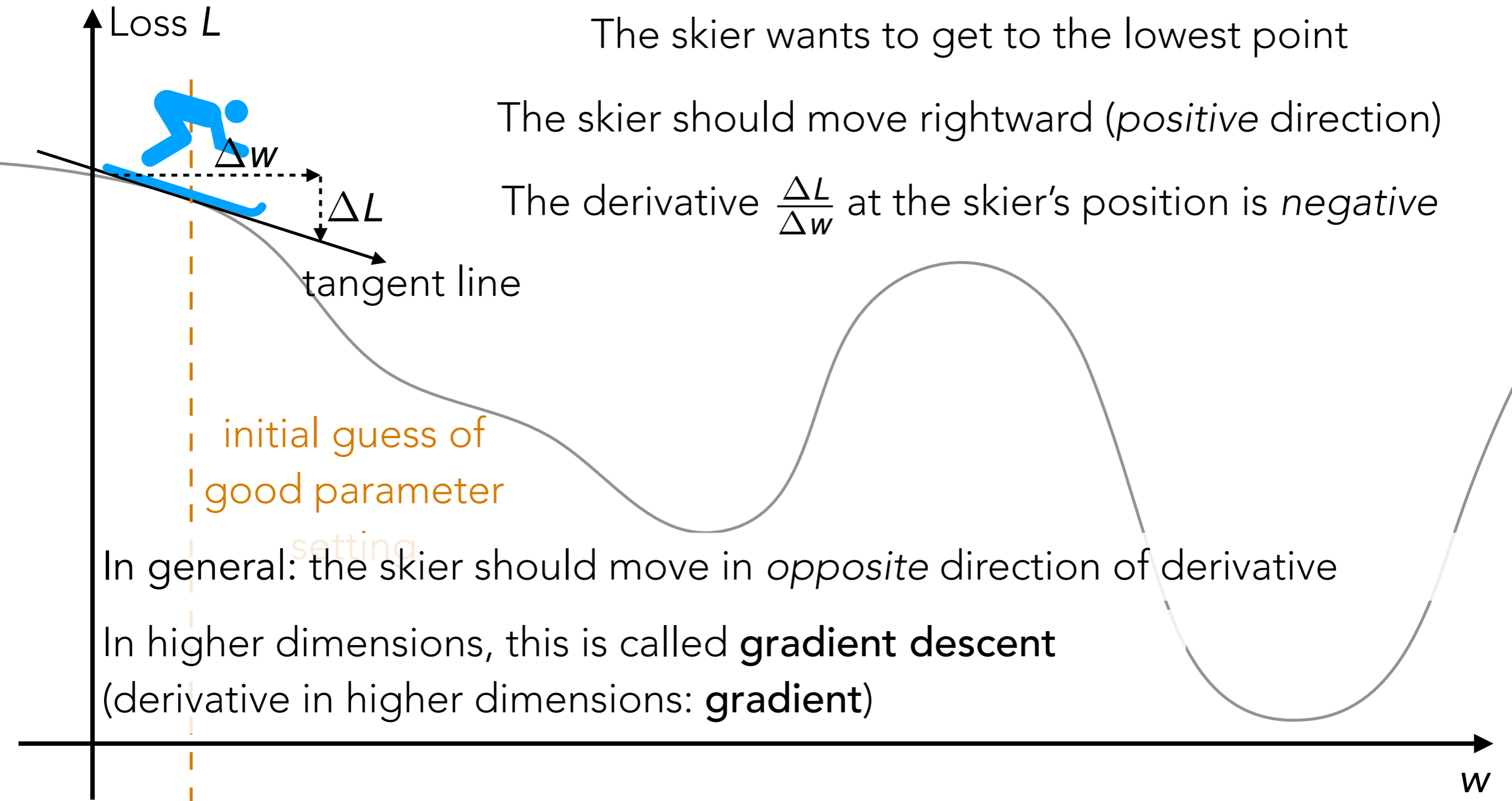
Training a Deep Net

Suppose the neural network has a single real number parameter w

The skier wants to get to the lowest point

The skier should move rightward (*positive* direction)

The derivative $\frac{\Delta L}{\Delta w}$ at the skier's position is *negative*



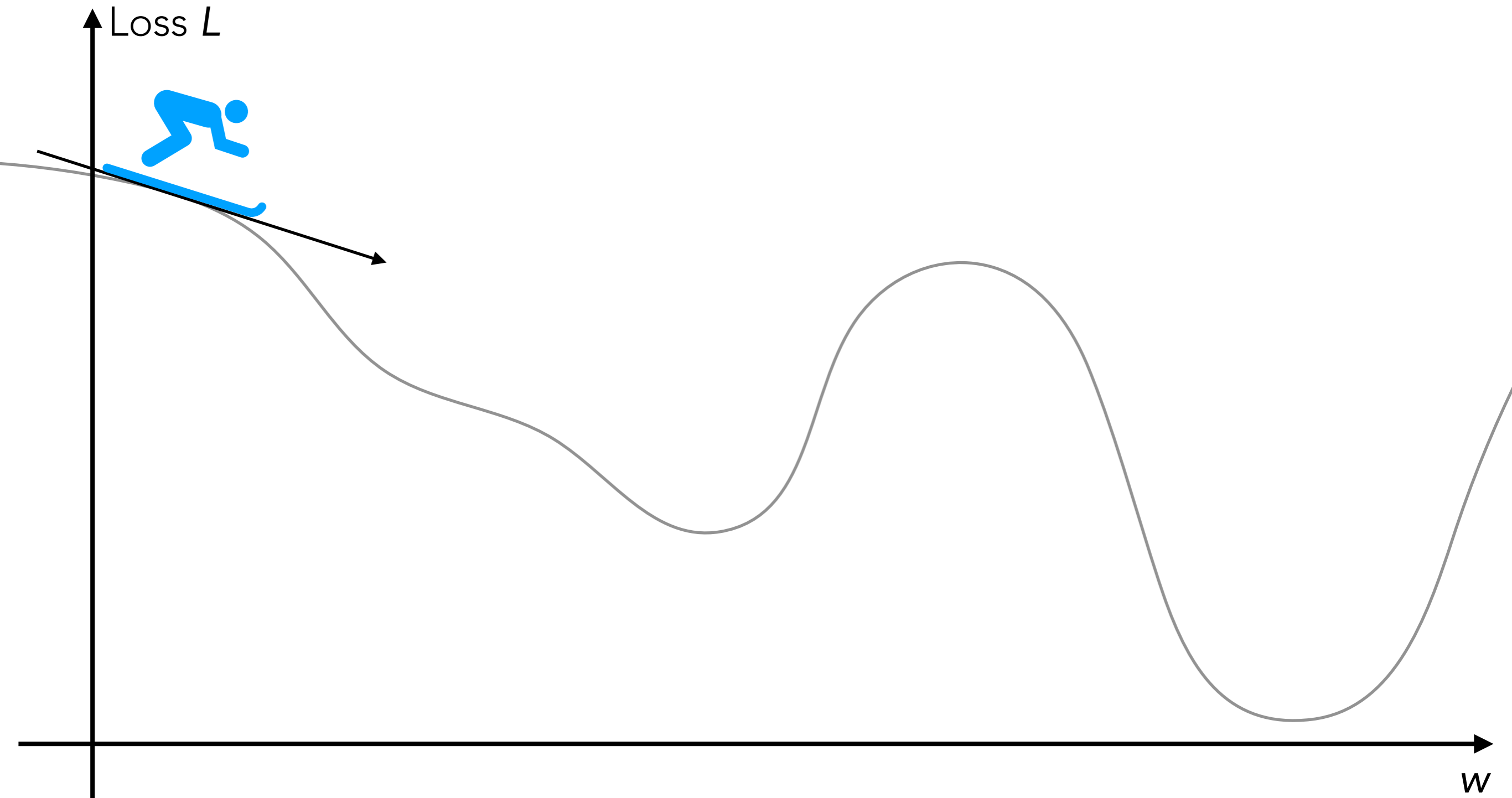
In general: the skier should move in *opposite* direction of derivative

In higher dimensions, this is called **gradient descent**

(derivative in higher dimensions: **gradient**)

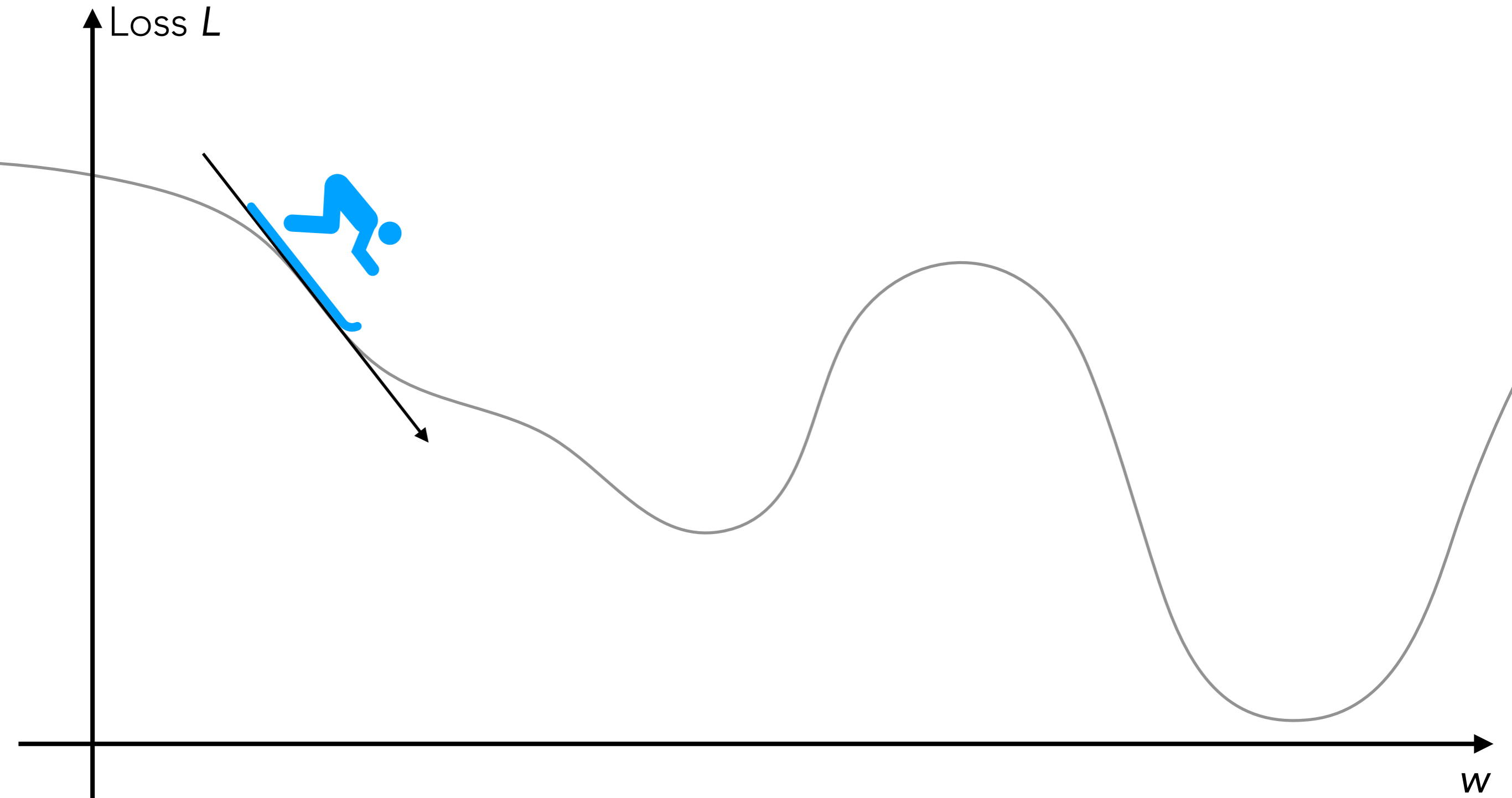
Training a Deep Net

Suppose the neural network has a single real number parameter w



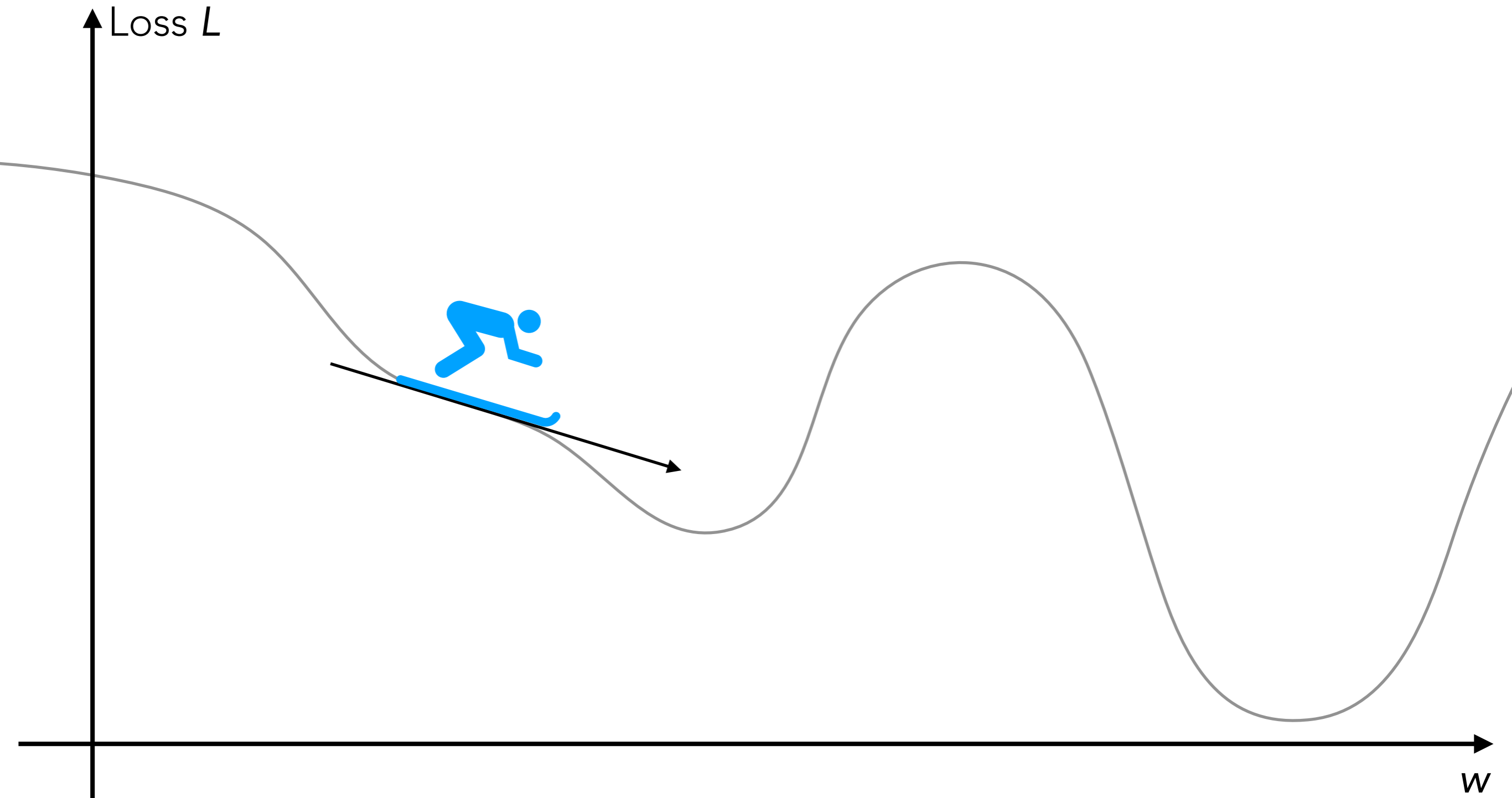
Training a Deep Net

Suppose the neural network has a single real number parameter w



Training a Deep Net

Suppose the neural network has a single real number parameter w



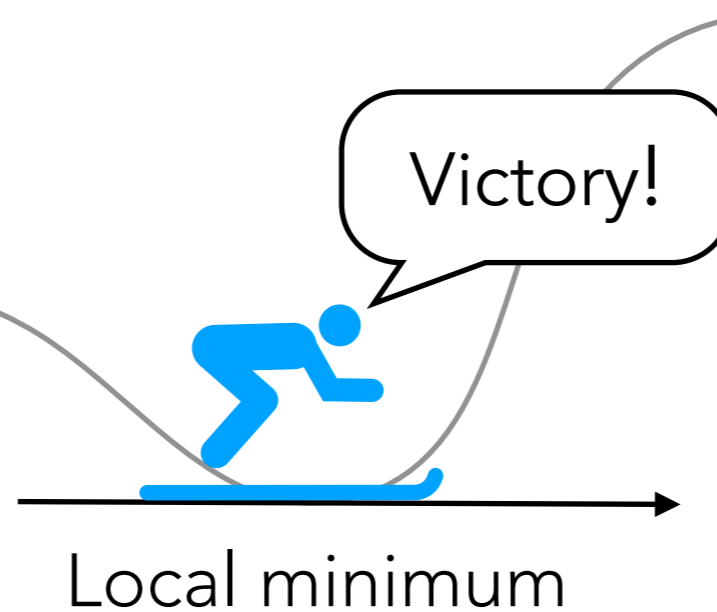
Training a Deep Net

Suppose the neural network has a single real number parameter w

In general: not obvious what error landscape looks like!
→ we wouldn't know there's a better solution beyond the hill

Popular optimizers
(e.g., Adam, RMSProp,
Lookahead) are variants
of gradient descent

The optimizer is the skier!

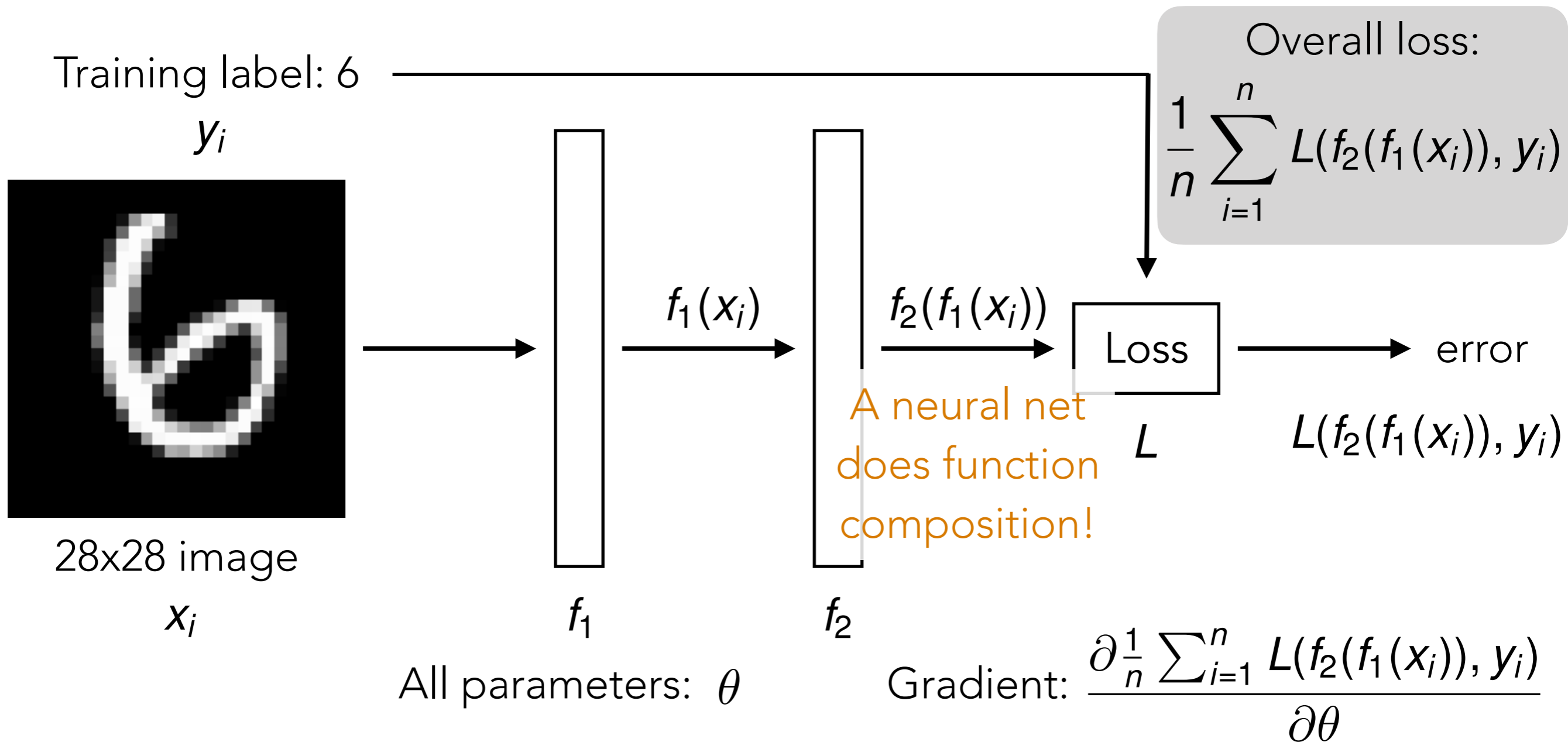


Better
solution

In very high-dimensional parameter spaces, local minima can be rare but we might get stuck in parts of the error landscape where the slope downwards is very gradual/not steep

w

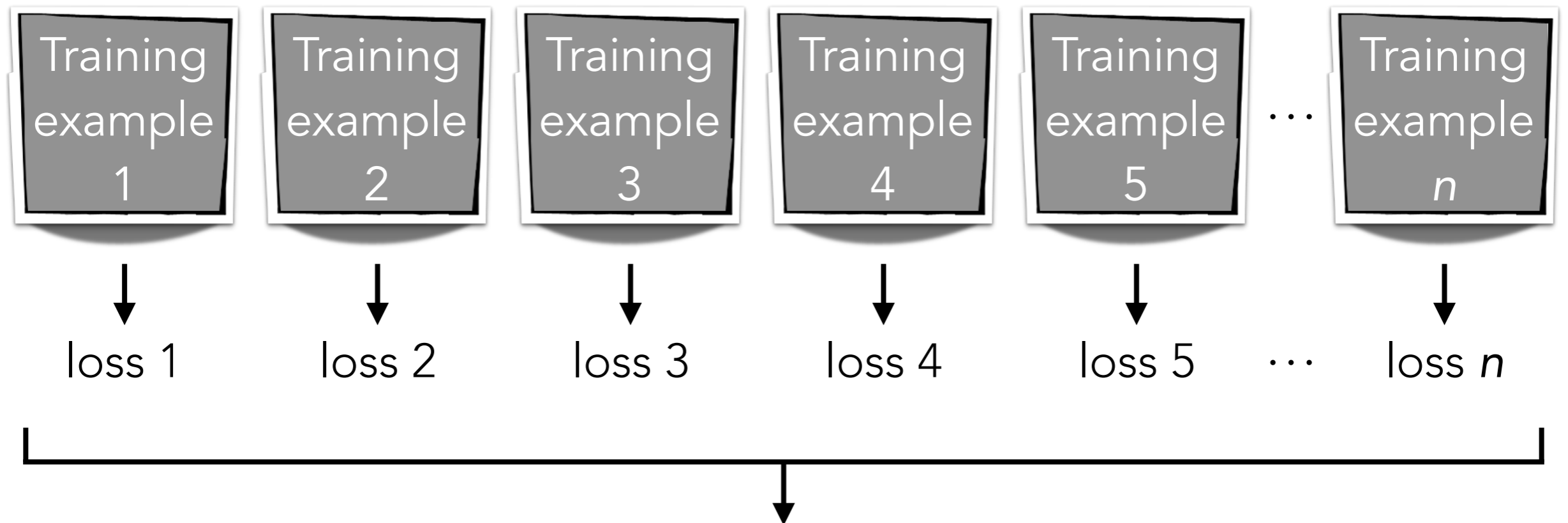
Handwritten Digit Recognition



Automatic differentiation is crucial in learning deep nets!

Careful derivative chain rule calculation: **back-propagation**

Gradient Descent

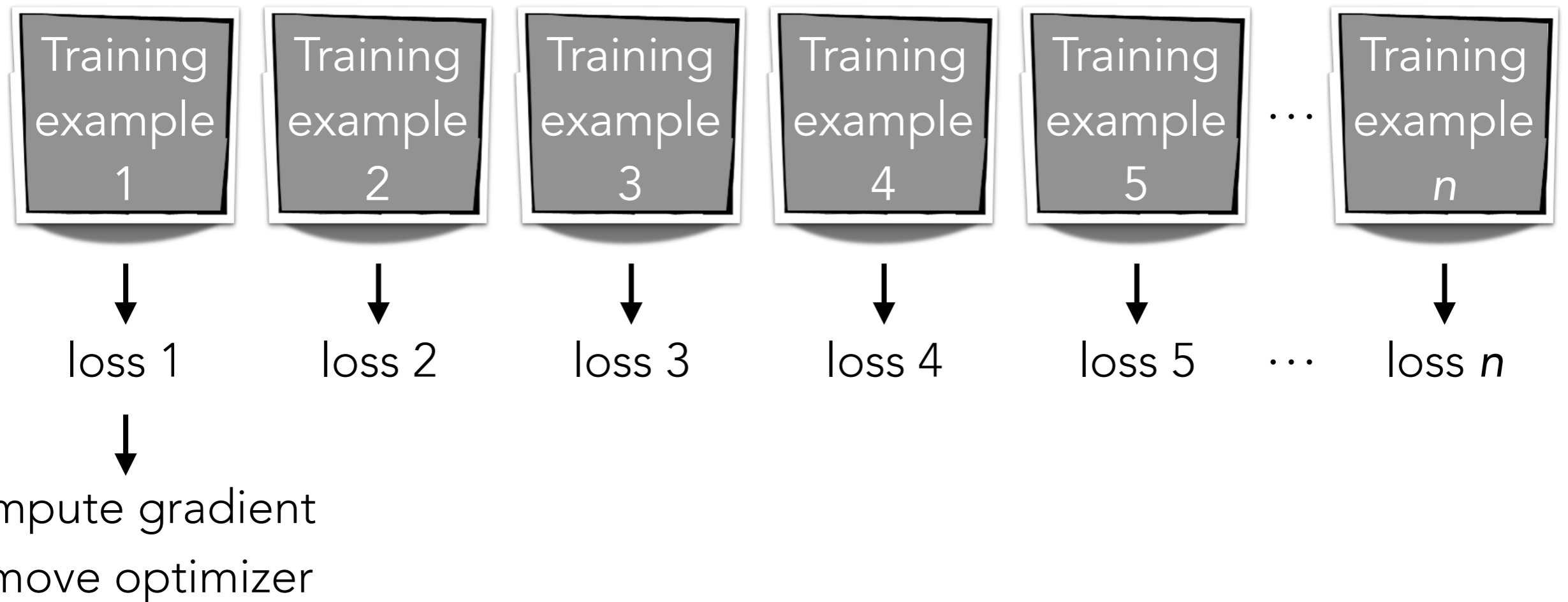


We have to compute lots of gradients to help the optimizer know where to go!

average loss
↓
compute gradient & move optimizer

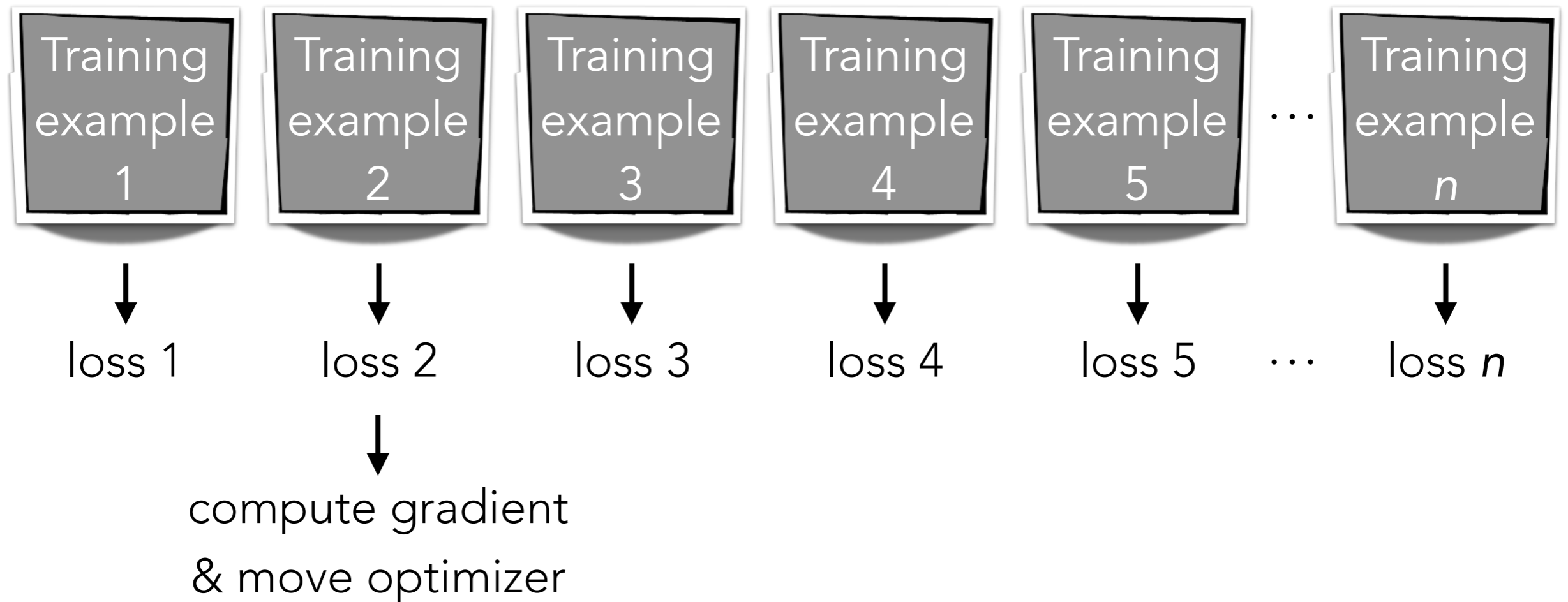
Computing gradients using all the training data seems really expensive!

Stochastic Gradient Descent (SGD)



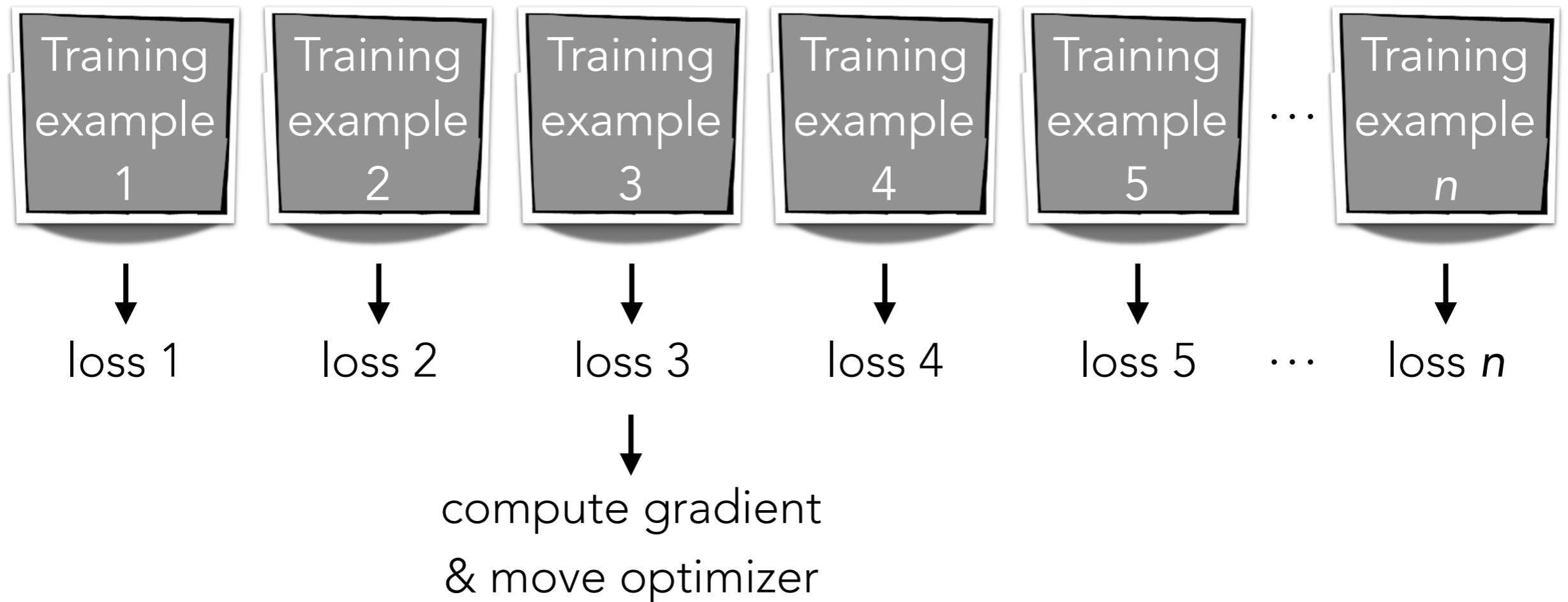
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the "full" gradient)

Stochastic Gradient Descent (SGD)



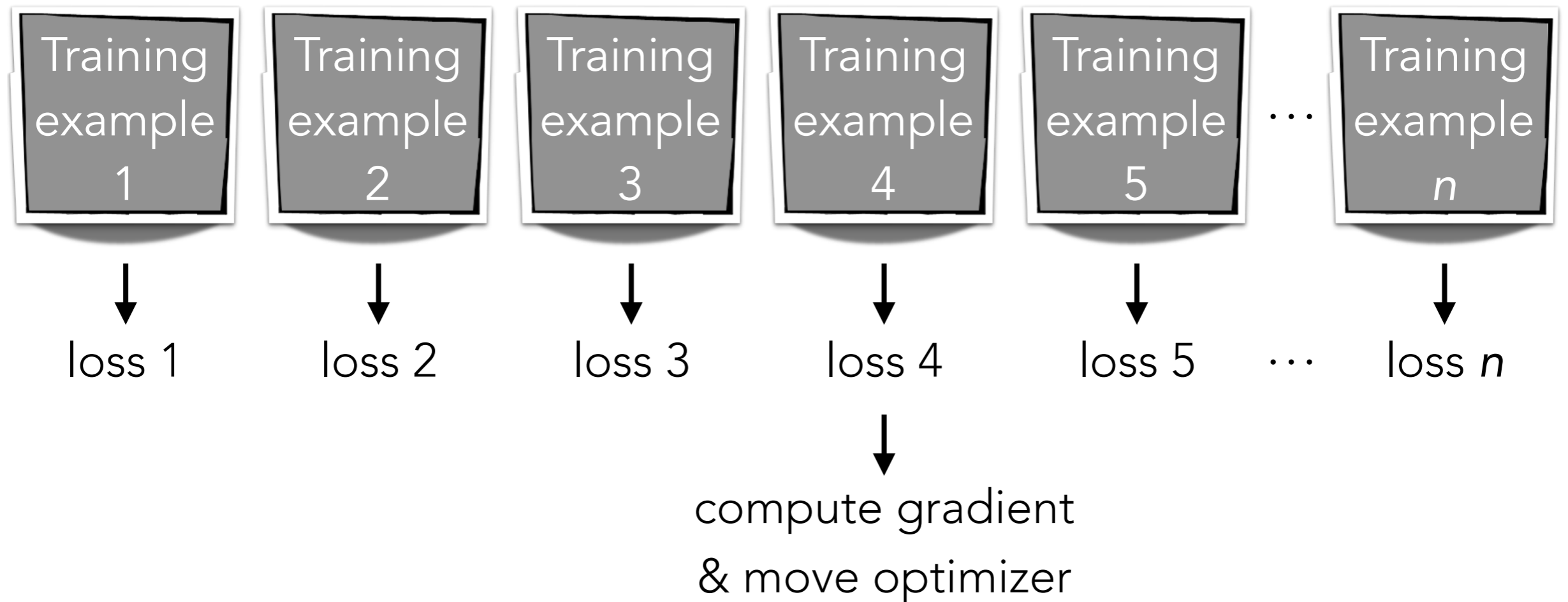
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the “full” gradient)

Stochastic Gradient Descent (SGD)



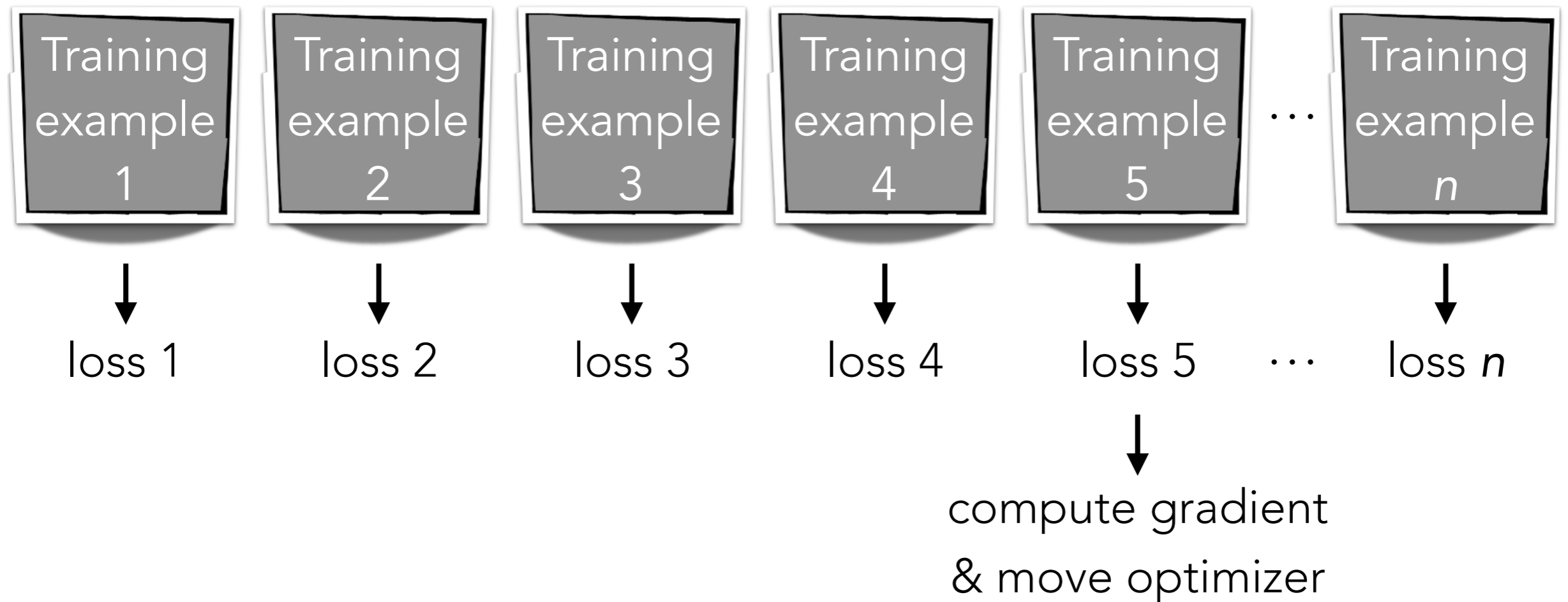
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the "full" gradient)

Stochastic Gradient Descent (SGD)



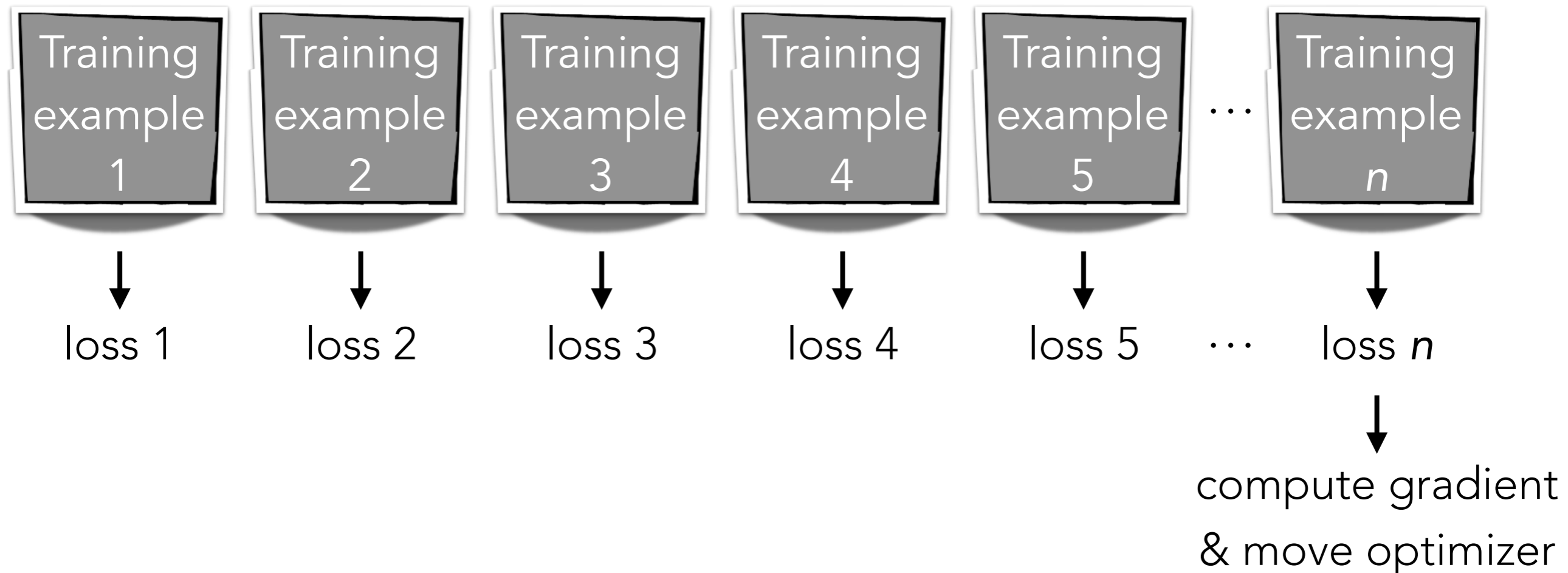
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the "full" gradient)

Stochastic Gradient Descent (SGD)



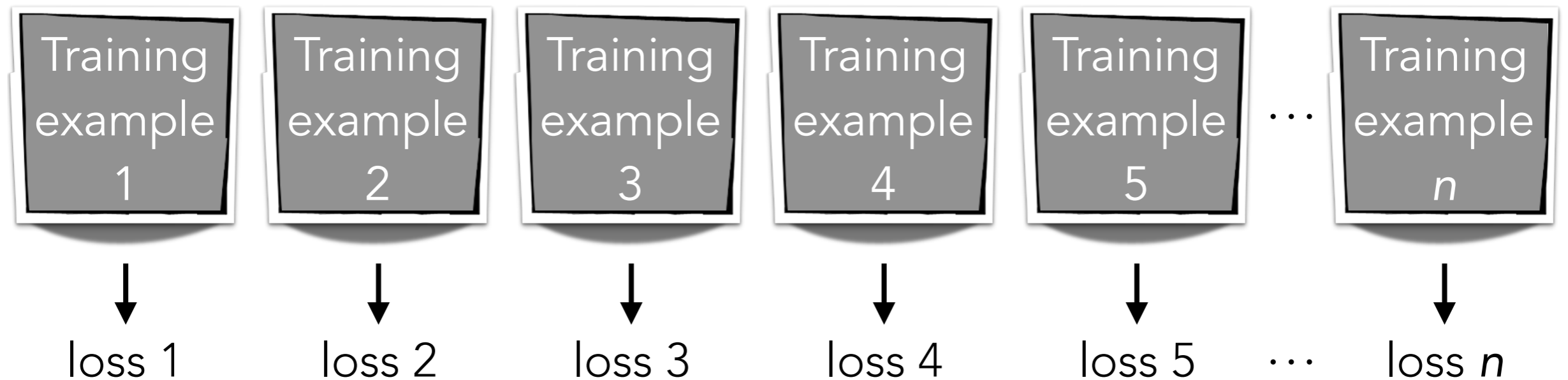
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the "full" gradient)

Stochastic Gradient Descent (SGD)



SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the "full" gradient)

Stochastic Gradient Descent (SGD)

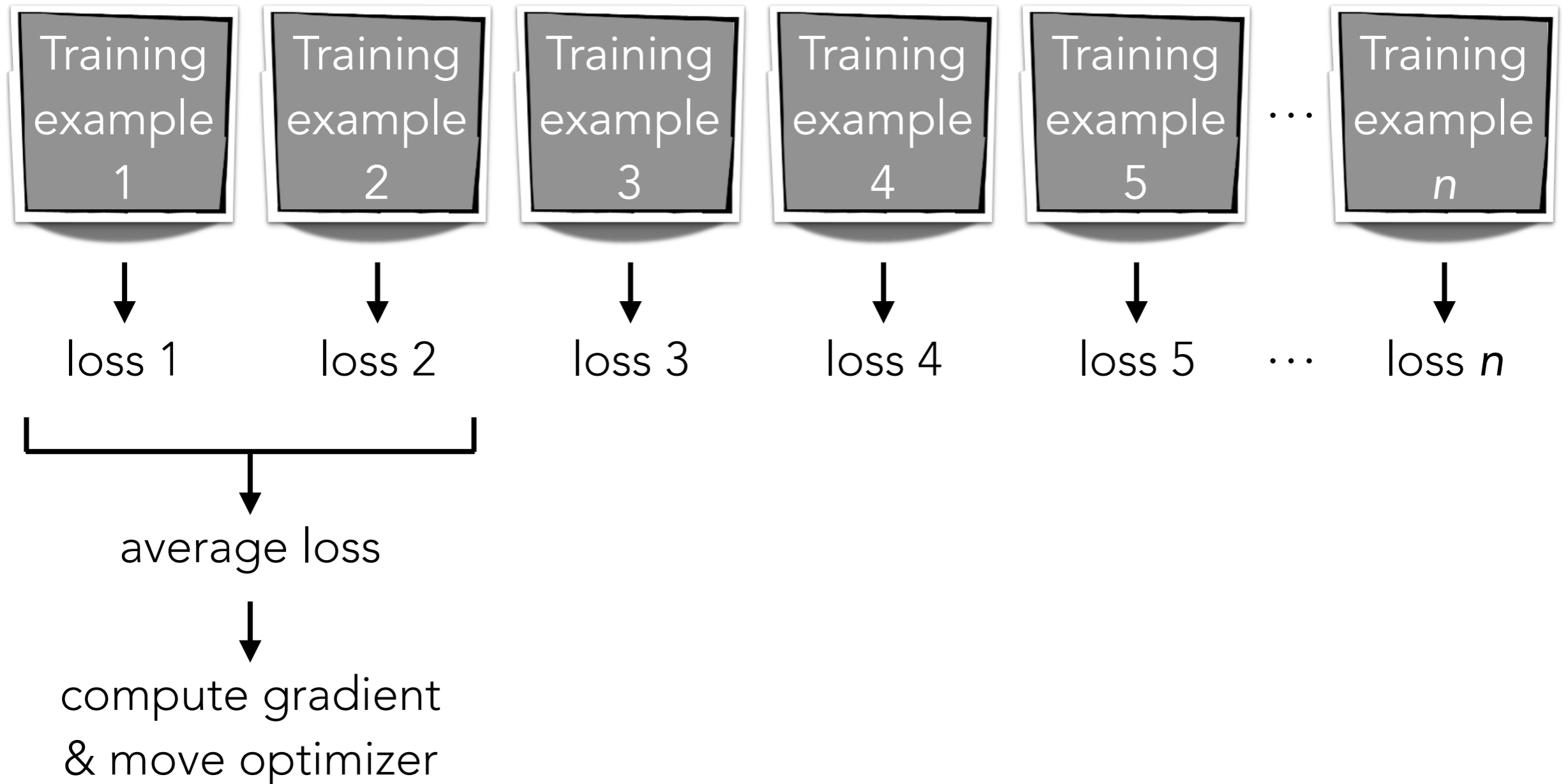


compute gradient
& move optimizer

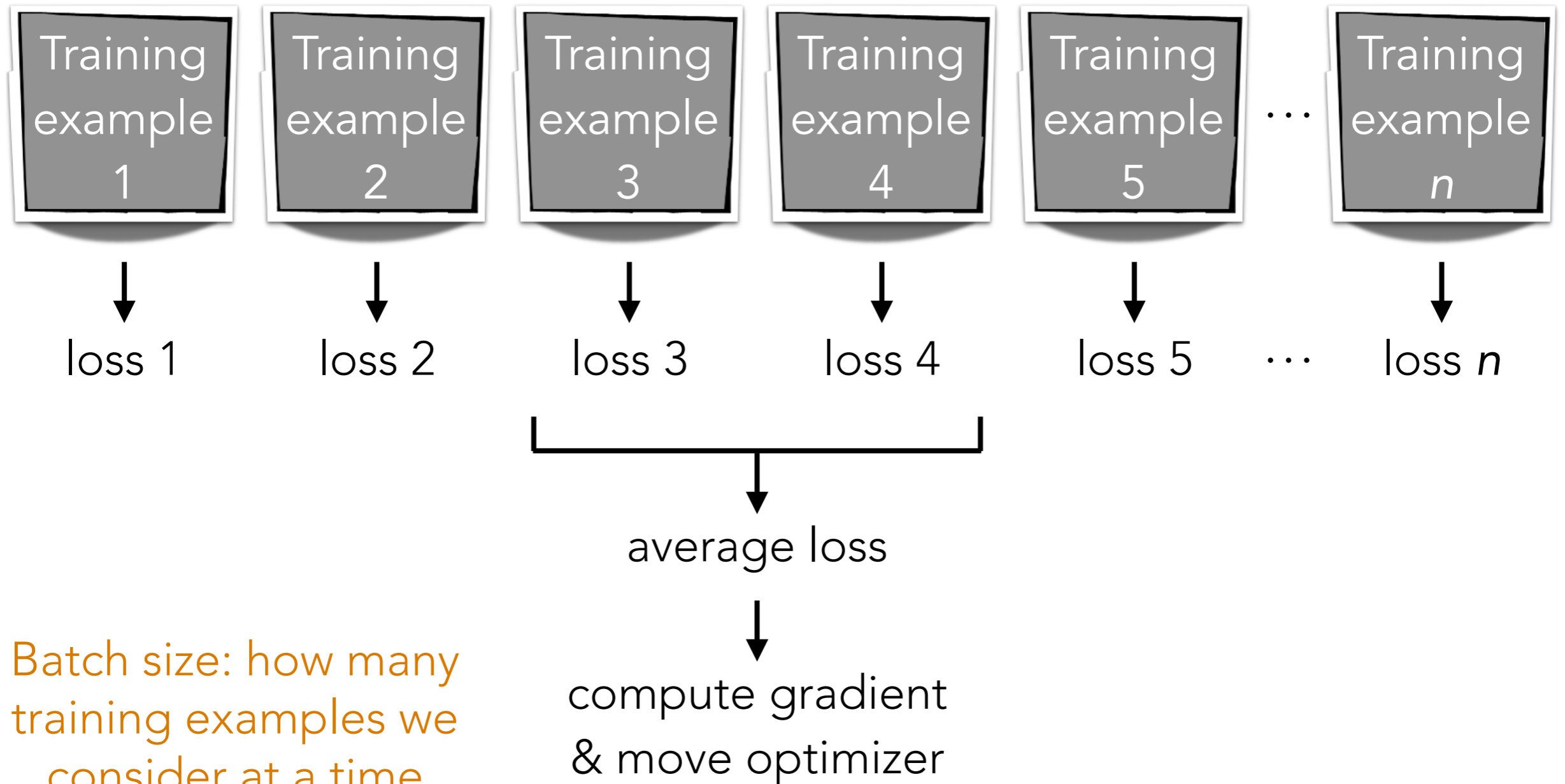
An epoch refers to 1 full pass through all
the training data

SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the "full" gradient)

Minibatch Gradient Descent



Minibatch Gradient Descent



Batch size: how many training examples we consider at a time (in this example: 2)

**Best optimizer? Best learning rate? Best
of epochs? Best batch size?**

Active area of research

Depends on problem, data, hardware, etc

Example: even with a GPU, you can get slow learning (slower than CPU!)
if you choose # epochs/batch size poorly!!!

UDA_pytorch_utils.py

A look at [UDA_pytorch_classifier_fit](#)