

Deep Learning and 95-865 Wrap-Up

nearly all slides by George Chen (CMU)

1 slide by Phillip Isola (OpenAI, UC Berkeley)

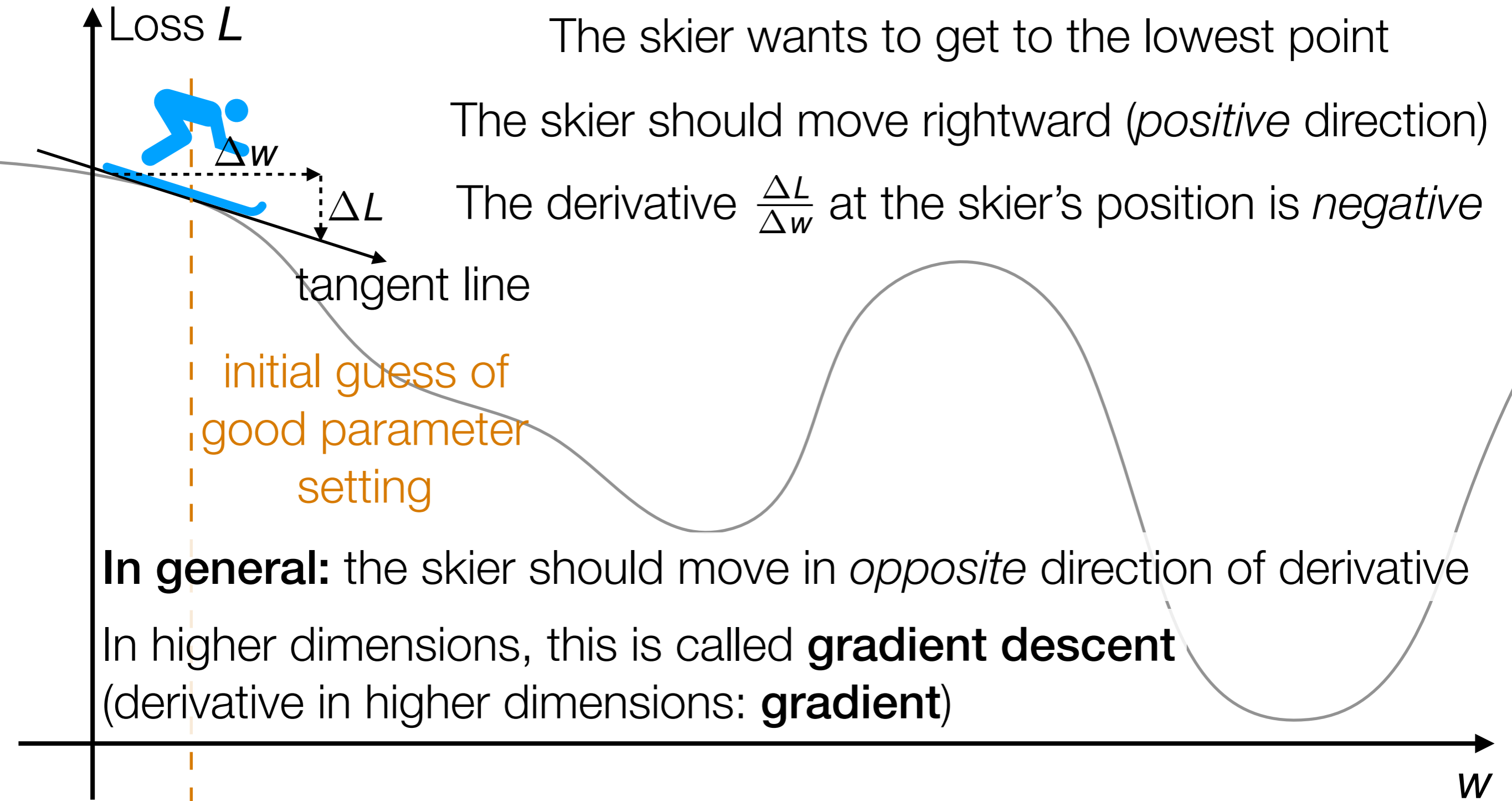
Today

- How learning a deep net works
- A bunch of deep learning topics we didn't cover
- Course wrap-up

Learning a Deep Net

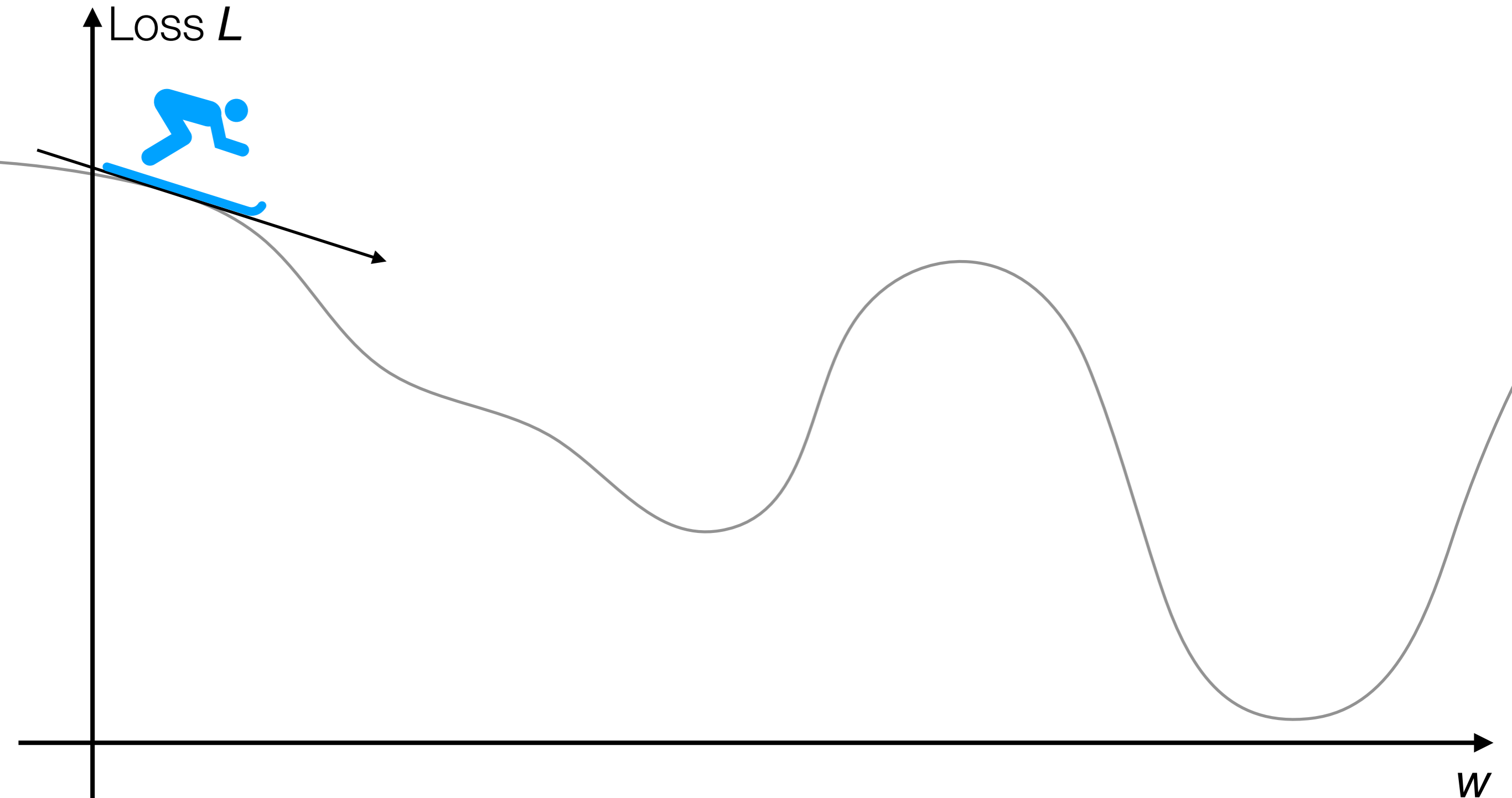
Gradient Descent

Suppose the neural network has a single real number parameter w



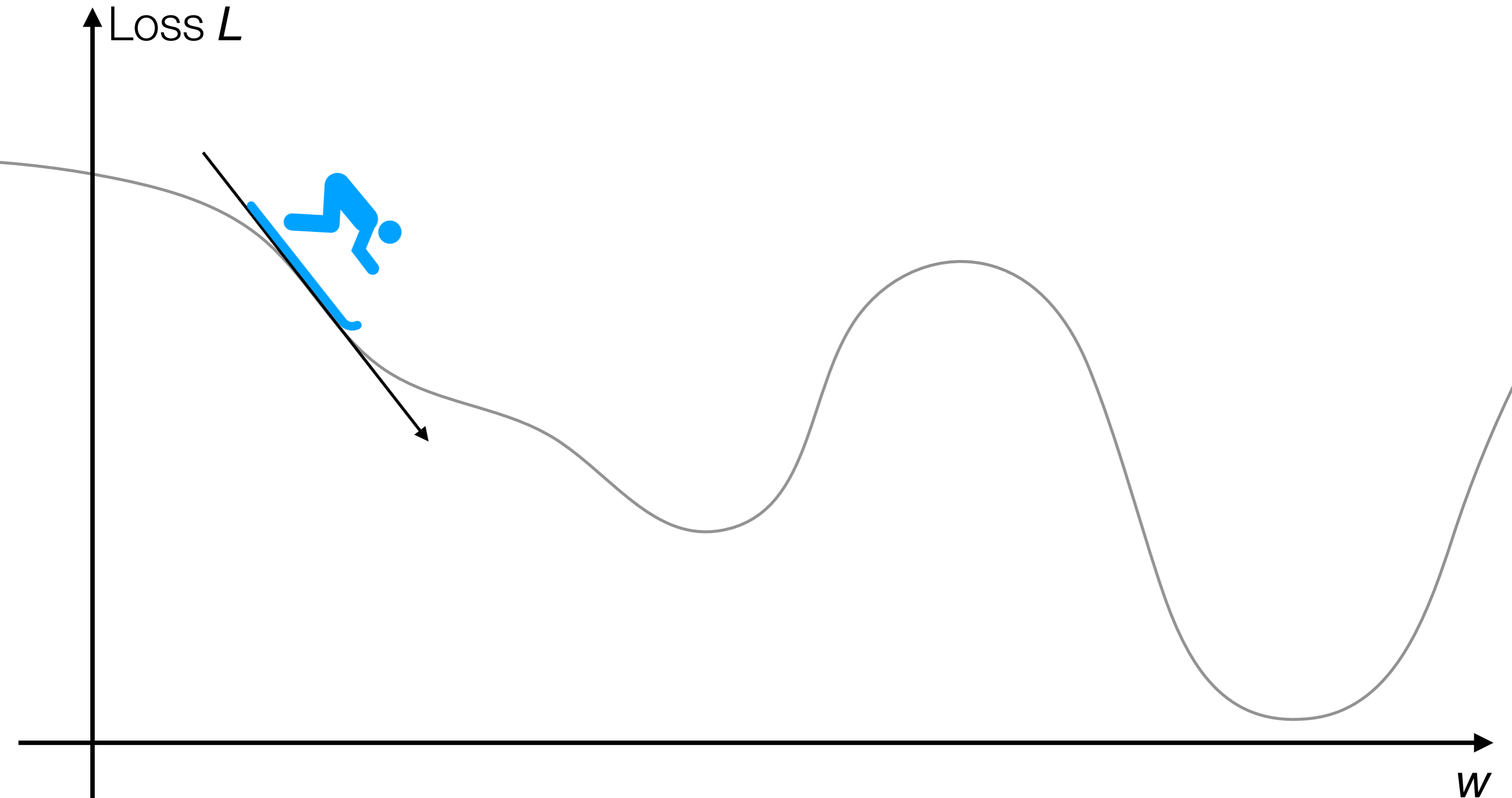
Gradient Descent

Suppose the neural network has a single real number parameter w



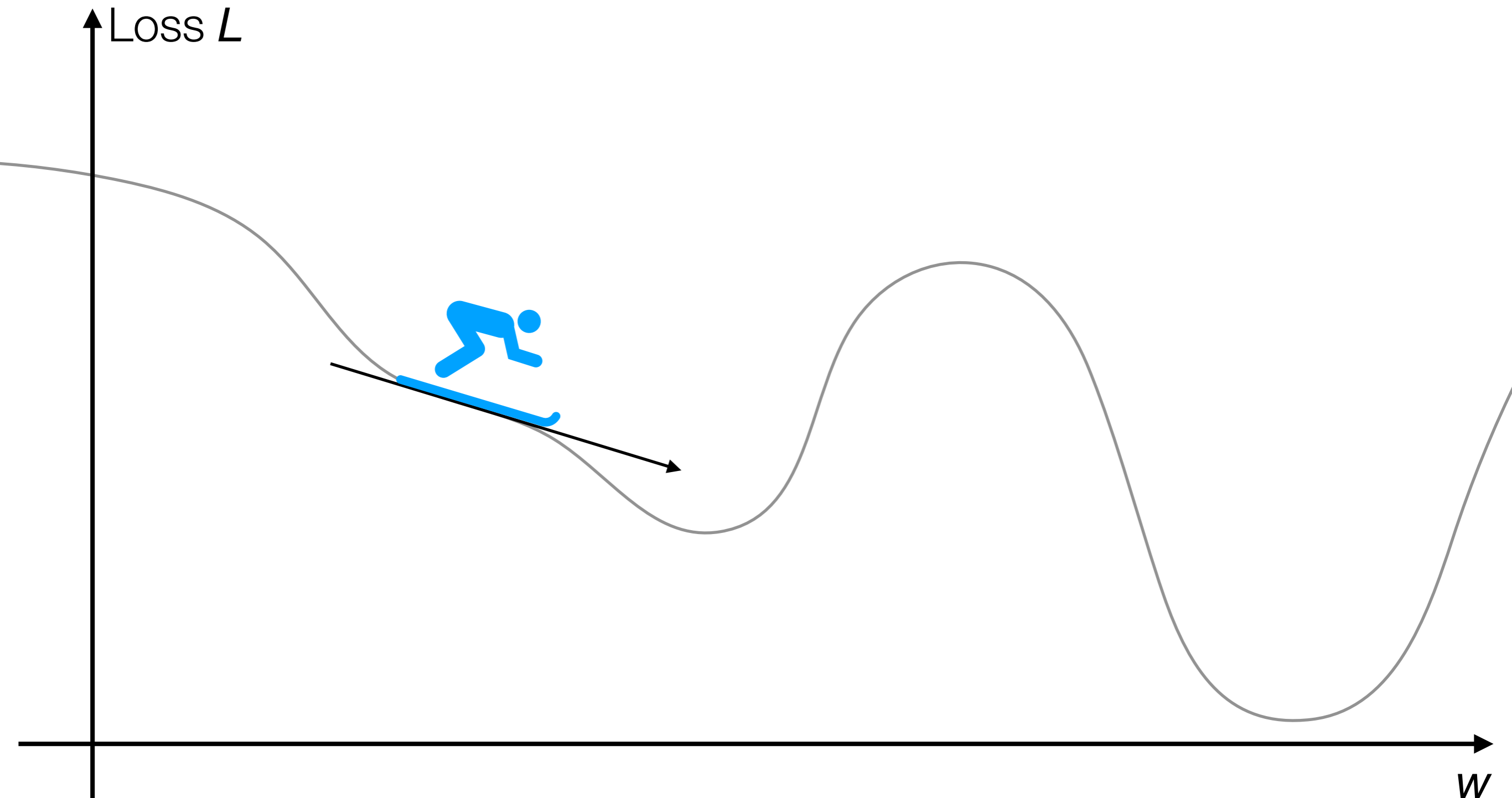
Gradient Descent

Suppose the neural network has a single real number parameter w



Gradient Descent

Suppose the neural network has a single real number parameter w

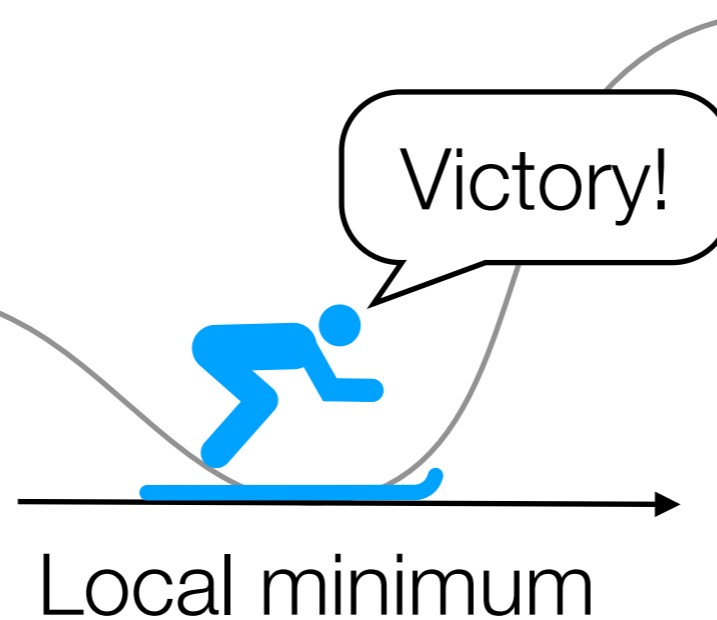


Gradient Descent

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., RMSprop,
ADAM, AdaGrad,
AdaDelta) are variants
of gradient descent

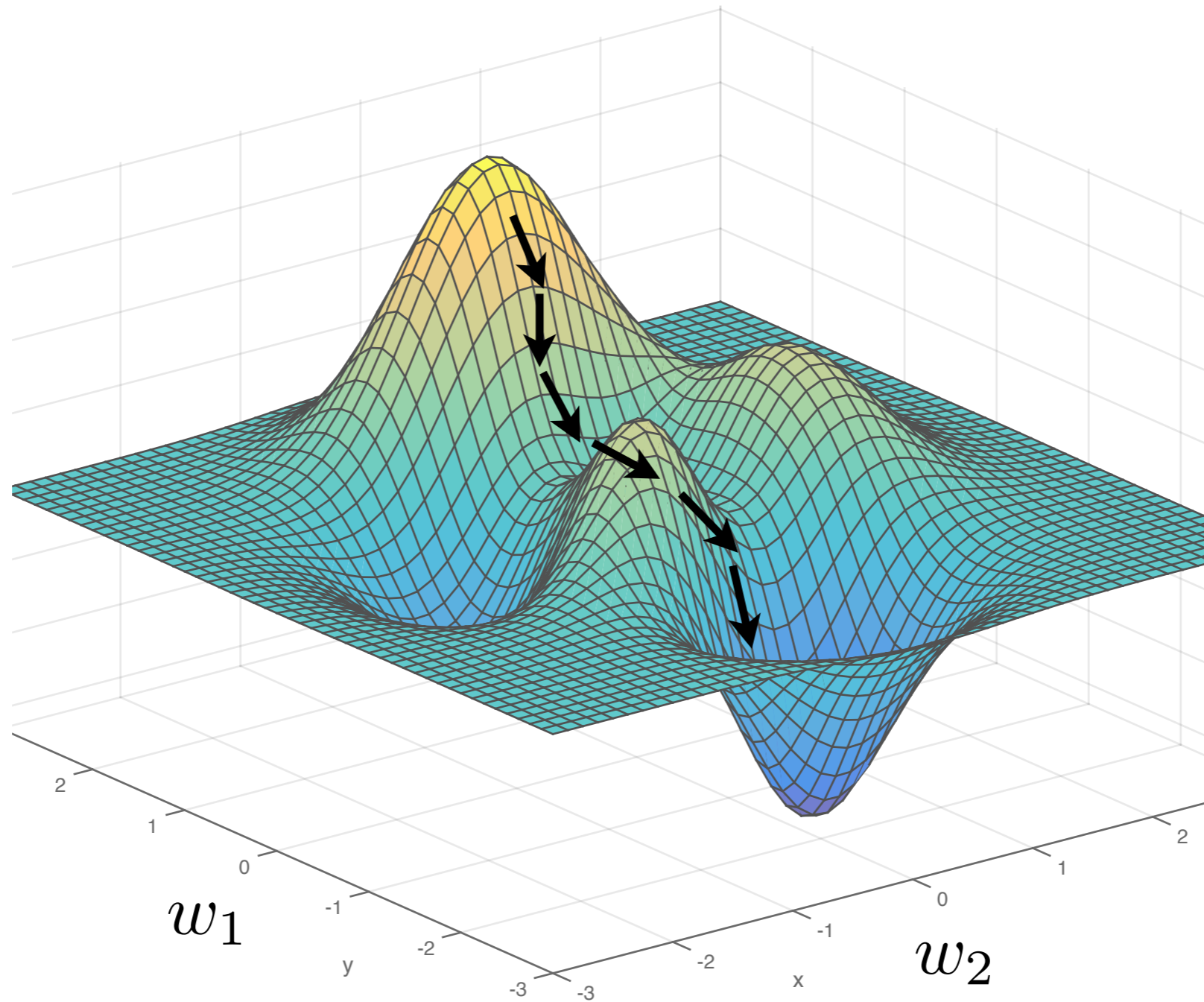


In practice: local minimum often good enough

Gradient Descent

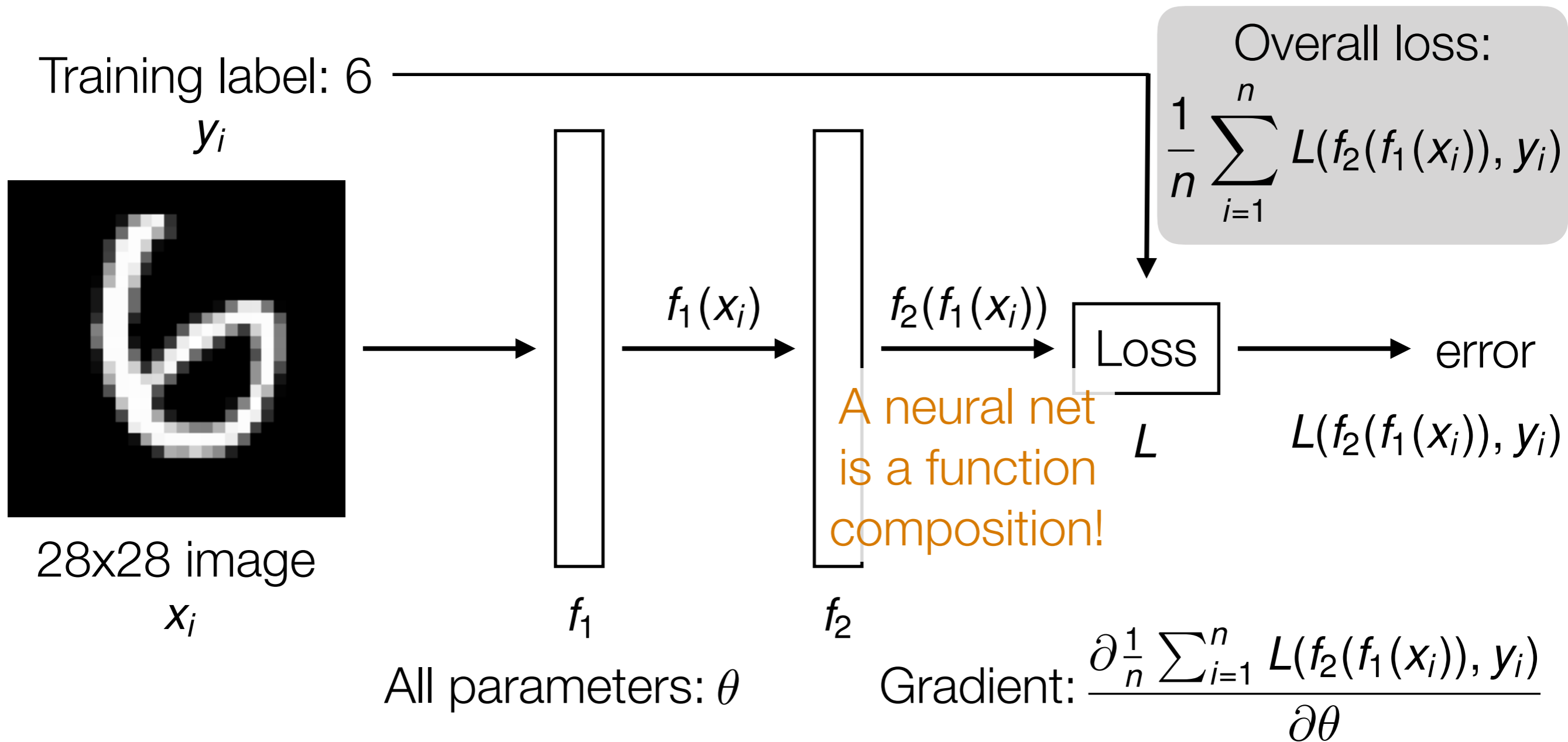
2D example

$L(\mathbf{w})$



Remark: In practice, deep nets often have $>$ *millions* of parameters, so *very* high-dimensional gradient descent

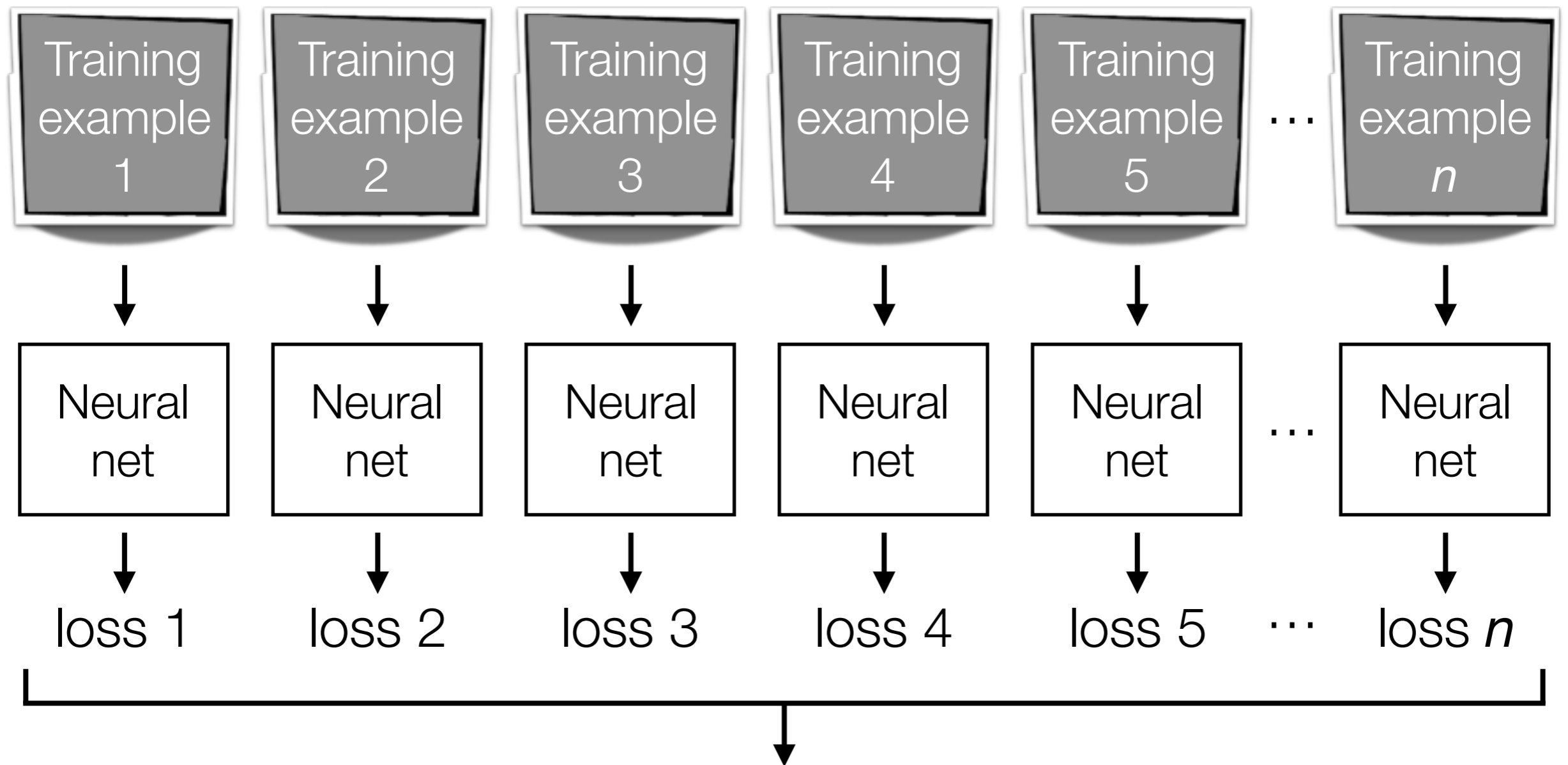
Handwritten Digit Recognition



Automatic differentiation is crucial in learning deep nets!

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

Gradient Descent

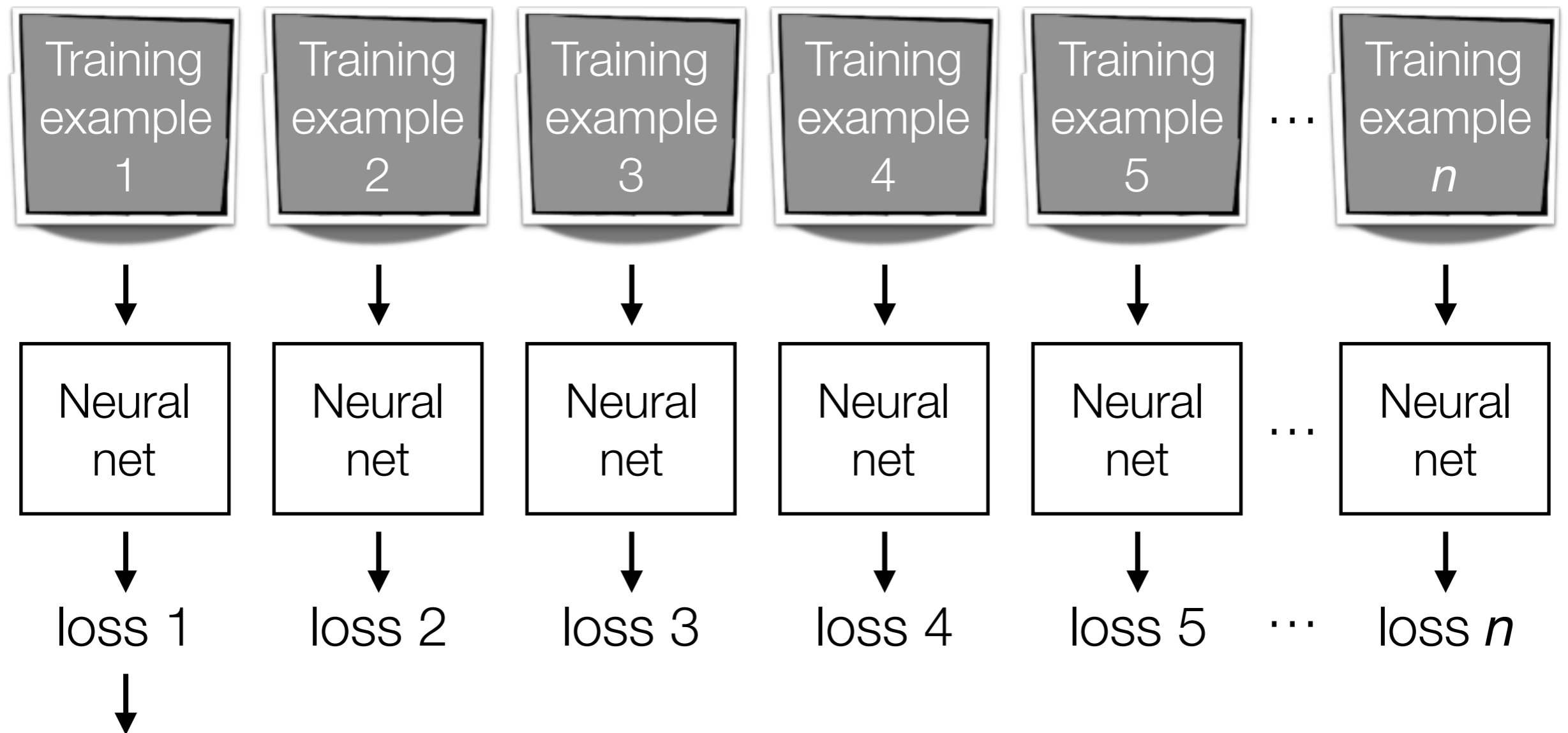


We have to compute lots of gradients to help the skier know where to go!

average loss
↓
compute gradient and move skier

Computing gradients using all the training data seems really expensive!

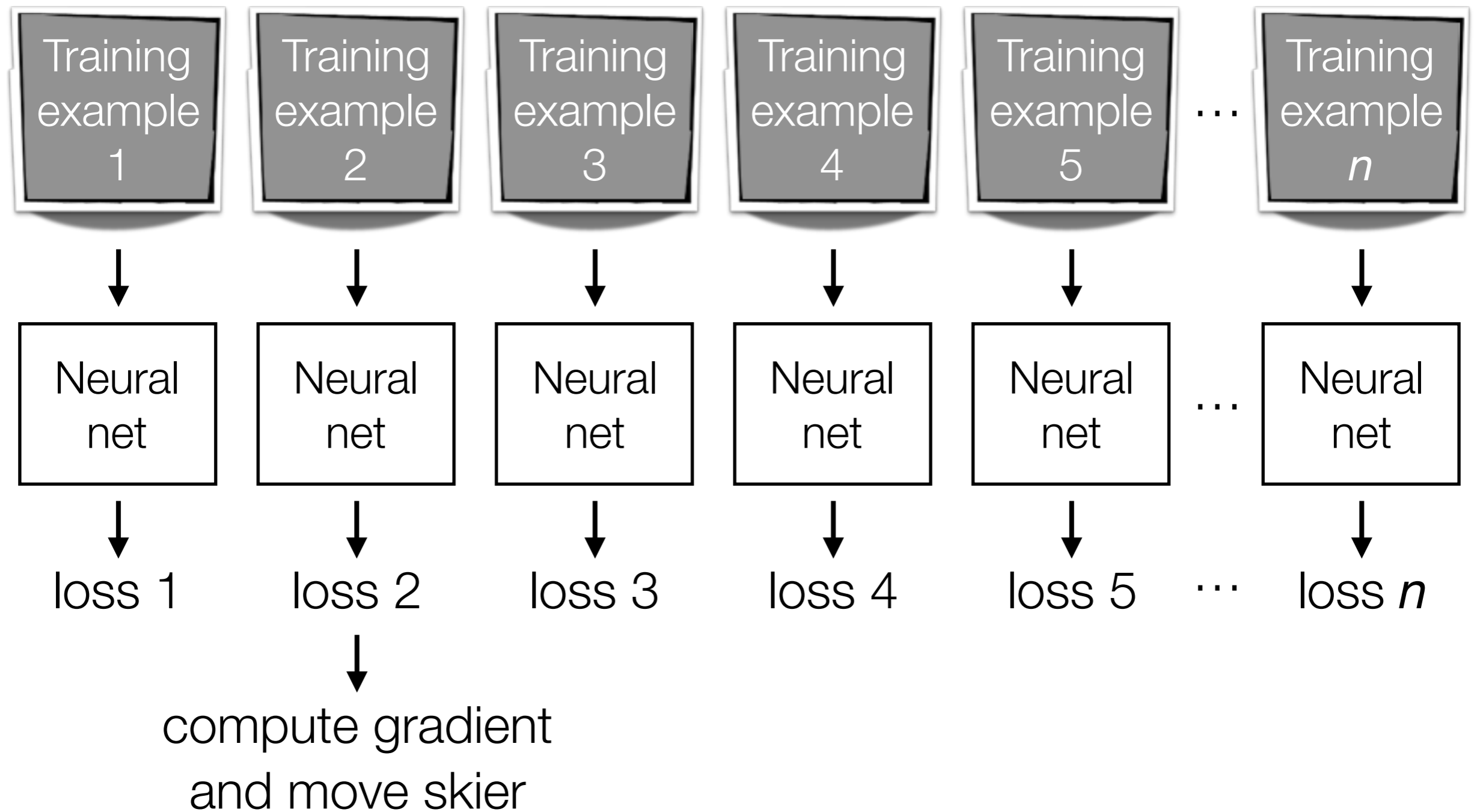
Stochastic Gradient Descent (SGD)



compute gradient
and move skier

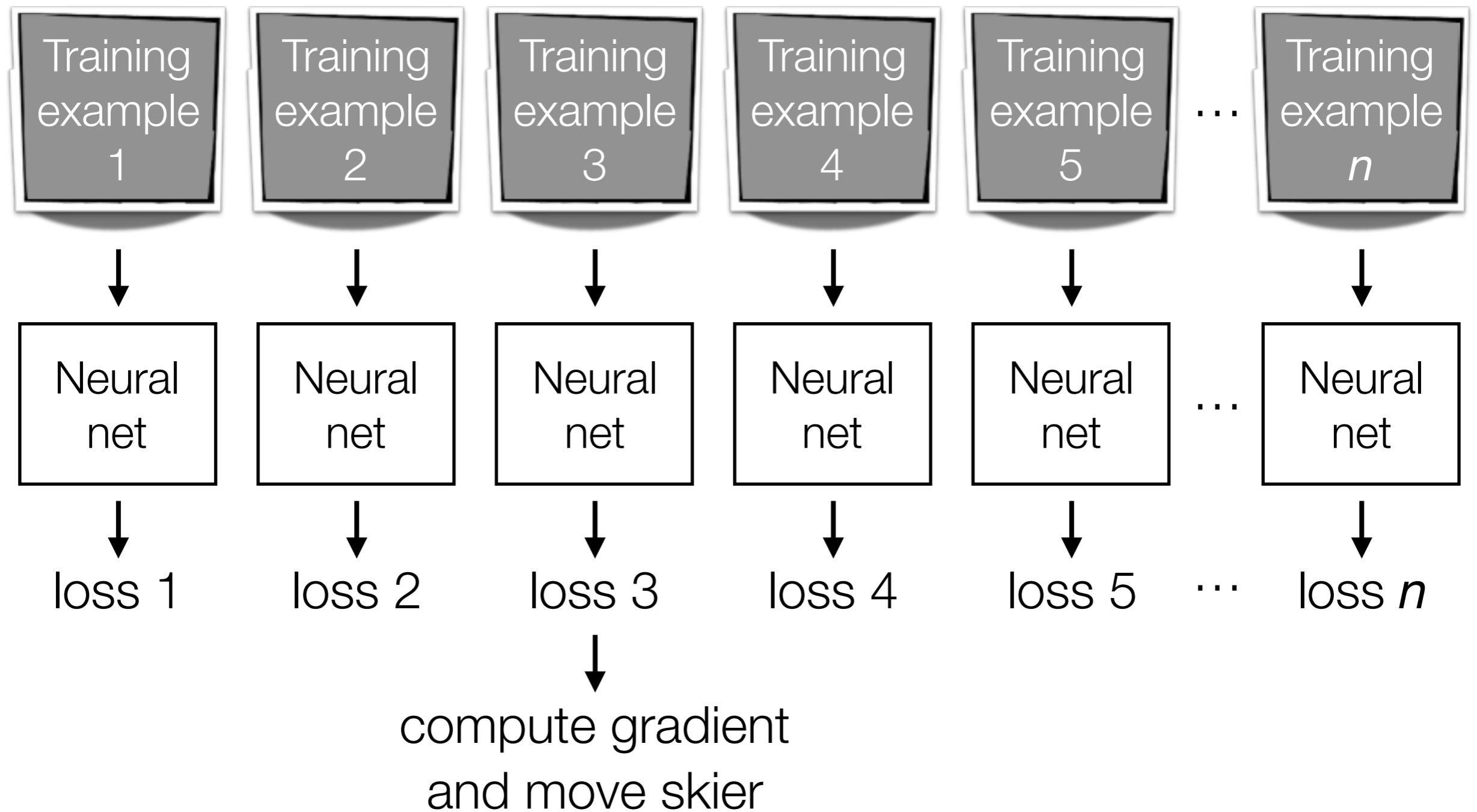
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the “full” gradient)

Stochastic Gradient Descent (SGD)



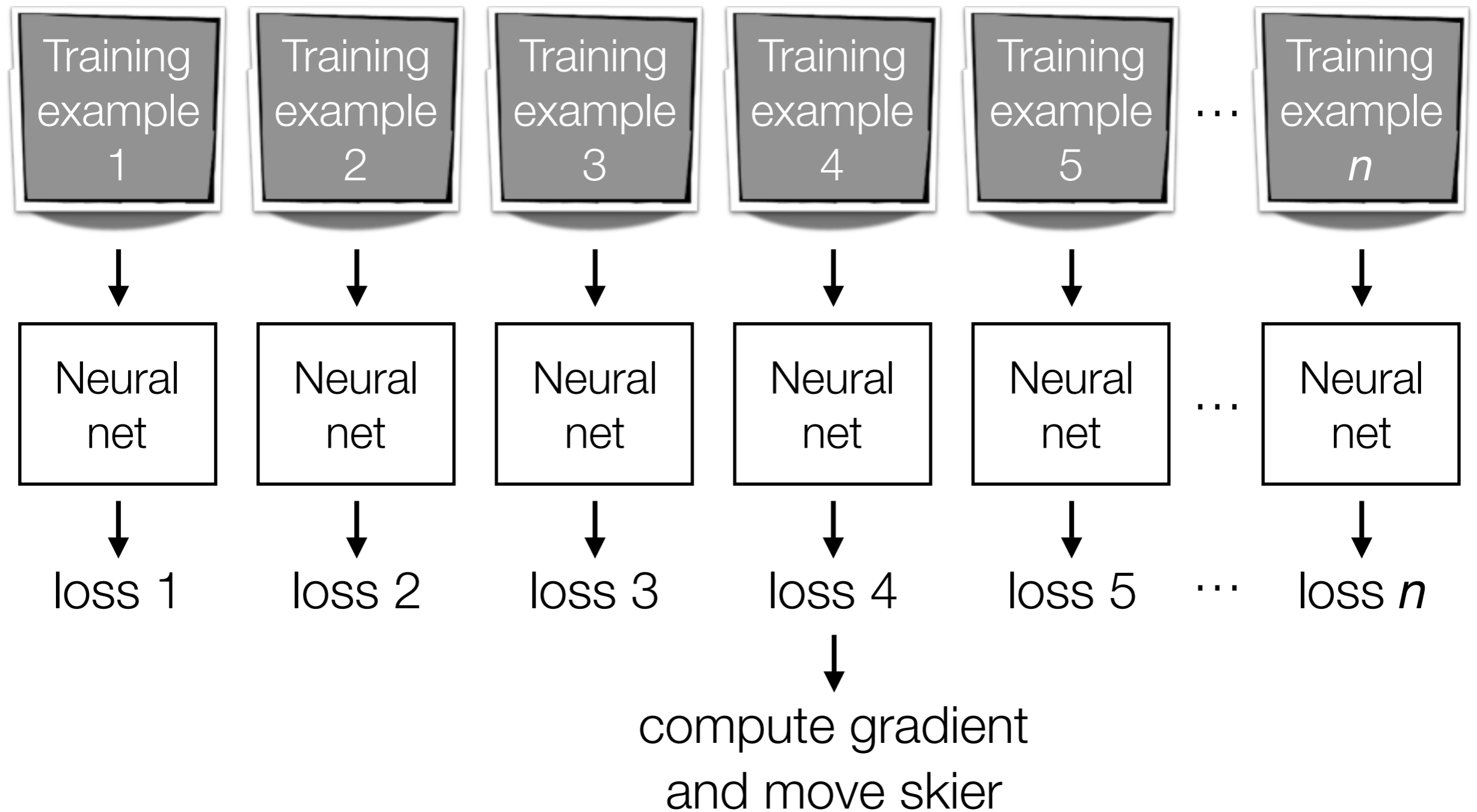
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the “full” gradient)

Stochastic Gradient Descent (SGD)



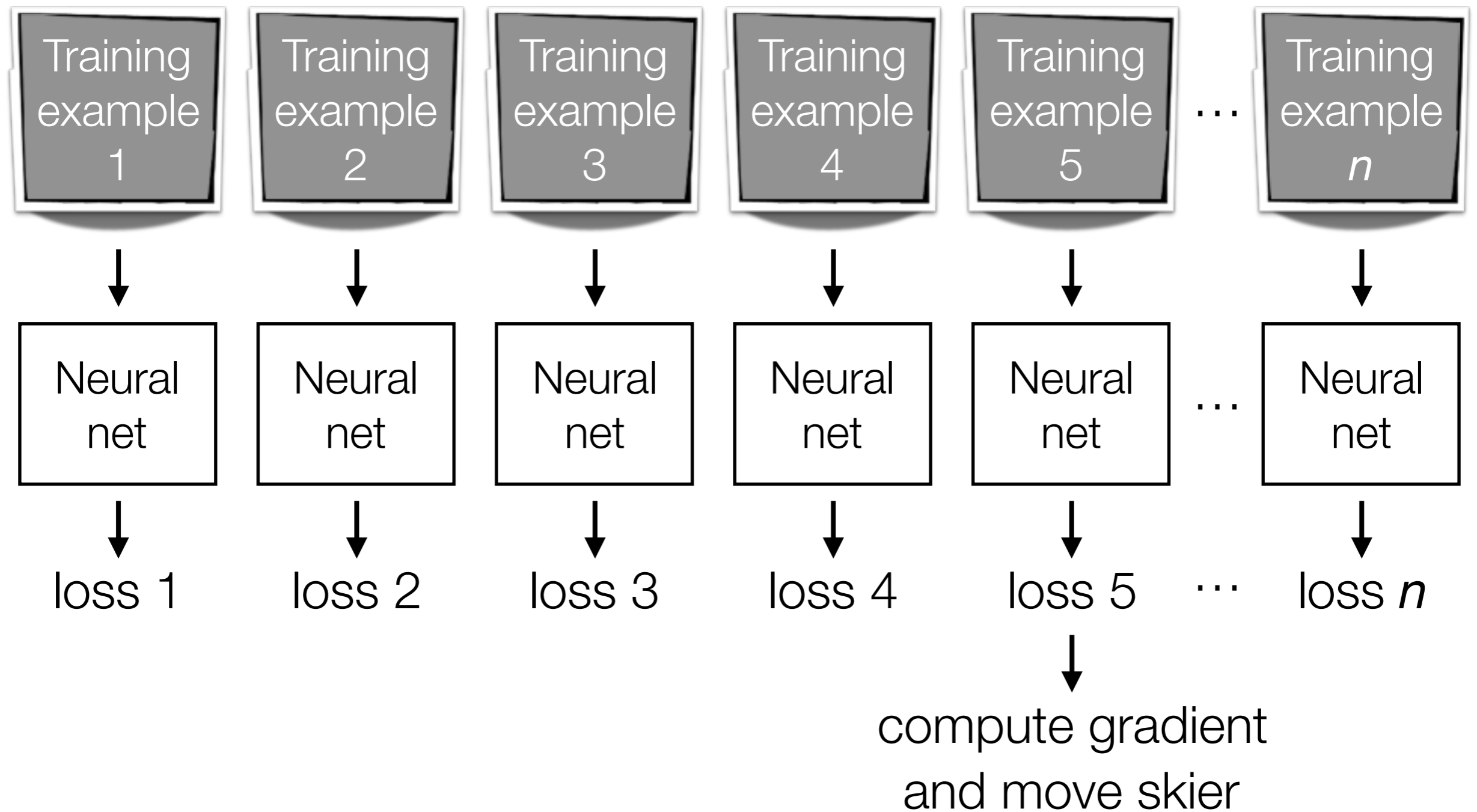
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the “full” gradient)

Stochastic Gradient Descent (SGD)



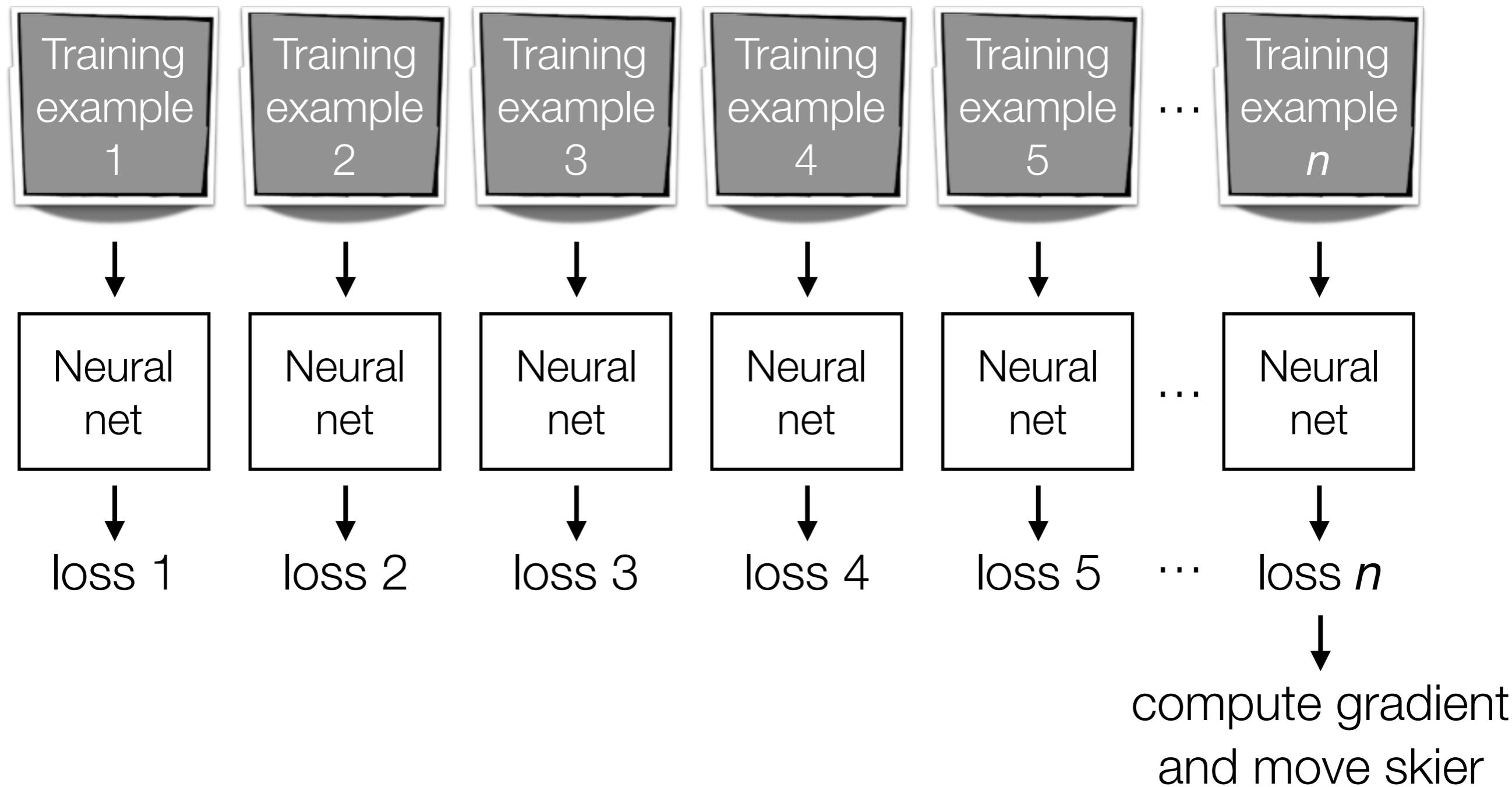
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the “full” gradient)

Stochastic Gradient Descent (SGD)



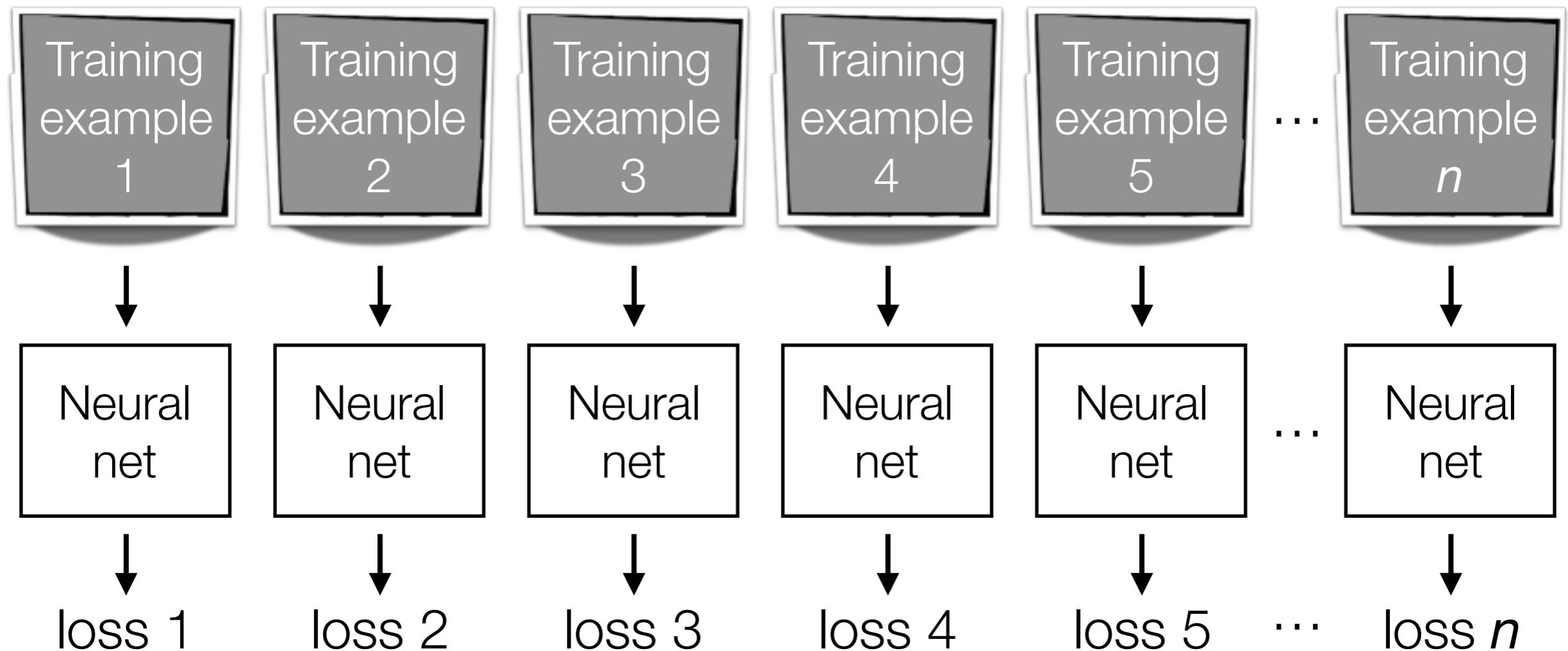
SGD: compute gradient using only 1 training example at a time
(can think of this gradient as a noisy approximation of the “full” gradient)

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)

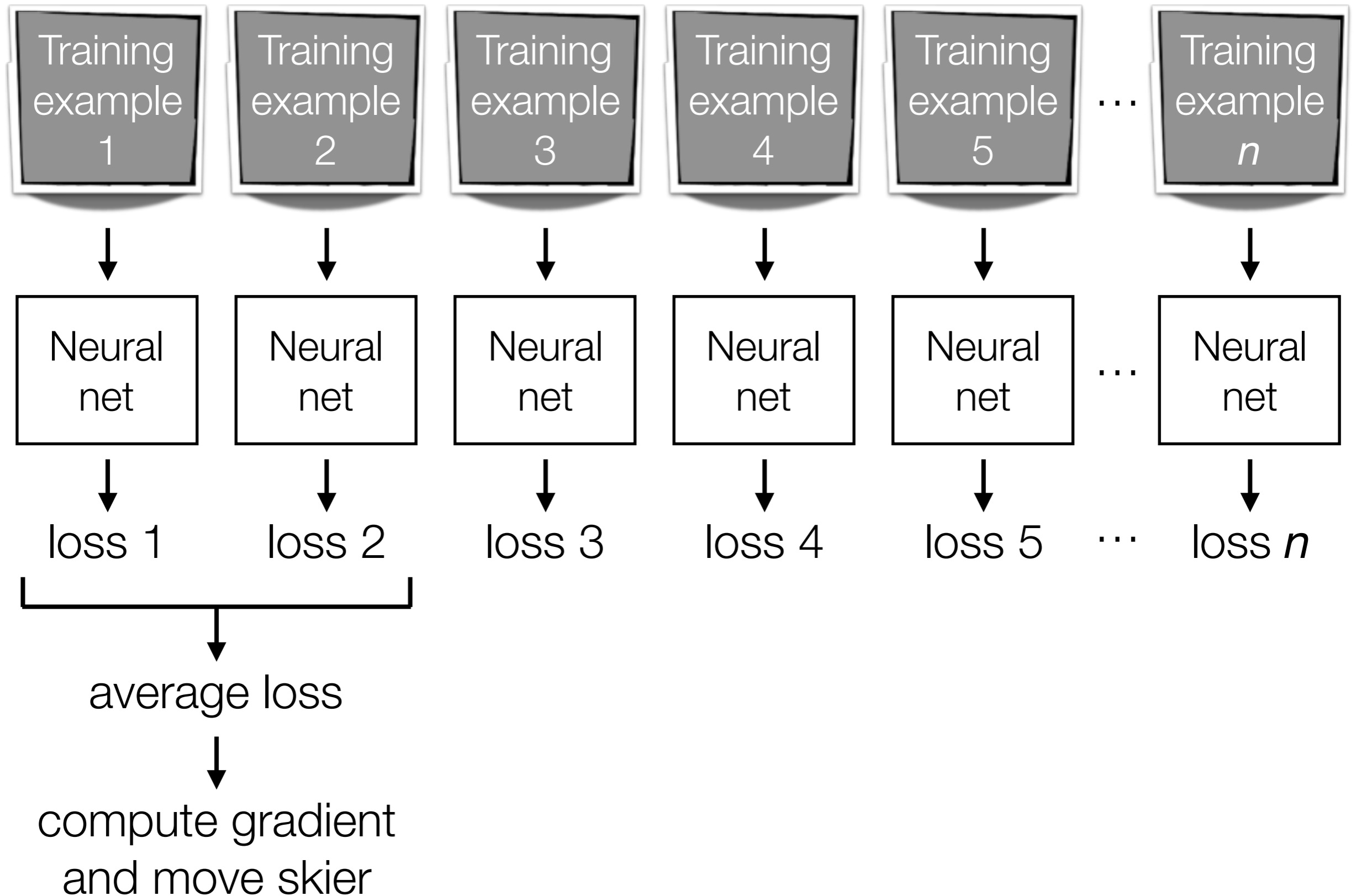


compute gradient
and move skier

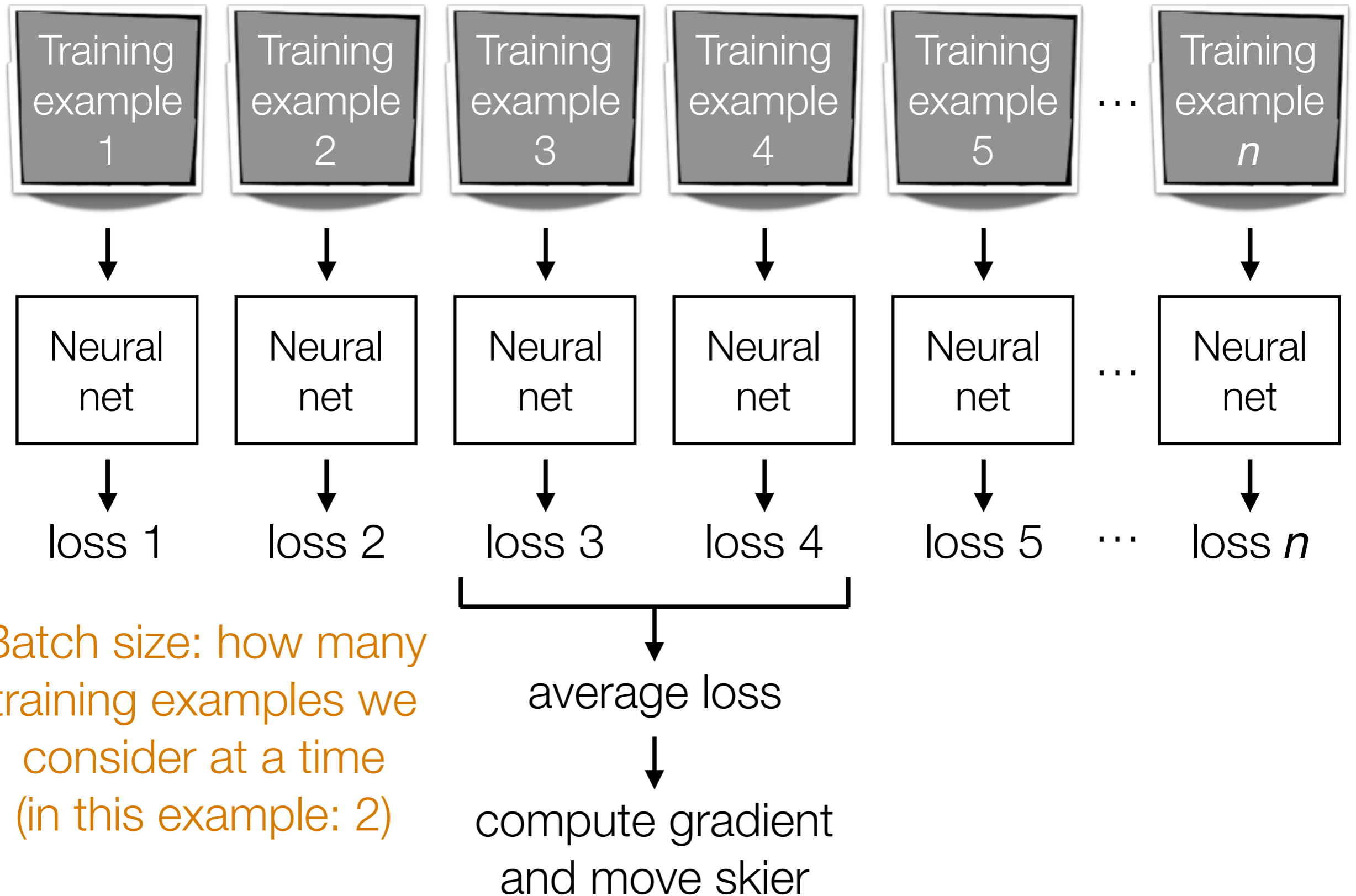
An epoch refers to 1 full pass
through all the training data

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

Mini-Batch Gradient Descent



Mini-Batch Gradient Descent



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

There's a lot more to deep learning that we didn't cover

Dealing with Small Datasets

Data augmentation: generate perturbed versions of your training data to get larger training dataset



Training image
Training label: cat



Mirrored
Still a cat!



Rotated & translated
Still a cat!

We just turned 1 training example in 3 training examples

Allowable perturbations depend on data
(e.g., for handwritten digits, rotating by 180 degrees would be bad: confuse 6's and 9's)

Dealing with Small Datasets

Fine tuning: if there's an existing pre-trained neural net, you could modify it for your problem that has a small dataset

Example: classify between Tesla's and Toyota's



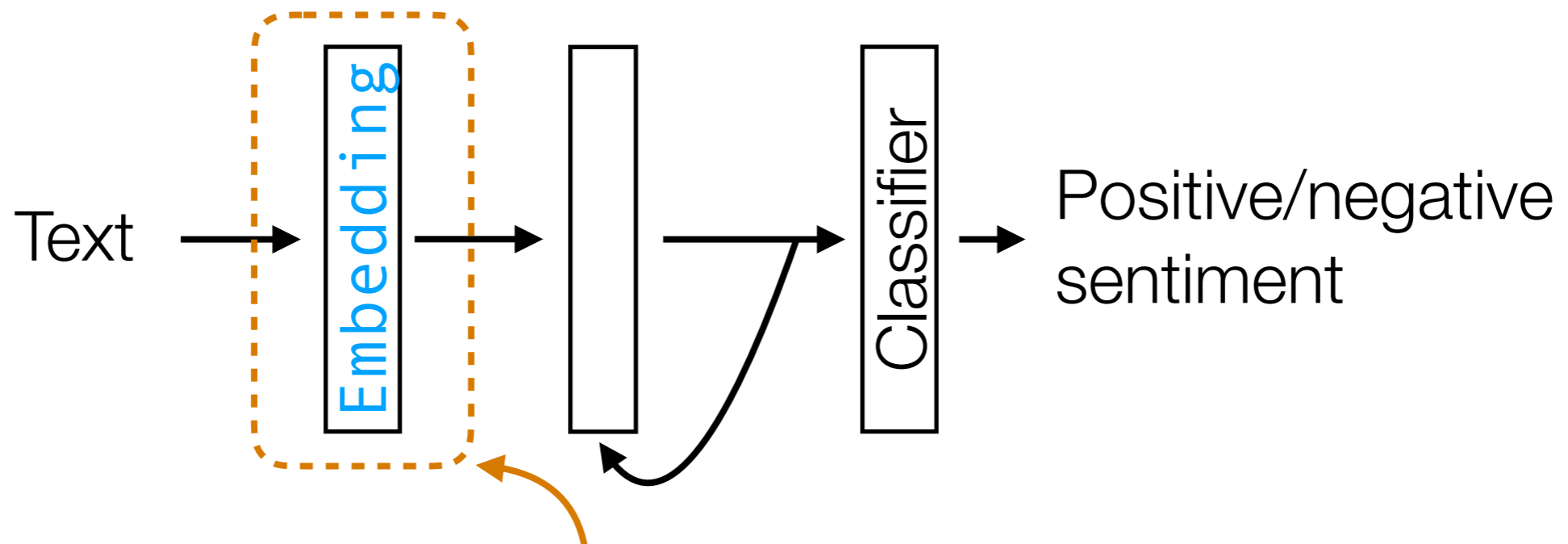
You collect photos from the internet of both, but your dataset size is small, on the order of 1000 images

Strategy: take existing pre-trained CNN for ImageNet classification and change final layer to do classification between Tesla's and Toyota's rather than classifying into 1000 objects

Dealing with Small Datasets

Fine tuning: if there's an existing pre-trained neural net, you could modify it for your problem that has a small dataset

Example: sentiment analysis RNN demo



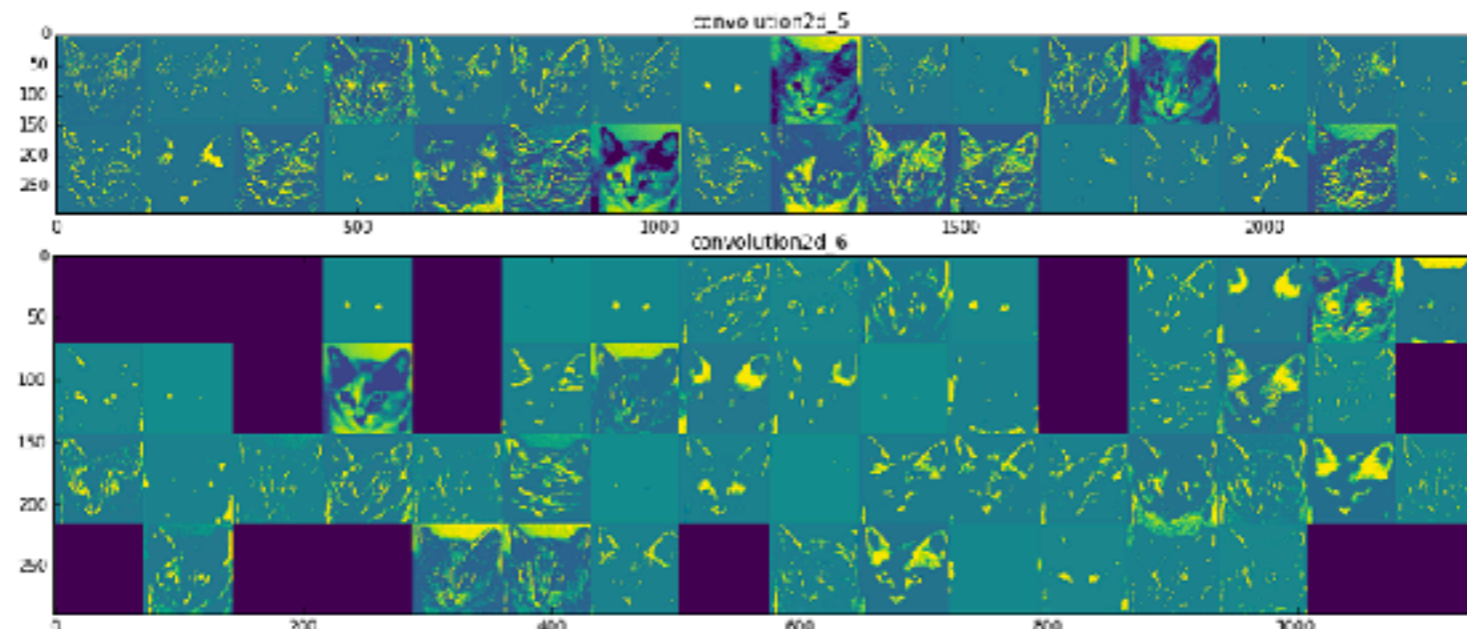
We fixed the weights here to come from GloVe and disabled training for this layer!

GloVe vectors pre-trained on massive dataset (Wikipedia + Gigaword)

IMDb review dataset is small in comparison

Visualizing What a Deep Net Learned

- Very straight-forward for CNNs
 - Plot filter outputs at different layers




- Plot regions that maximally activate an output neuron



Self-Supervised Learning

Even without labels, we can set up a prediction task!

Example: word embeddings like word2vec, GloVe


The opioid epidemic or opioid crisis is the rapid increase in the use of prescription and non-prescription opioid drugs in the United States and Canada in the 2010s.

Predict context of each word!


Training data point: epidemic

“Training label”: the, opioid, or, opioid

Self-Supervised Learning

Even without labels, we can set up a prediction task!

Example: word embeddings like word2vec, GloVe

The opioid epidemic or opioid crisis is the rapid increase in the use of prescription and non-prescription opioid drugs in the United States and Canada in the 2010s.

Predict context of each word!

Training data point: or


“Training label”: opioid, epidemic, opioid, crisis

Self-Supervised Learning

Even without labels, we can set up a prediction task!

Example: word embeddings like word2vec, GloVe

The opioid epidemic or opioid crisis is the rapid increase in the use of prescription and non-prescription opioid drugs in the United States and Canada in the 2010s.


A diagram illustrating word embeddings. The words 'epidemic', 'or', and 'crisis' are enclosed in dashed green boxes. Above these boxes, two magenta curved arrows point from 'epidemic' to 'or' and from 'or' to 'crisis', indicating a relationship between the words.

Predict context of each word!

Training data point: opioid

“Training label”: epidemic, or, crisis, is

There are “positive” examples of what context words are for “opioid”

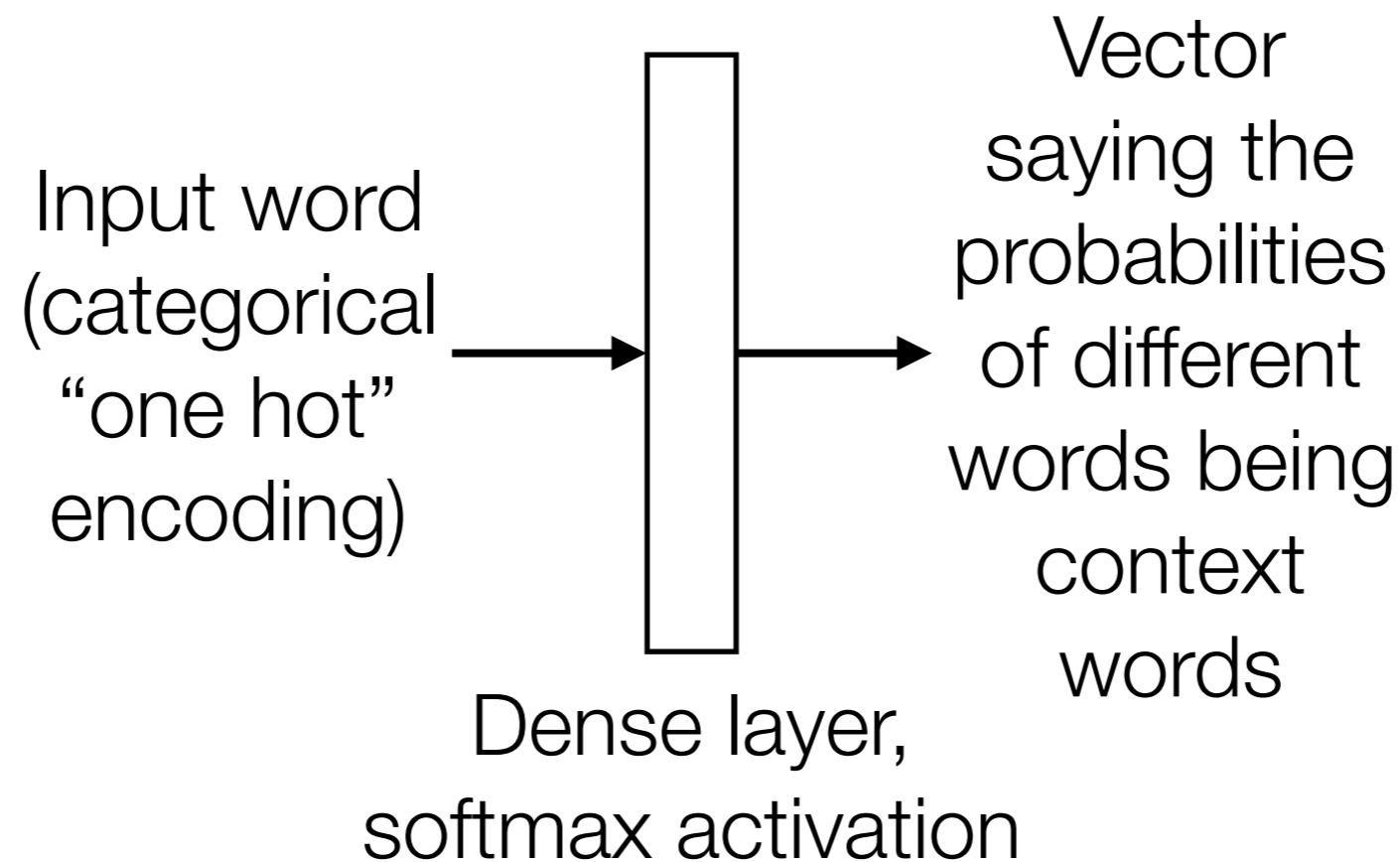
A magenta arrow points from the text 'There are “positive” examples of what context words are for “opioid”' to the training label 'epidemic, or, crisis, is'.

Also provide “negative” examples of words that are *not* likely to be context words (e.g., randomly sample words elsewhere in document)

Self-Supervised Learning

Even without labels, we can set up a prediction task!

Example: word embeddings like word2vec, GloVe



This actually relates to PMI!

Weight matrix: (# words in vocab) by (# neurons)

Dictionary word i has "word embedding" given by row i of weight matrix

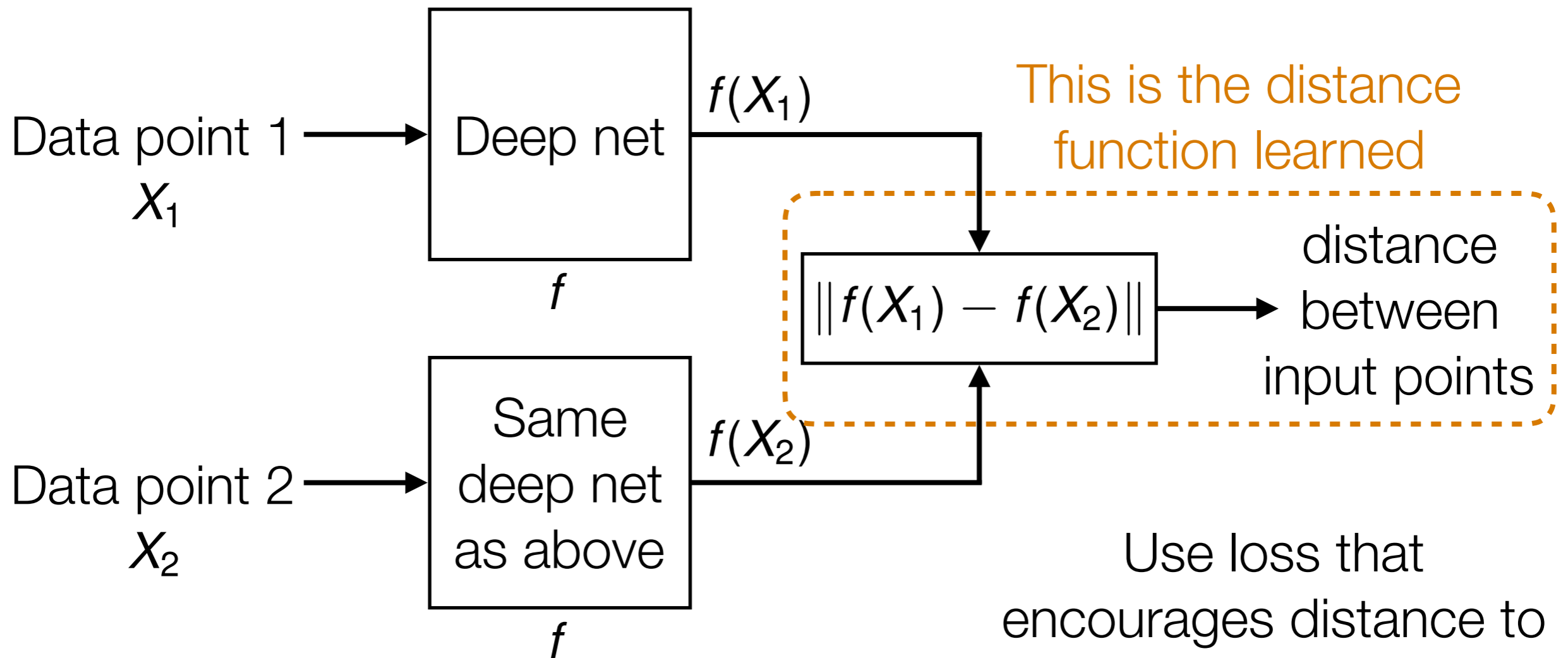
Self-Supervised Learning

Even without labels, we can set up a prediction task!

- Key idea: predict part of the training data from other parts of the training data
- No actual training labels required — we are defining what the training labels are just using the unlabeled training data
- This is an *unsupervised* method that sets up a *supervised prediction* task

Learning Distances with Siamese Nets

Using labeled data, we can learn a distance function



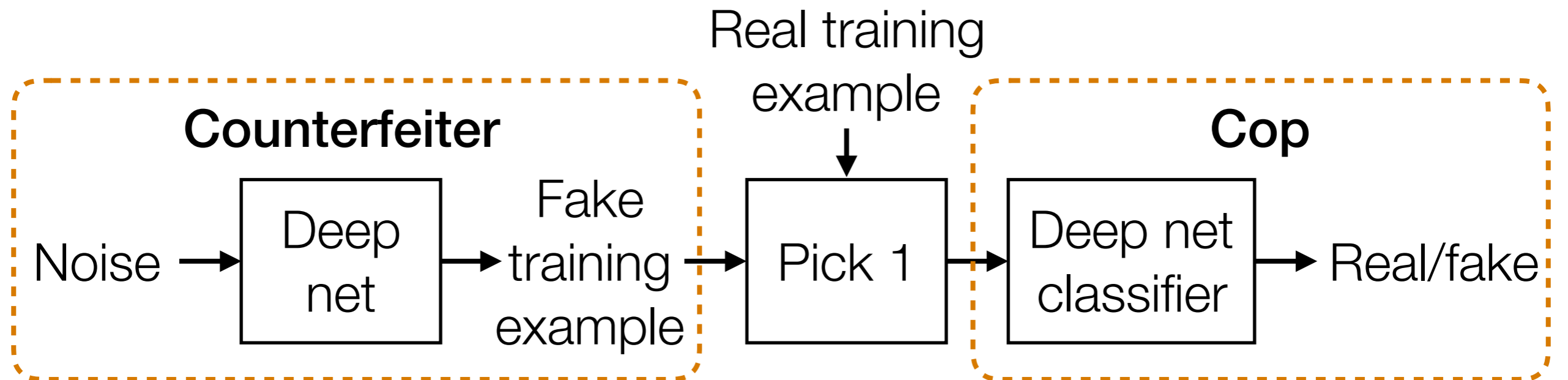
Use loss that encourages distance to be small for data points with same label and large otherwise

Note: we are learning the function f

Generate Fake Data that Look Real

Unsupervised approach: generate data that look like training data

Example: Generative Adversarial Network (GAN)



Counterfeiter tries to get better at tricking the cop

Cop tries to get better at telling which examples are real vs fake

Terminology: counterfeiter is the **generator**, cop is the **discriminator**

Other approaches: variational autoencoders, pixelRNNs/pixelCNNs

Generate Fake Data that Look Real



Fake celebrities generated by NVIDIA using GANs
(Karras et al Oct 27, 2017)

Google DeepMind's WaveNet makes fake audio that sounds like
whoever you want using pixelRNNs (Oord et al 2016)

Generate Fake Data that Look Real

Monet ↔ Photos



Monet → photo

Zebras ↔ Horses



zebra → horse

Summer ↔ Winter



summer → winter



photo → Monet



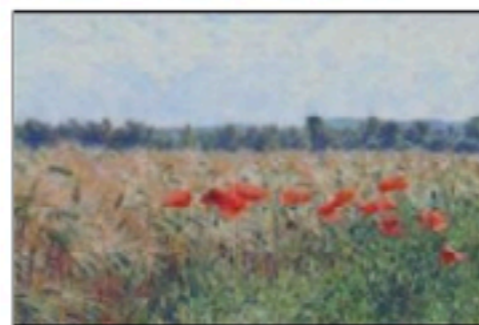
horse → zebra



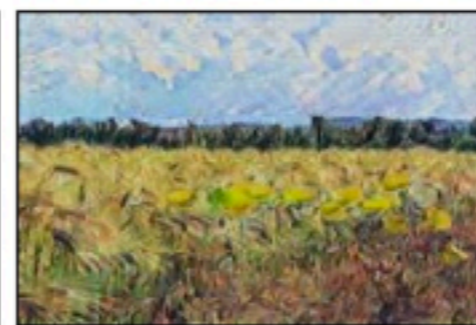
winter → summer



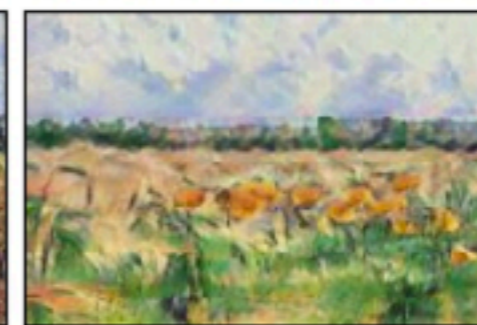
Photograph



Monet



Van Gogh



Cezanne

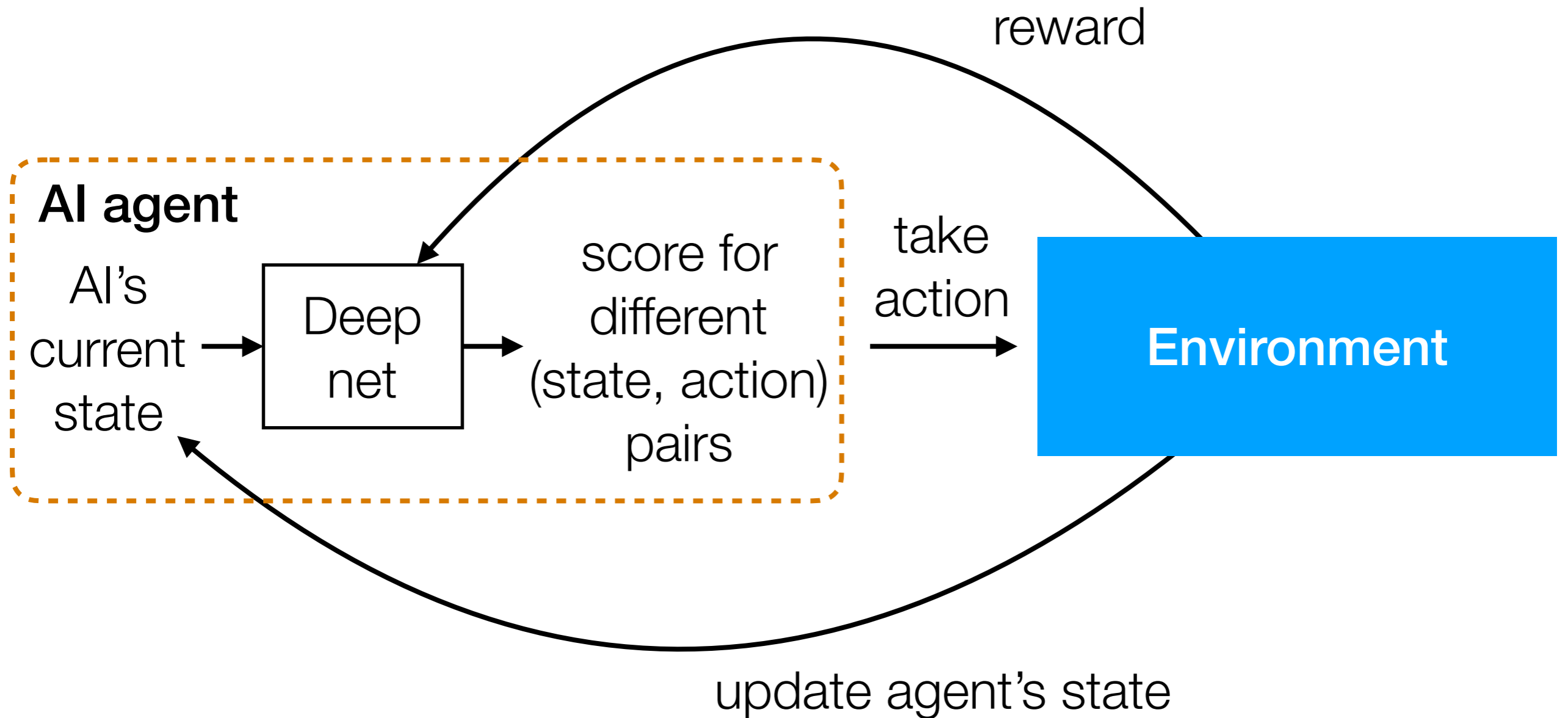


Ukiyo-e

Image-to-image translation results from UC Berkeley using GANs
(Isola et al 2017, Zhu et al 2017)

Deep Reinforcement Learning

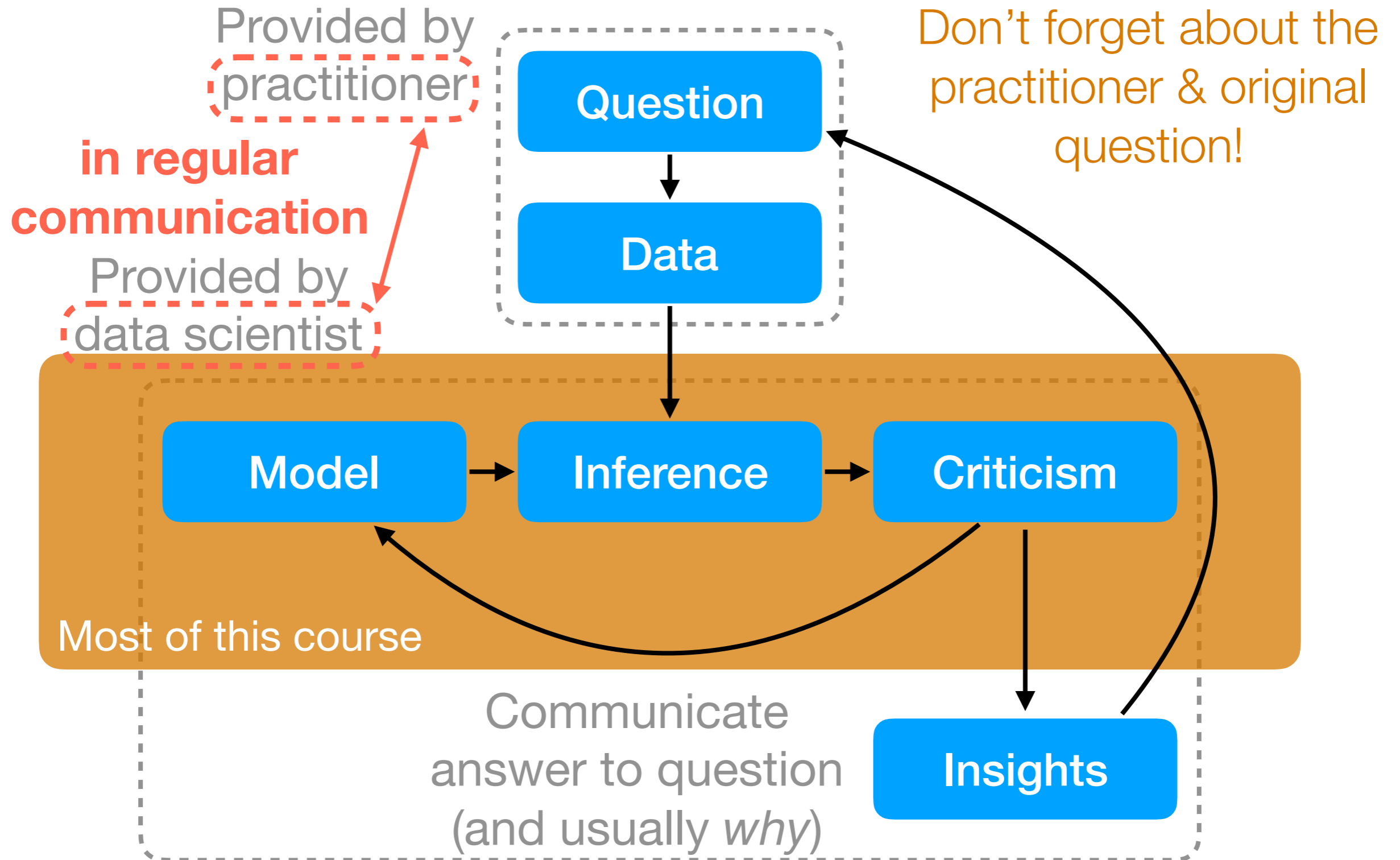
The machinery behind AlphaGo and similar systems



The Future of Deep Learning

- Deep learning currently is still limited in what it can do — the layers do simple operations and have to be differentiable
 - How do we make deep nets that generalize better?
- Still lots of engineering and expert knowledge used to design some of the best systems (e.g., AlphaGo)
 - How do we get away with using less expert knowledge?
- How do we do lifelong learning?

95-865



95-865 Some Parting Thoughts

- Remember to **visualize different steps of your data analysis pipeline** — very helpful when you're still debugging
- Very often there are *tons* of models that you could try
 - Come up with **quantitative metrics** that make sense for your problem, and use these metrics to **evaluate models with a prediction task on held-out data**
- Often times you won't have labels!
 - Manually obtain labels (either you do it or crowdsource)
 - Set up self-supervised learning task

Thanks for being a beta tester!