# Advanced Deep Learning

## DATA.ML.230

# Fitting models and regularization

# Fitting models

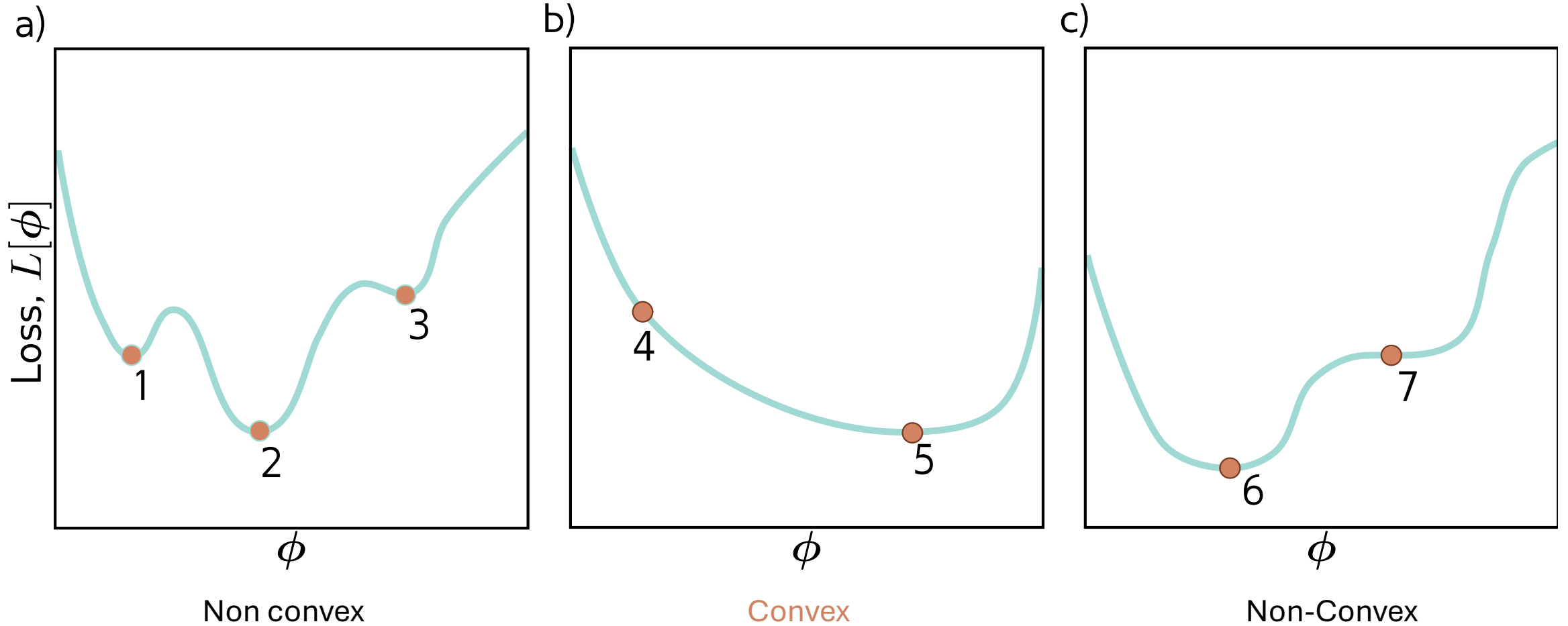**Step 1.** Compute the derivatives of the loss with respect to the parameters:

$$\frac{\partial L}{\partial \phi} = \begin{bmatrix} \frac{\partial L}{\partial \phi_0} \\ \frac{\partial L}{\partial \phi_1} \\ \vdots \\ \frac{\partial L}{\partial \phi_N} \end{bmatrix}.$$

**Step 2.** Update the parameters according to the rule:

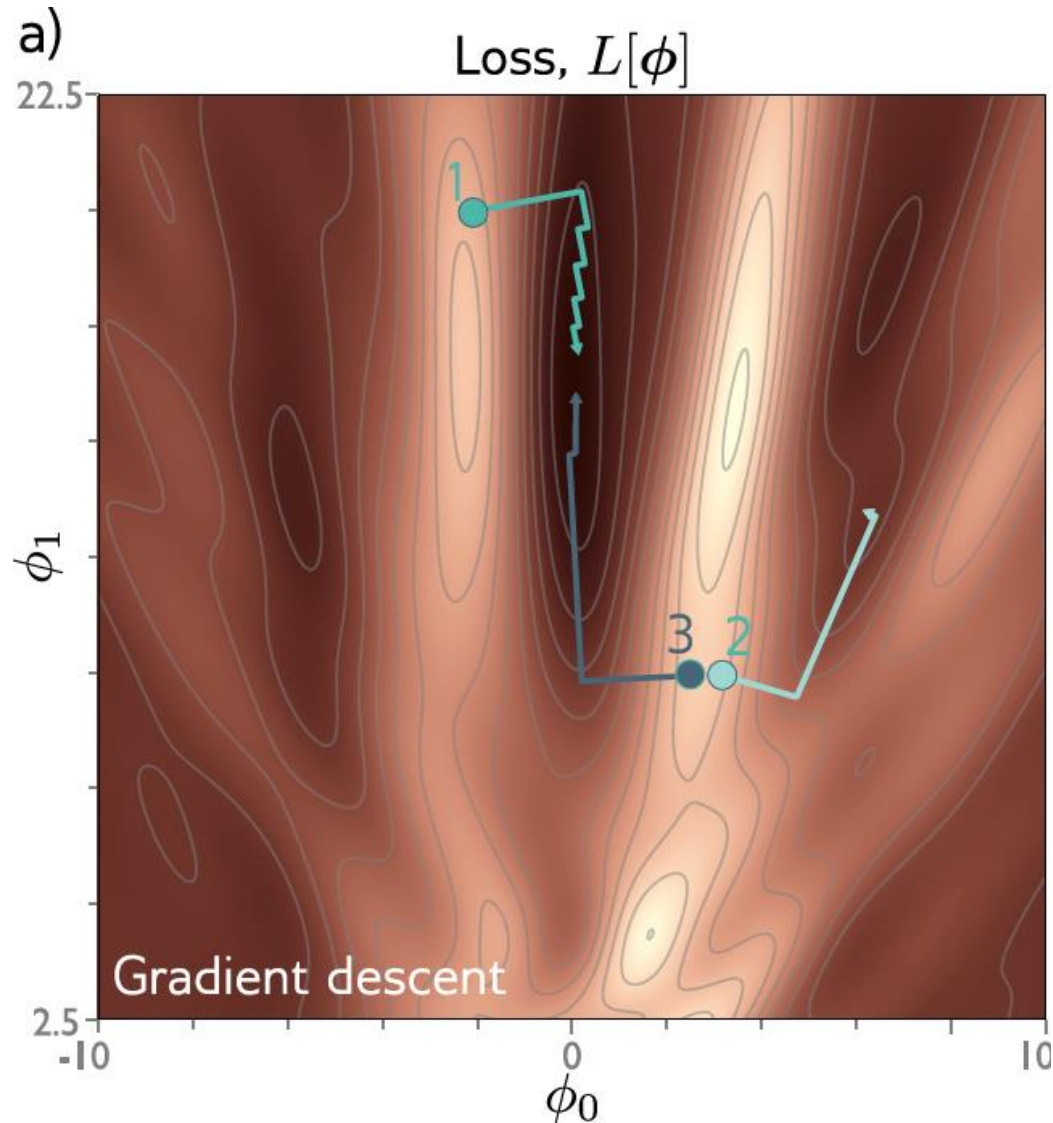$$\phi \longleftarrow \phi - \alpha \frac{\partial L}{\partial \phi},$$

where the positive scalar $\alpha$ determines the magnitude of the change.

# Convexity



Test for convexity is that 2$^{nd}$ derivative is positive everywhere

# Stochastic Gradient Descent



a)

Loss, $L[\phi]$
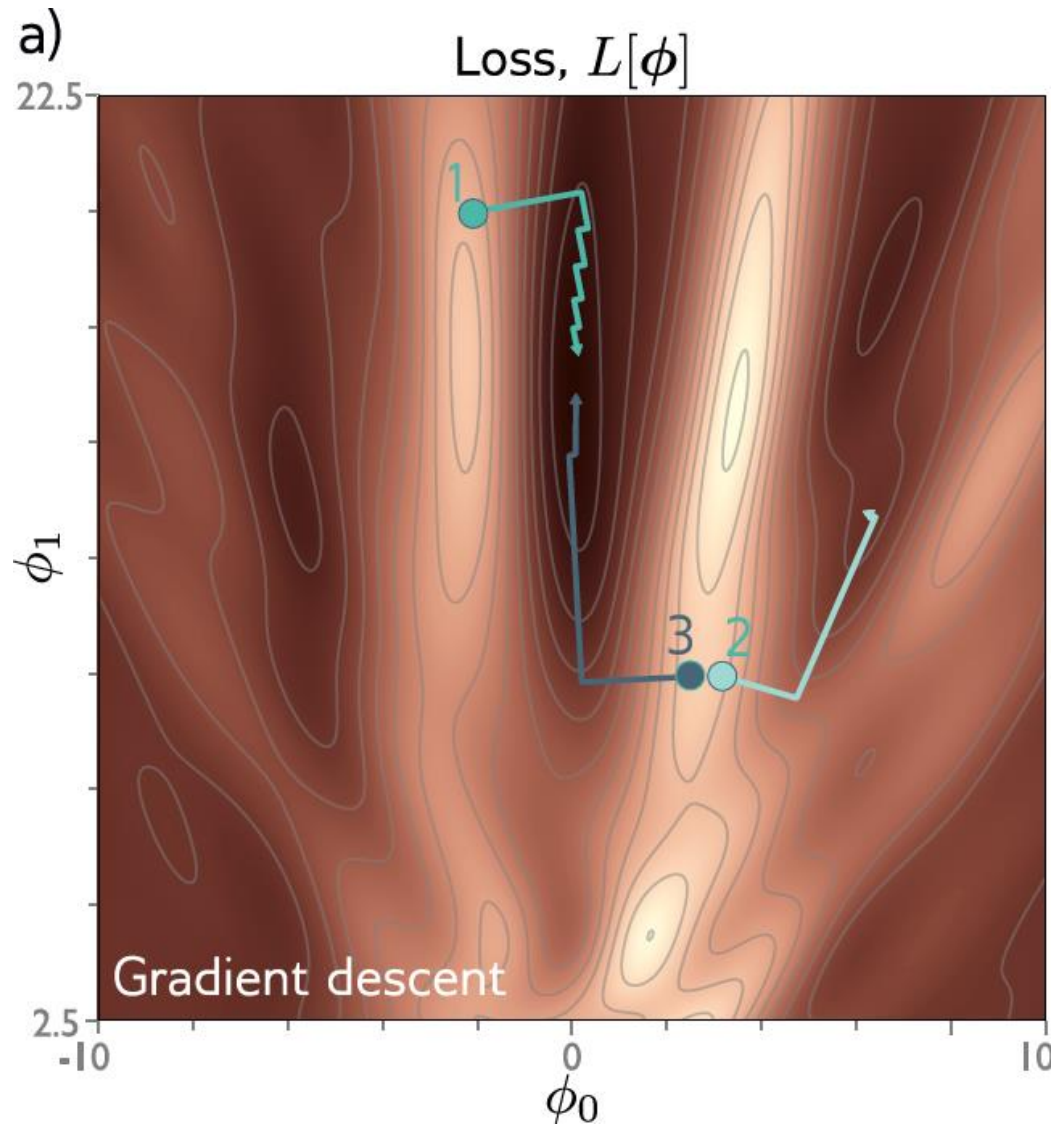
Gradient descent

IDEA: add noise

- Stochastic gradient descent
- Compute gradient based on only a subset of points – a mini-batch
- Work through dataset sampling without replacement
- One pass though the data is called an epoch

# Stochastic Gradient Descent



a) Loss, $L[\phi]$
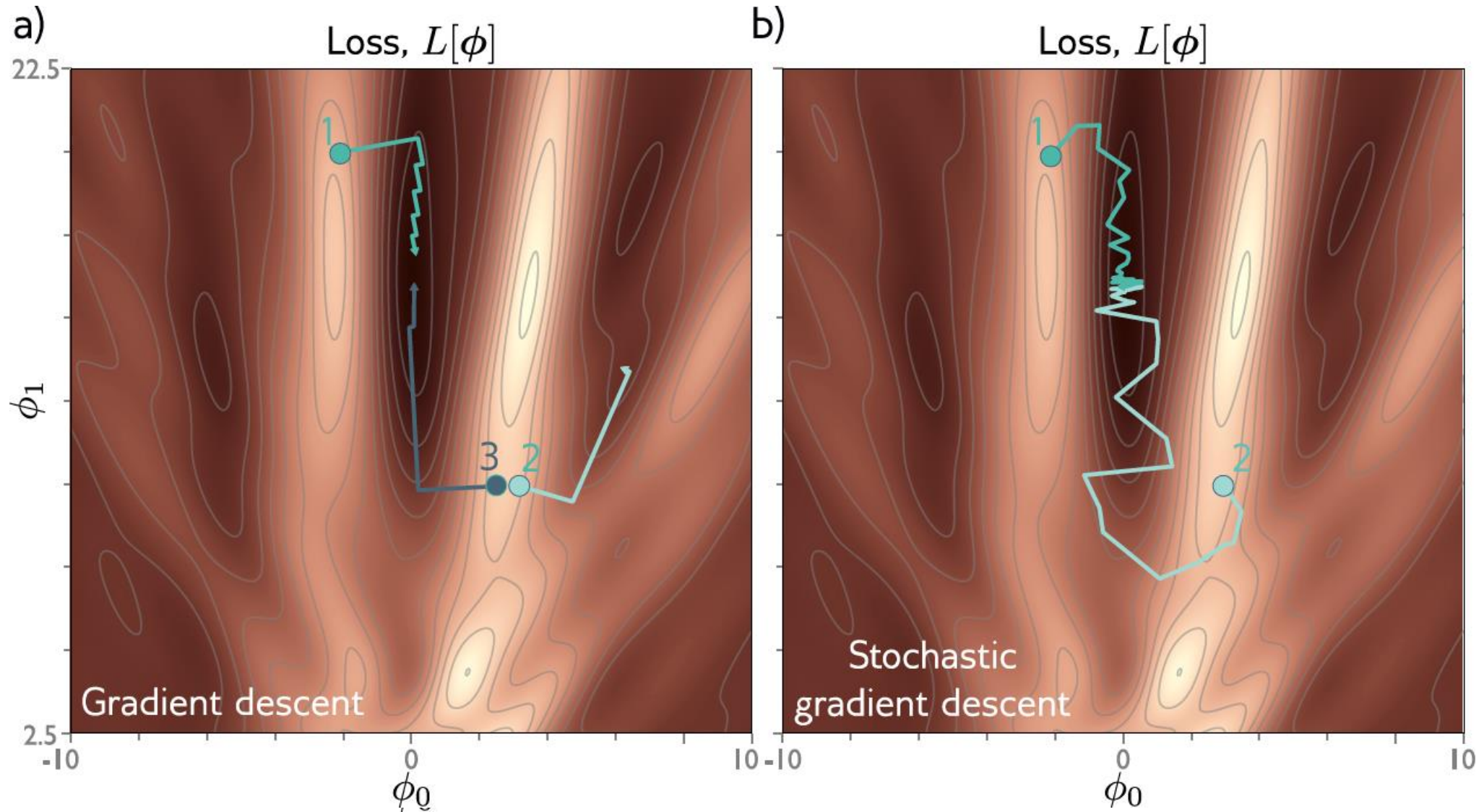
Gradient descent

Before (full batch descent)

$$\phi_{t+1} \longleftarrow \phi_t - \alpha \sum_{i=1}^{I} \frac{\partial \ell_i[\phi_t]}{\partial \phi},$$

After (SGD)

$$\phi_{t+1} \longleftarrow \phi_t - \alpha \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\phi_t]}{\partial \phi},$$

Fixed learning rate α

# Stochastic Gradient Descent

# Properties of SGD

- Can escape from local minima
- Adds noise, but still sensible updates as based on part of data
- Uses all data equally
- Less computationally expensive
- Seems to find better solutions

- Doesn't converge in traditional sense
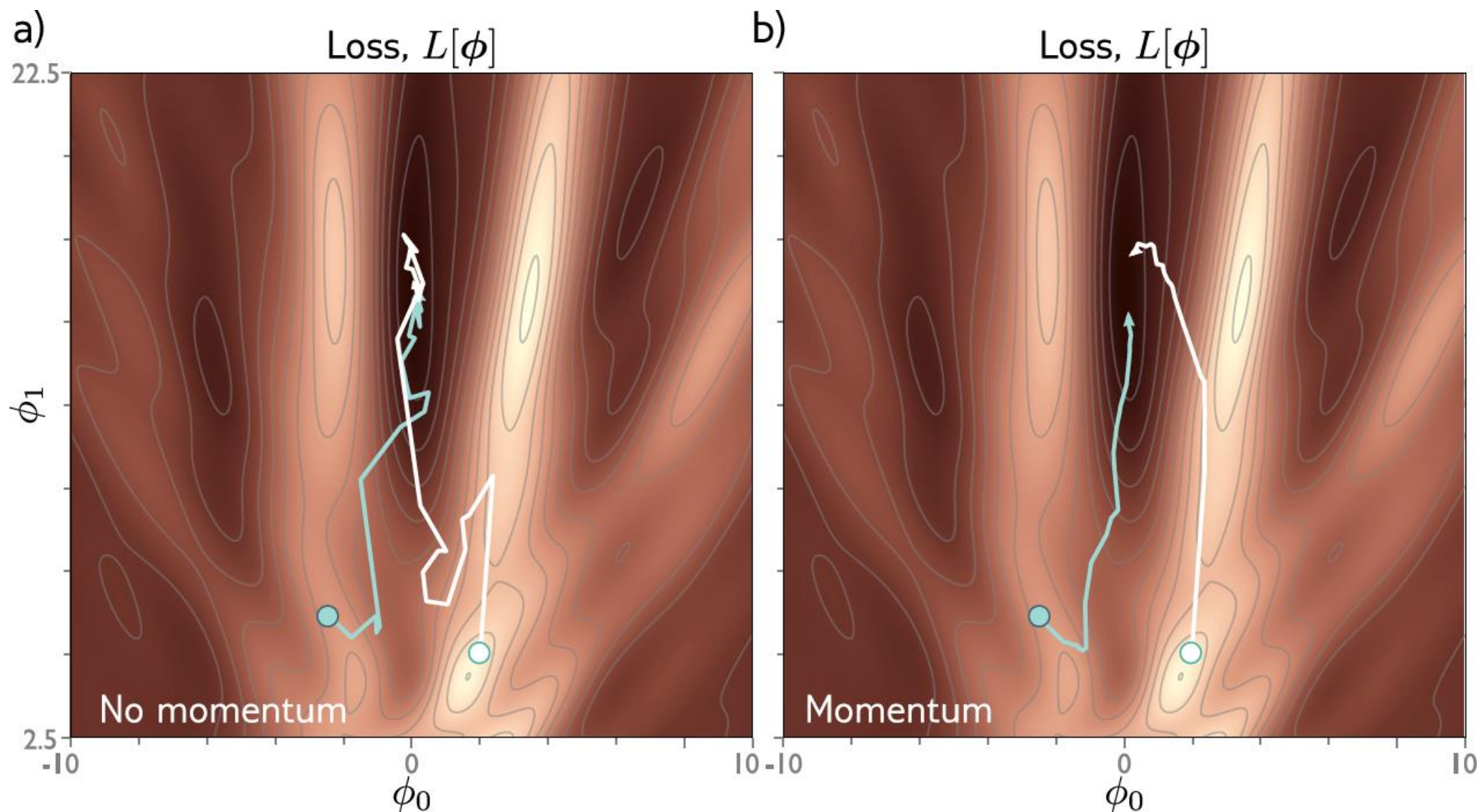- Learning rate schedule – decrease learning rate over time

# Momentum

- Weighted sum of the current gradient and the previous gradient:

$$\mathbf{m}_{t+1} \leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\boldsymbol{\phi}_t]}{\partial \boldsymbol{\phi}}$$

$$\boldsymbol{\phi}_{t+1} \leftarrow \boldsymbol{\phi}_t - \alpha \cdot \mathbf{m}_{t+1}$$

# Momentum



a) Loss, $L[\phi]$ — No momentum

b) Loss, $L[\phi]$ — Momentum

# Nesterov accelerated momentum

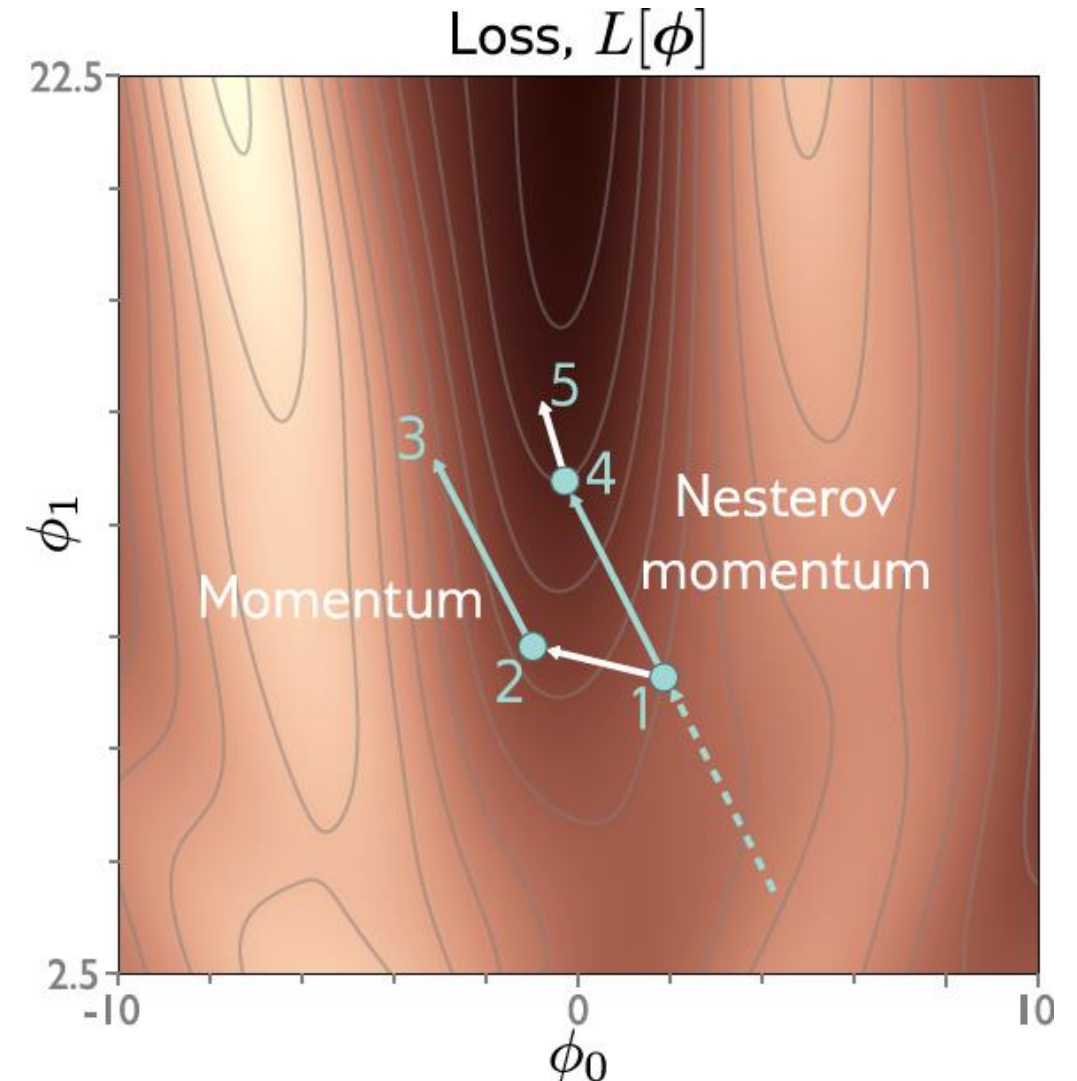- Momentum is kind of like a prediction of where we are going

$$\mathbf{m}_{t+1} \leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\boldsymbol{\phi}_t]}{\partial \boldsymbol{\phi}}$$

$$\boldsymbol{\phi}_{t+1} \leftarrow \boldsymbol{\phi}_t - \alpha \cdot \mathbf{m}_{t+1}$$
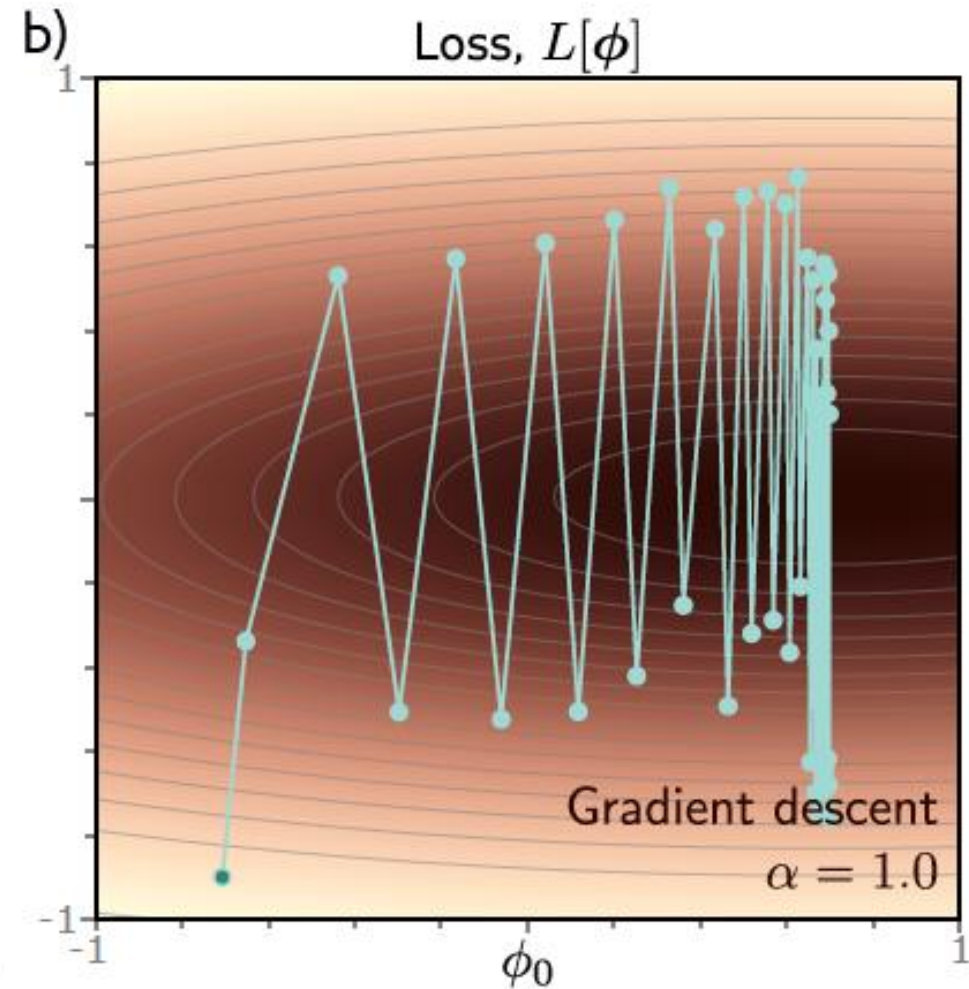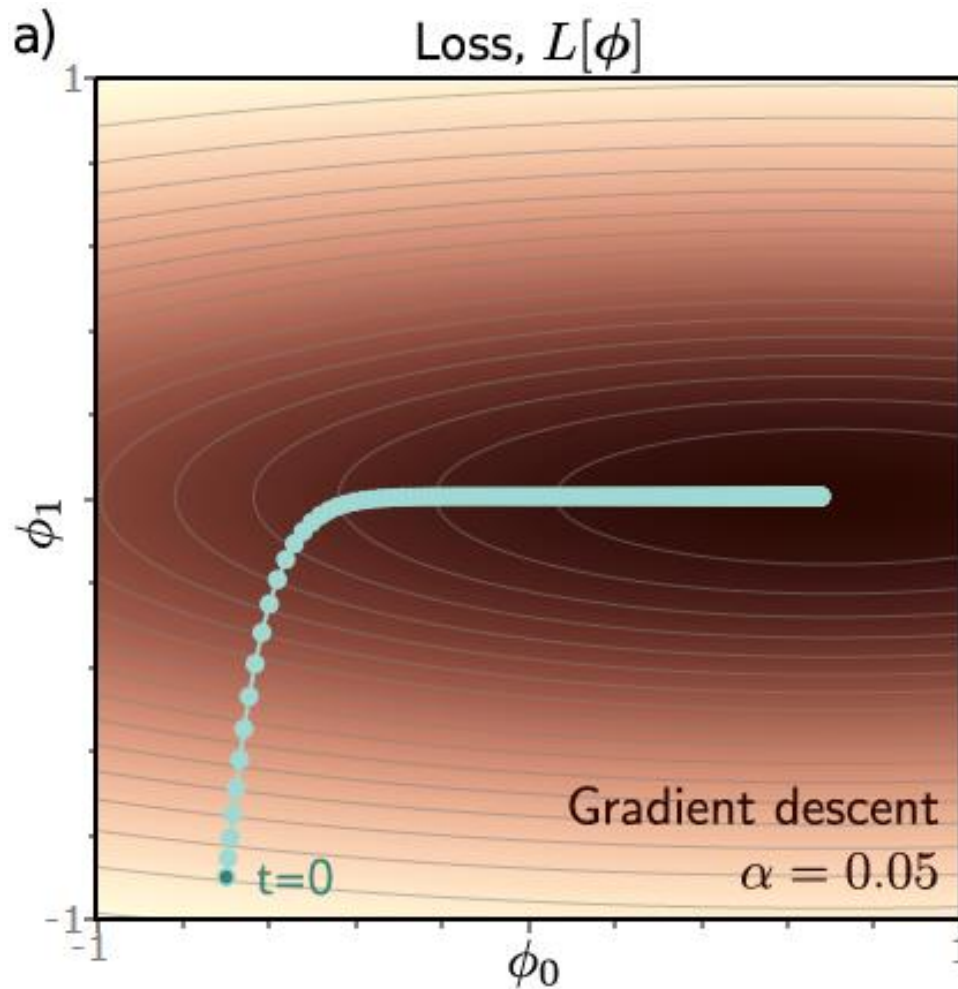
- Move in the predicted direction, THEN, measure the gradient

$$\mathbf{m}_{t+1} \leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta) \sum_{i \in \mathcal{B}_t} \frac{\partial \ell_i[\boldsymbol{\phi}_t - \alpha \cdot \mathbf{m}_t]}{\partial \boldsymbol{\phi}}$$

$$\boldsymbol{\phi}_{t+1} \leftarrow \boldsymbol{\phi}_t - \alpha \cdot \mathbf{m}_{t+1}$$



Loss, $L[\boldsymbol{\phi}]$

# Adaptive moment estimation (Adam)

# Normalized gradients

- Measure mean and pointwise squared gradient

$$\mathbf{m}_{t+1} \leftarrow \frac{\partial L[\phi_t]}{\partial \phi}$$

$$\mathbf{v}_{t+1} \leftarrow \frac{\partial L[\phi_t]^2}{\partial \phi}$$

- Normalize:

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}$$

# Normalized gradients

- Measure mean and pointwise squared gradient

$$\mathbf{m}_{t+1} \leftarrow \frac{\partial L[\phi_t]}{\partial \phi}$$

$$\mathbf{v}_{t+1} \leftarrow \frac{\partial L[\phi_t]^2}{\partial \phi}$$

$$\mathbf{m}_{t+1} = \begin{bmatrix} 3.0 \\ -2.0 \\ 5.0 \end{bmatrix}$$

$$\mathbf{v}_{t+1} = \begin{bmatrix} 9.0 \\ 4.0 \\ 25.0 \end{bmatrix}$$

- Normalize:

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}$$

$$\frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon} = \begin{bmatrix} 1.0 \\ -1.0 \\ 1.0 \end{bmatrix}$$
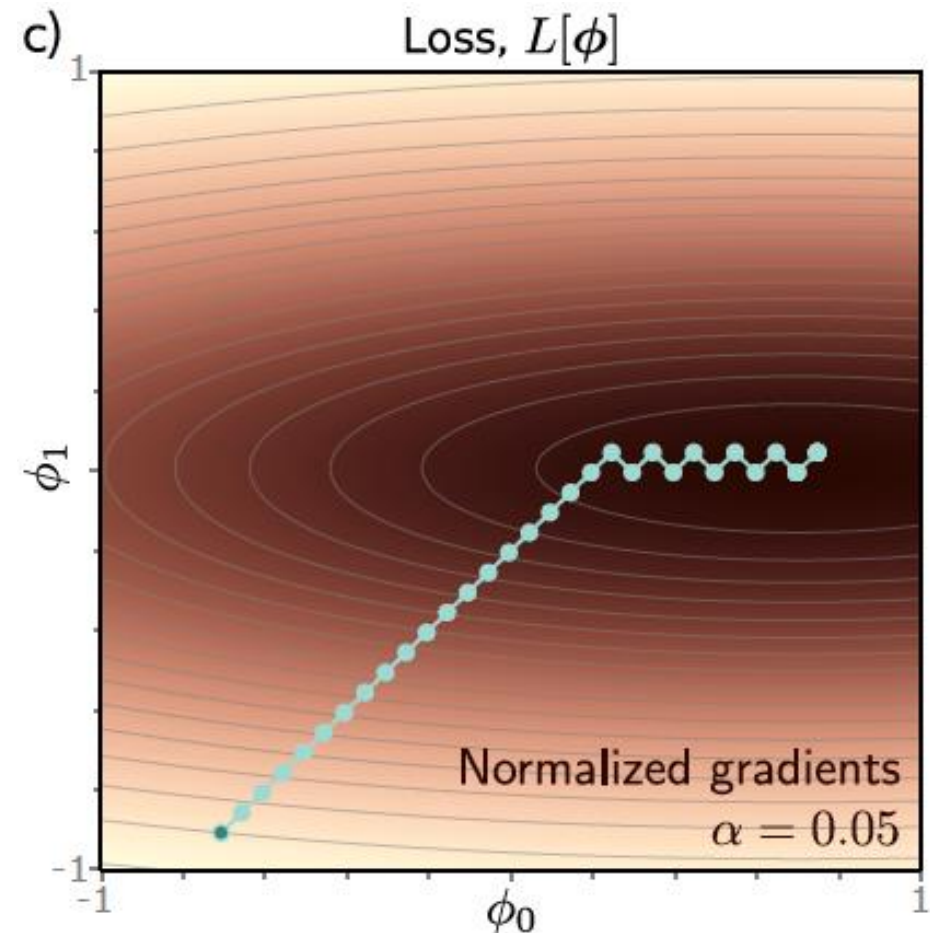
# Normalized gradients

- Measure mean and pointwise squared gradient

$$\mathbf{m}_{t+1} \leftarrow \frac{\partial L[\phi_t]}{\partial \phi}$$

$$\mathbf{v}_{t+1} \leftarrow \frac{\partial L[\phi_t]^2}{\partial \phi}$$

- Normalize:

$$\phi_{t+1} \leftarrow \phi_t - \alpha \cdot \frac{\mathbf{m}_{t+1}}{\sqrt{\mathbf{v}_{t+1}} + \epsilon}$$



c) Loss, $L[\phi]$

Normalized gradients
$\alpha = 0.05$

# Adaptive moment estimation (Adam)

- Compute mean and pointwise squared gradients with momentum

$$\mathbf{m}_{t+1} \leftarrow \beta \cdot \mathbf{m}_t + (1 - \beta)\frac{\partial L[\boldsymbol{\phi}_t]}{\partial \boldsymbol{\phi}}$$

$$\mathbf{v}_{t+1} \leftarrow \gamma \cdot \mathbf{v}_t + (1 - \gamma)\left(\frac{\partial L[\boldsymbol{\phi}_t]}{\partial \boldsymbol{\phi}}\right)^2$$
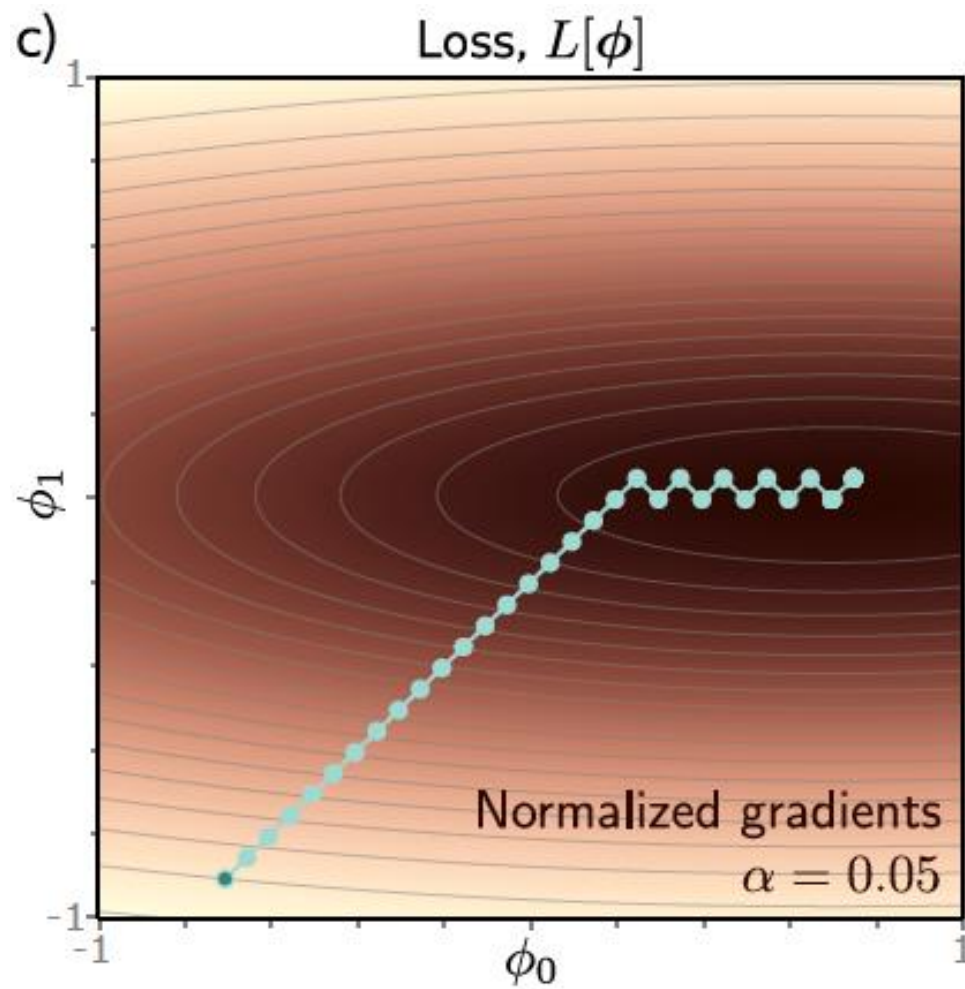
- Moderate near start of the sequence

$$\tilde{\mathbf{m}}_{t+1} \leftarrow \frac{\mathbf{m}_{t+1}}{1 - \beta^{t+1}}$$

$$\tilde{\mathbf{v}}_{t+1} \leftarrow \frac{\mathbf{v}_{t+1}}{1 - \gamma^{t+1}}$$

- Update the parameters

$$\boldsymbol{\phi}_{t+1} \leftarrow \boldsymbol{\phi}_t - \alpha \cdot \frac{\tilde{\mathbf{m}}_{t+1}}{\sqrt{\tilde{\mathbf{v}}_{t+1}} + \epsilon}$$

# Adaptive moment estimation (Adam)



c) Loss, $L[\phi]$

Normalized gradients
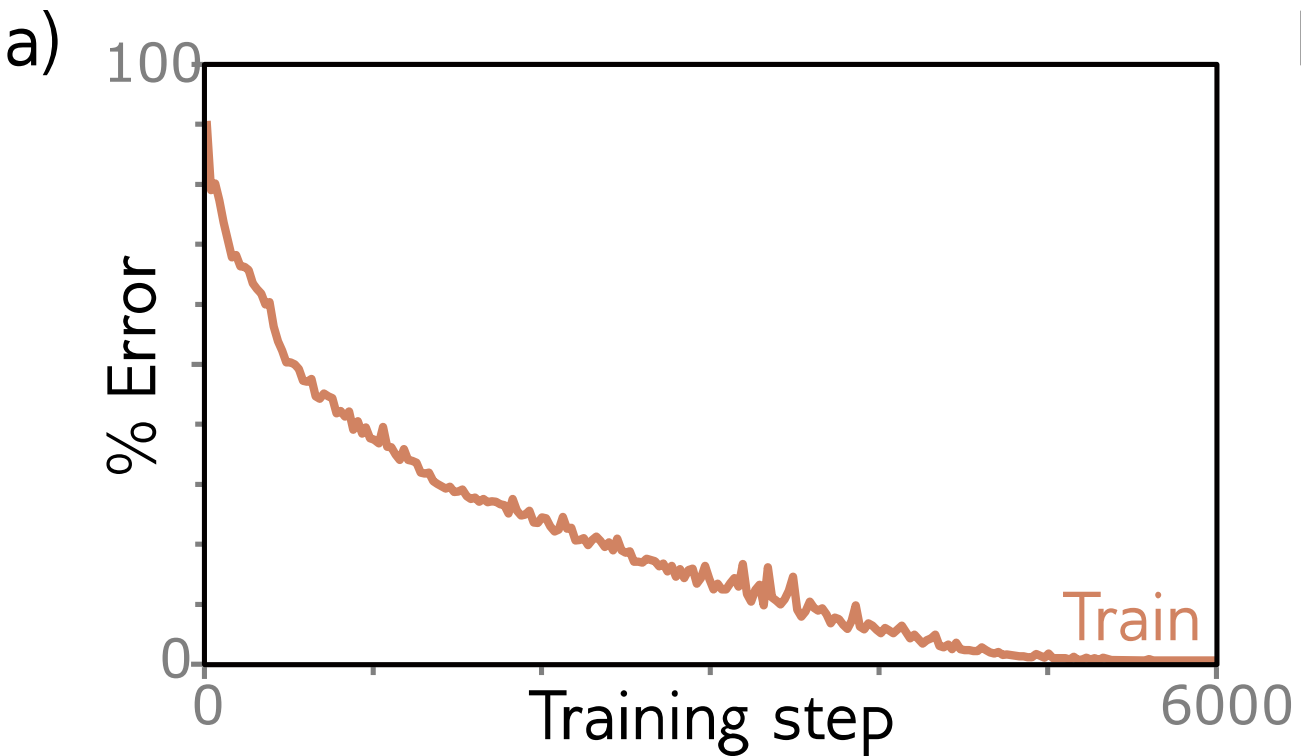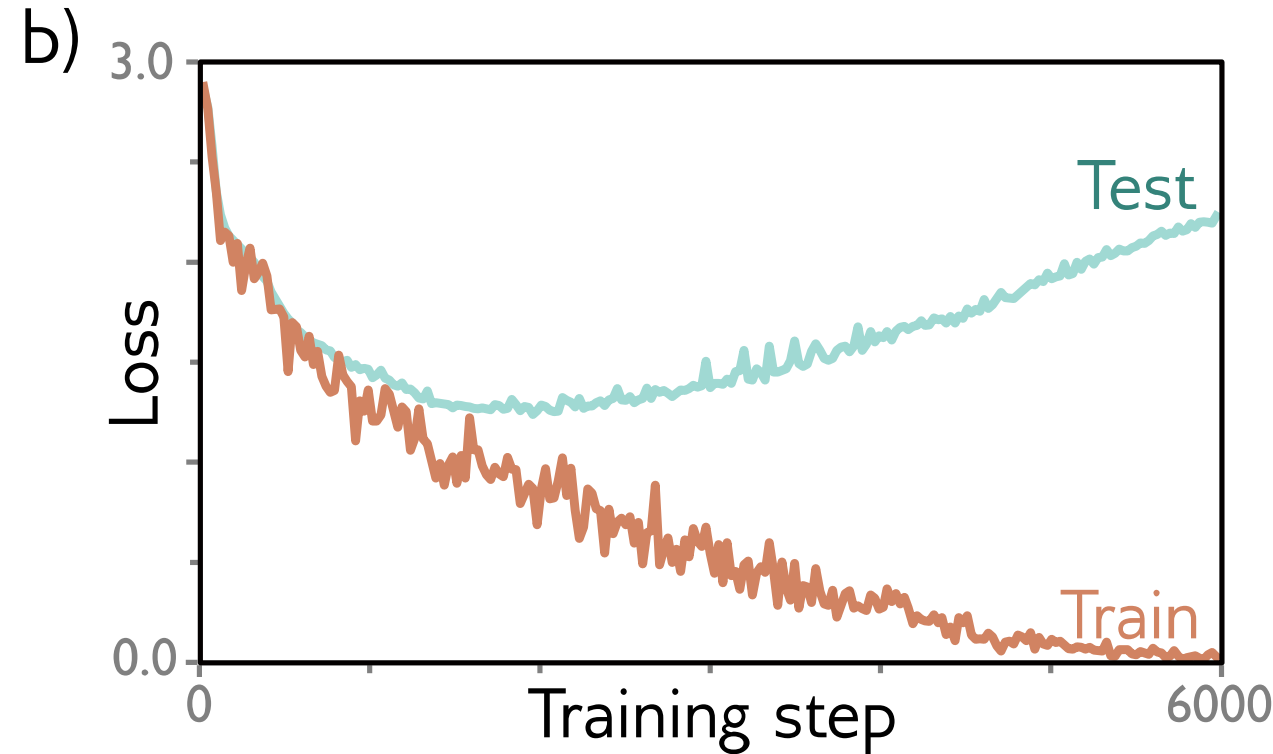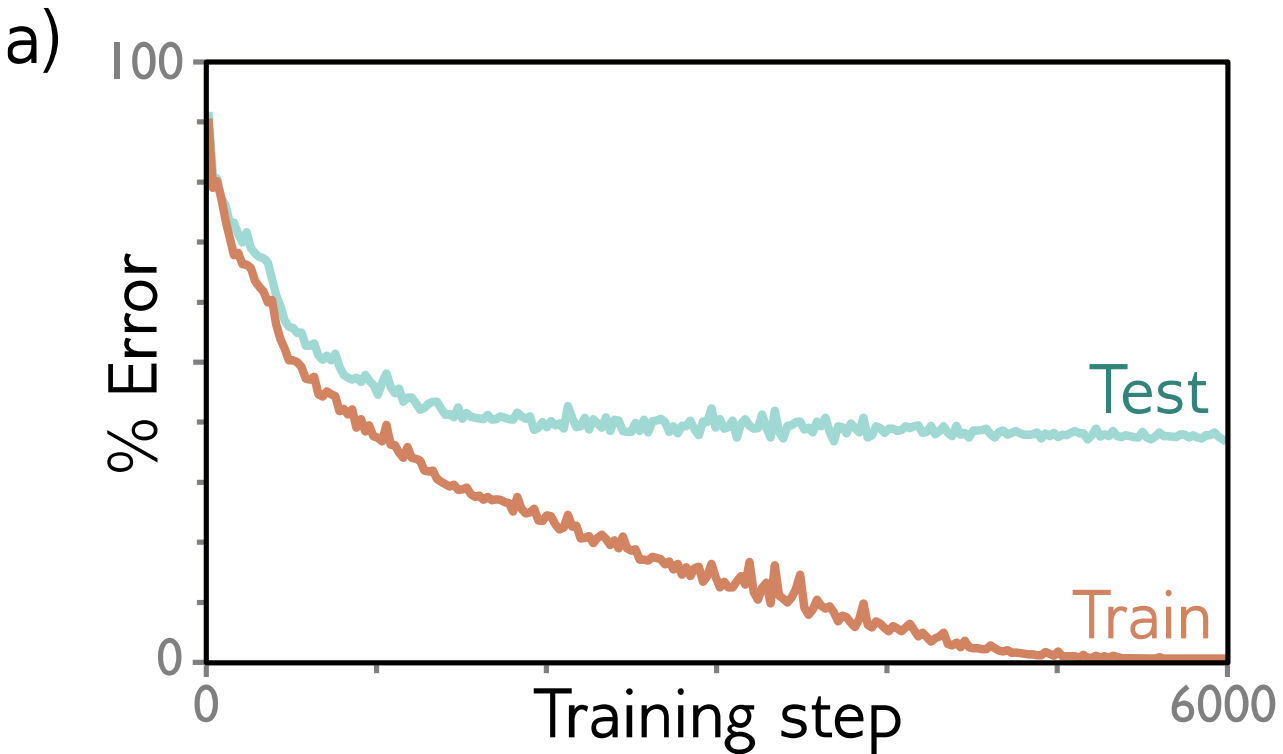$\alpha = 0.05$

d) Loss, $L[\phi]$

Adam
$\alpha = 0.05, \beta = 0.9, \gamma = 0.99$

# Hyperparameters

- Choice of learning algorithm
- Learning rate
- Momentum
- Neural network structure

# Measuring performance



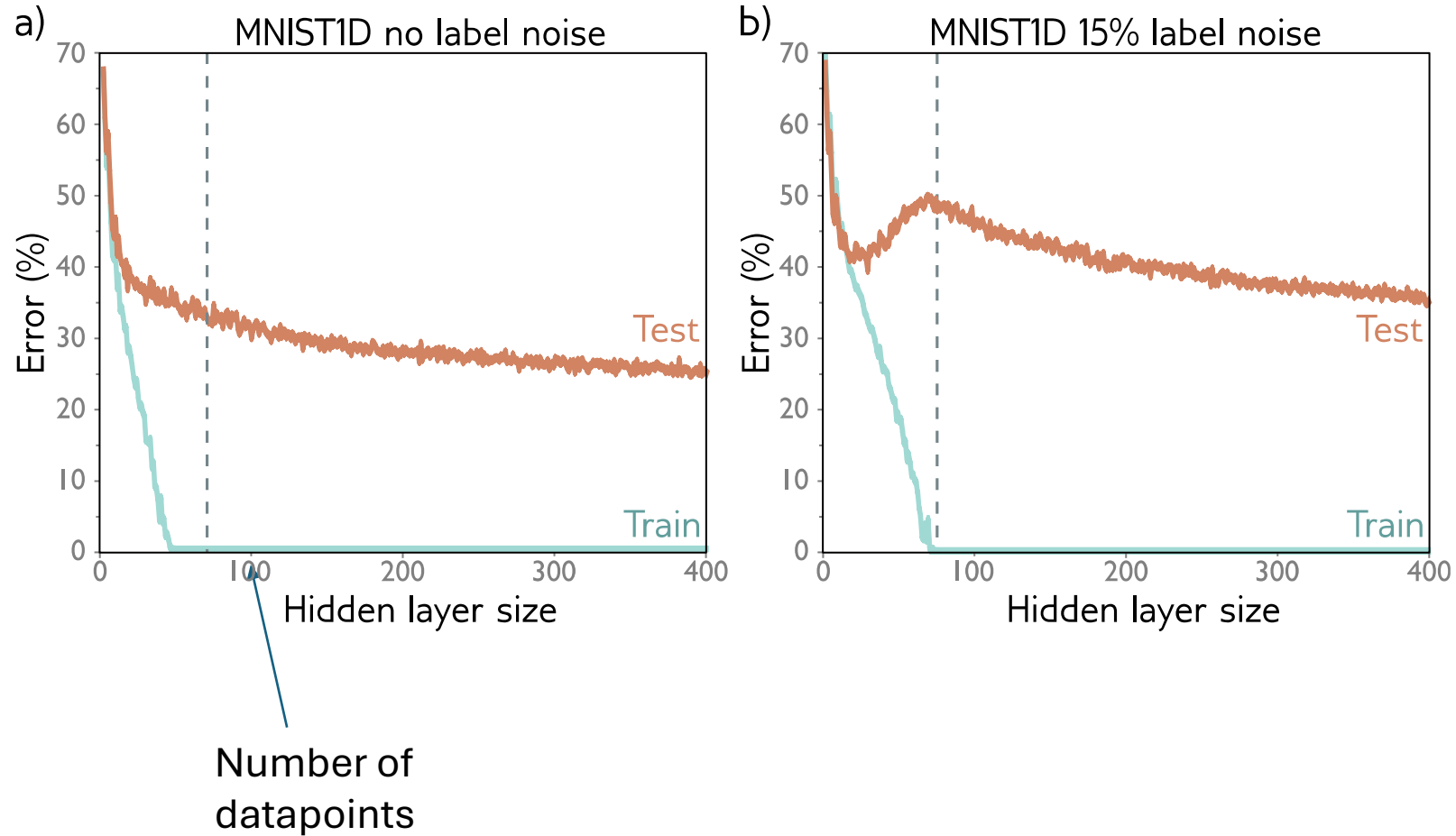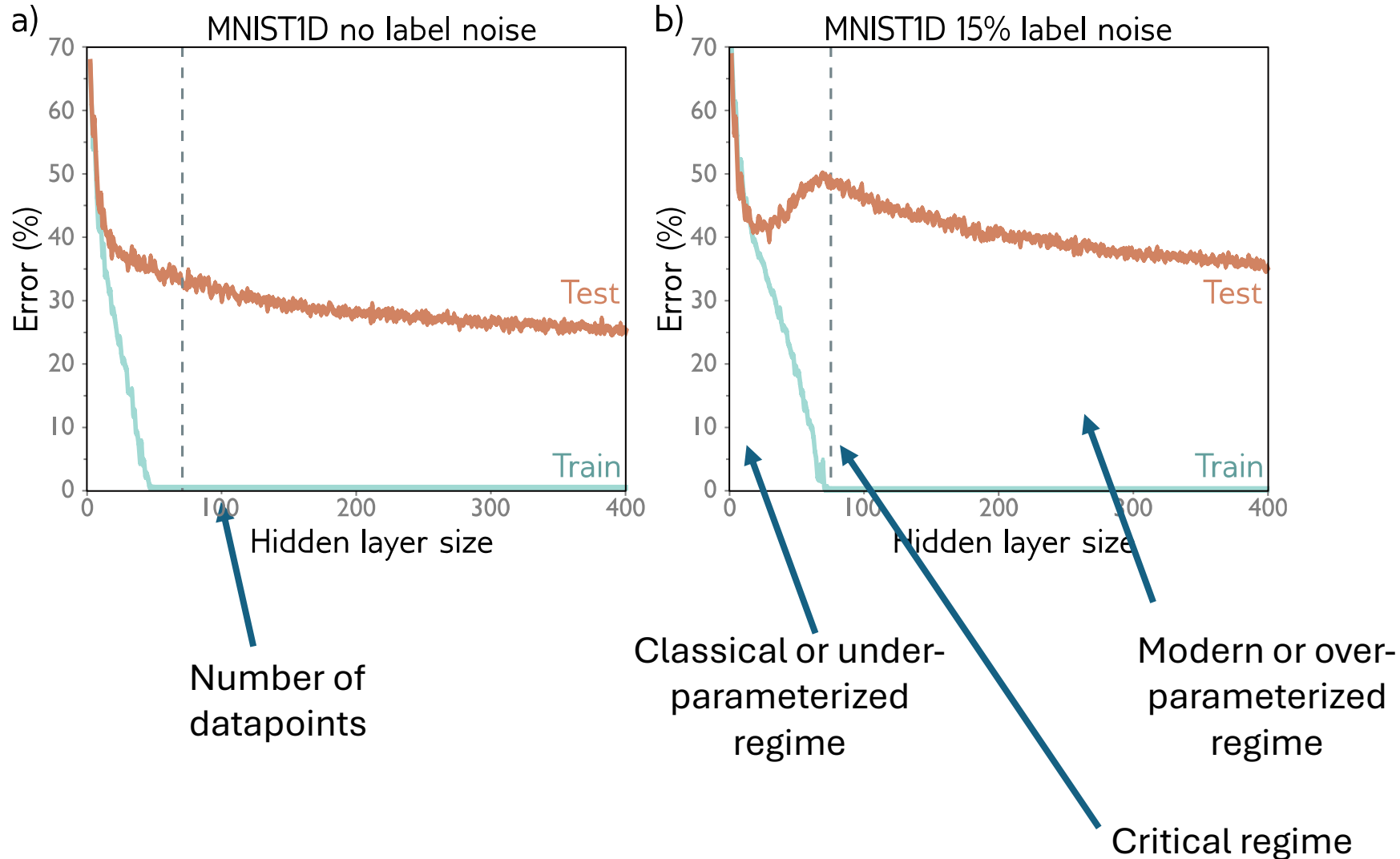a) A plot with y-axis labeled "% Error" ranging from 0 to 100 and x-axis labeled "Training step" ranging from 0 to 6000, showing a decreasing orange curve labeled "Train".

b) A plot with y-axis labeled "Loss" ranging from 0.0 to 3.0 and x-axis labeled "Training step" ranging from 0 to 6000, showing a decreasing orange curve labeled "Train".

# Measuring performance



Need to use separate test data

# Double descent



a) MNIST1D no label noise
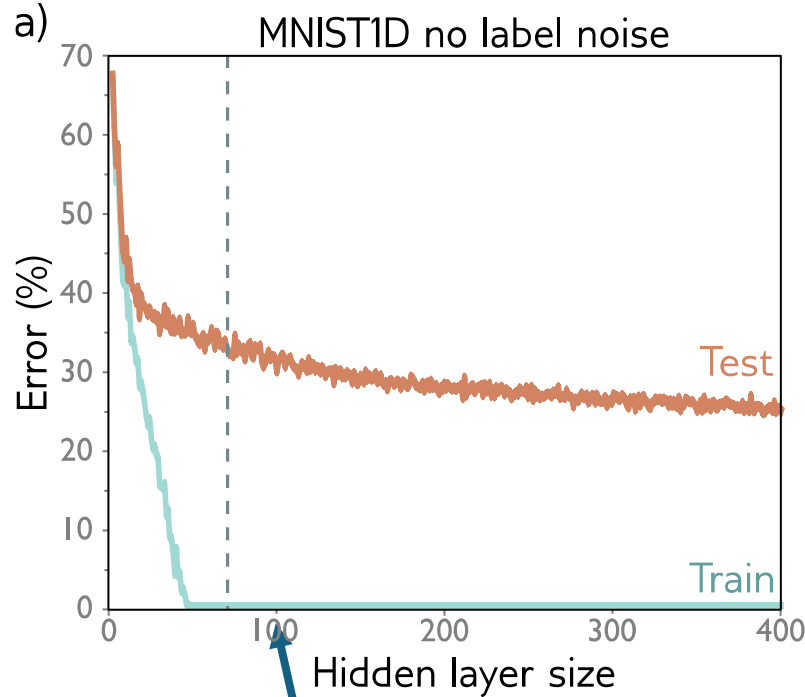
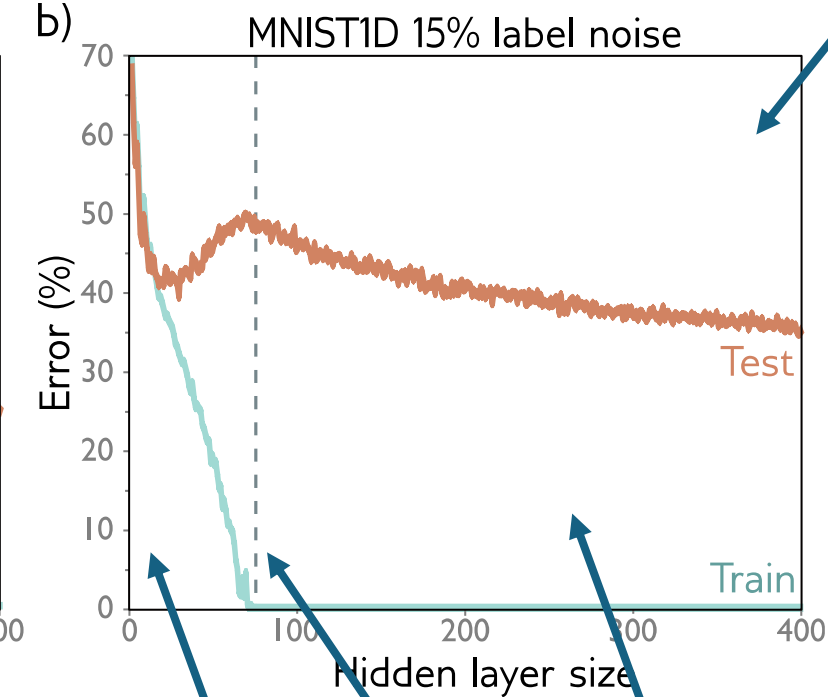b) MNIST1D 15% label noise

# Double descent

# Double descent

- Note that errors for the train data is very close to zero.
- Whatever is happening isn't happening at training data points
- Must be happening between the data points??



a) MNIST1D no label noise

b) MNIST1D 15% label noise

Number of datapoints

Classical or under-parameterized regime
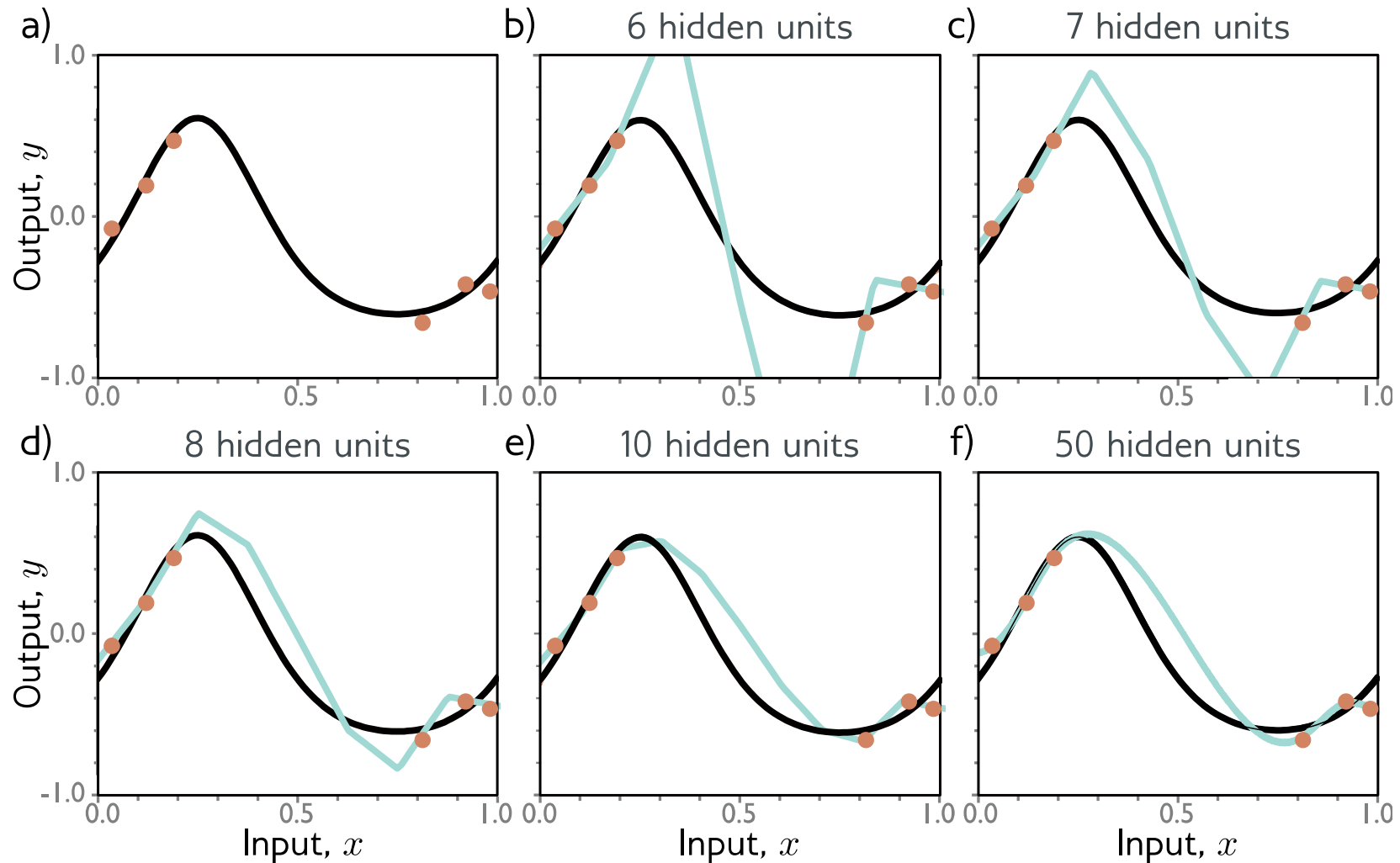
Modern or over-parameterized regime

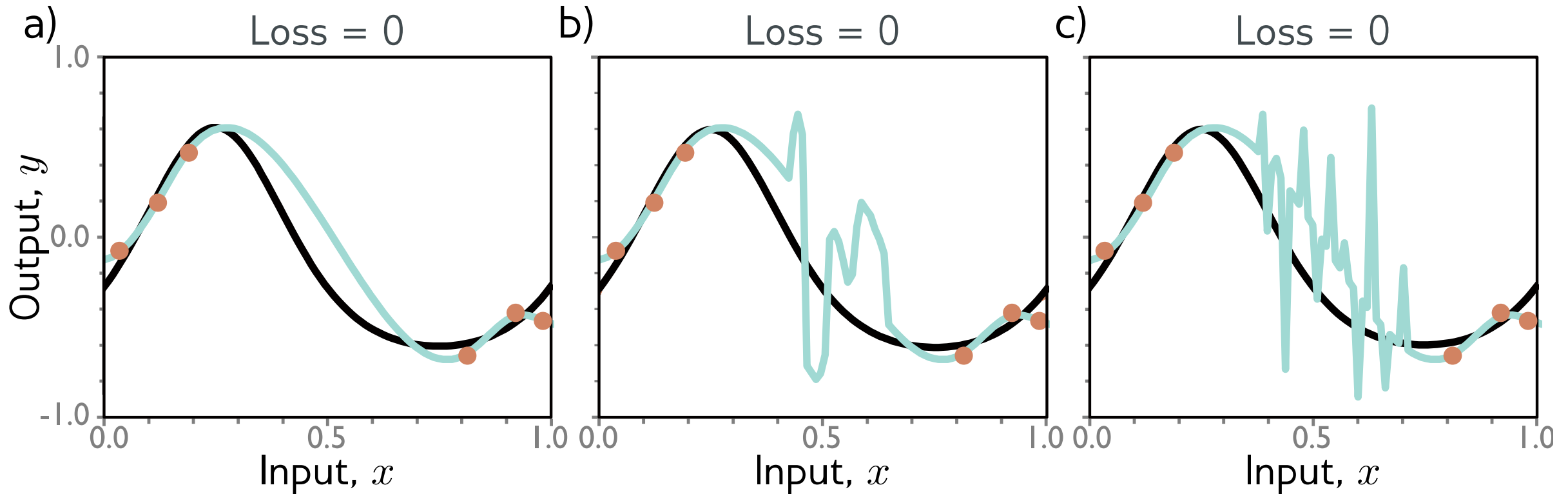Critical regime

# Double descent

Potential explanation:
- can make smoother functions with more hidden units
- being smooth between the datapoints is a reasonable thing to do

But why?

# Double descent



- All of these solutions are equivalent in terms of loss.
- Why should the model choose the smooth solution?
- Tendency of model to choose one solution over another is inductive bias
- Any factor that biases a solution toward a subset of equivalent solutions is known as a regularizer
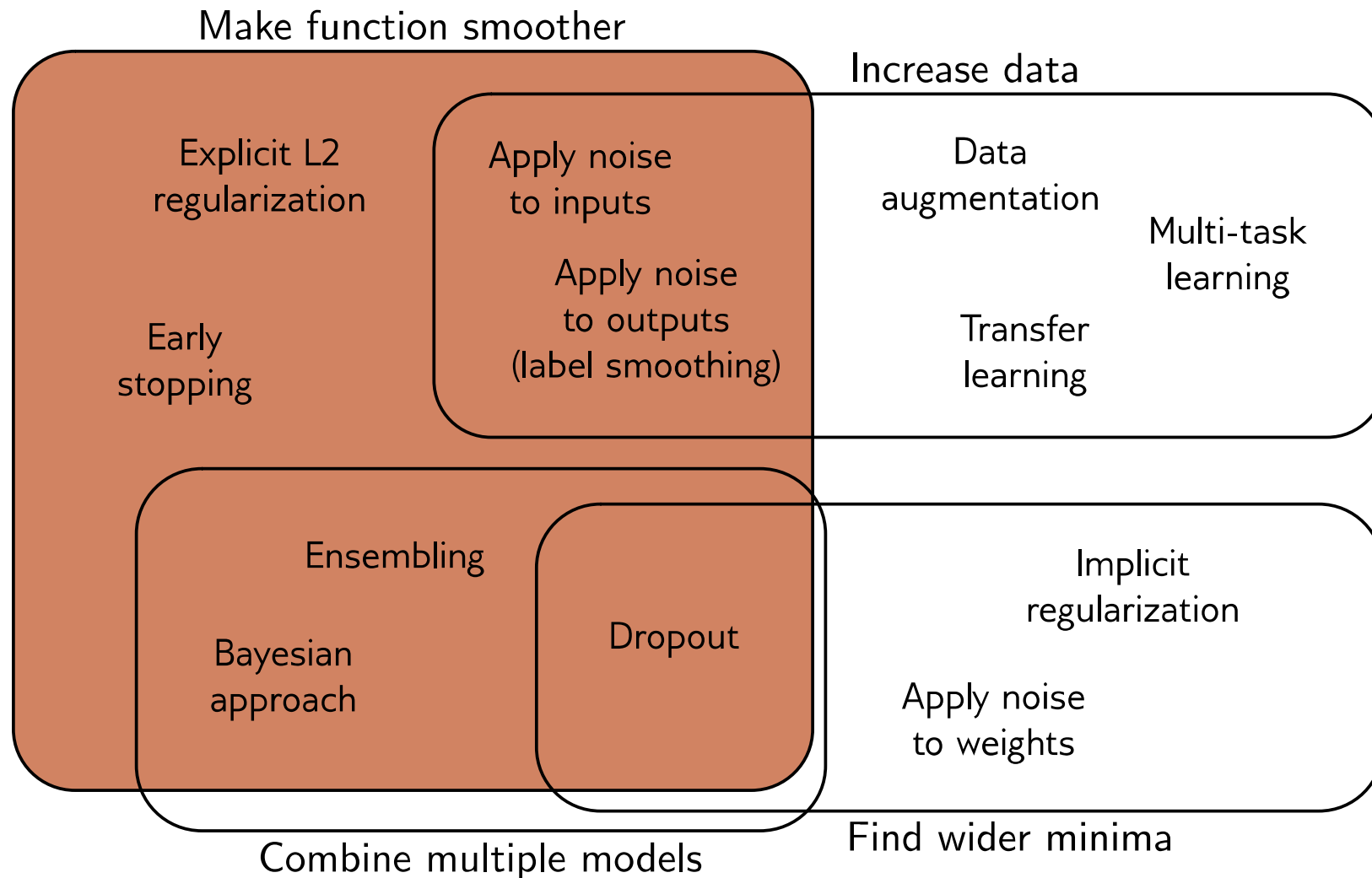
# Choosing hyperparameters

- We don't know the bias or the variance

- We don't know how much capacity to add

- How do we choose capacity in practice?
  - Or model structure
  - Or training algorithm
  - Or learning rate

- Third data set – validation set
  - Train models with different hyperparameters on training set
  - Choose best hyperparameters with validation set
  - Test once with test set

# Regularization

- Why is there a generalization gap between training and test data?
  - Overfitting (model describes statistical peculiarities)
  - Model unconstrained in areas where there are no training examples
- Regularization = methods to reduce the generalization gap
- Technically means adding terms to loss function
- But colloquially means any method (hack) to reduce gap

# Regularization overview

Make function smoother

Increase data

Explicit L2
regularization

Apply noise
to inputs

Data
augmentation

Multi-task
learning

Apply noise
to outputs
(label smoothing)

Transfer
learning

Early
stopping

Ensembling

Implicit
regularization

Bayesian
approach

Dropout

Apply noise
to weights

Combine multiple models

Find wider minima

# Explicit regularization

- Standard loss function:   $\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \big[ \mathrm{L}[\phi] \big]$

$$= \underset{\phi}{\operatorname{argmin}} \left[ \sum_{i=1}^{I} \ell_i[\mathbf{x}_i, \mathbf{y}_i] \right]$$

# Explicit regularization

- Standard loss function: $\hat{\phi} = \underset{\phi}{\operatorname{argmin}}\left[\mathrm{L}[\phi]\right]$

$$= \underset{\phi}{\operatorname{argmin}}\left[\sum_{i=1}^{I} \ell_i[\mathbf{x}_i, \mathbf{y}_i]\right]$$
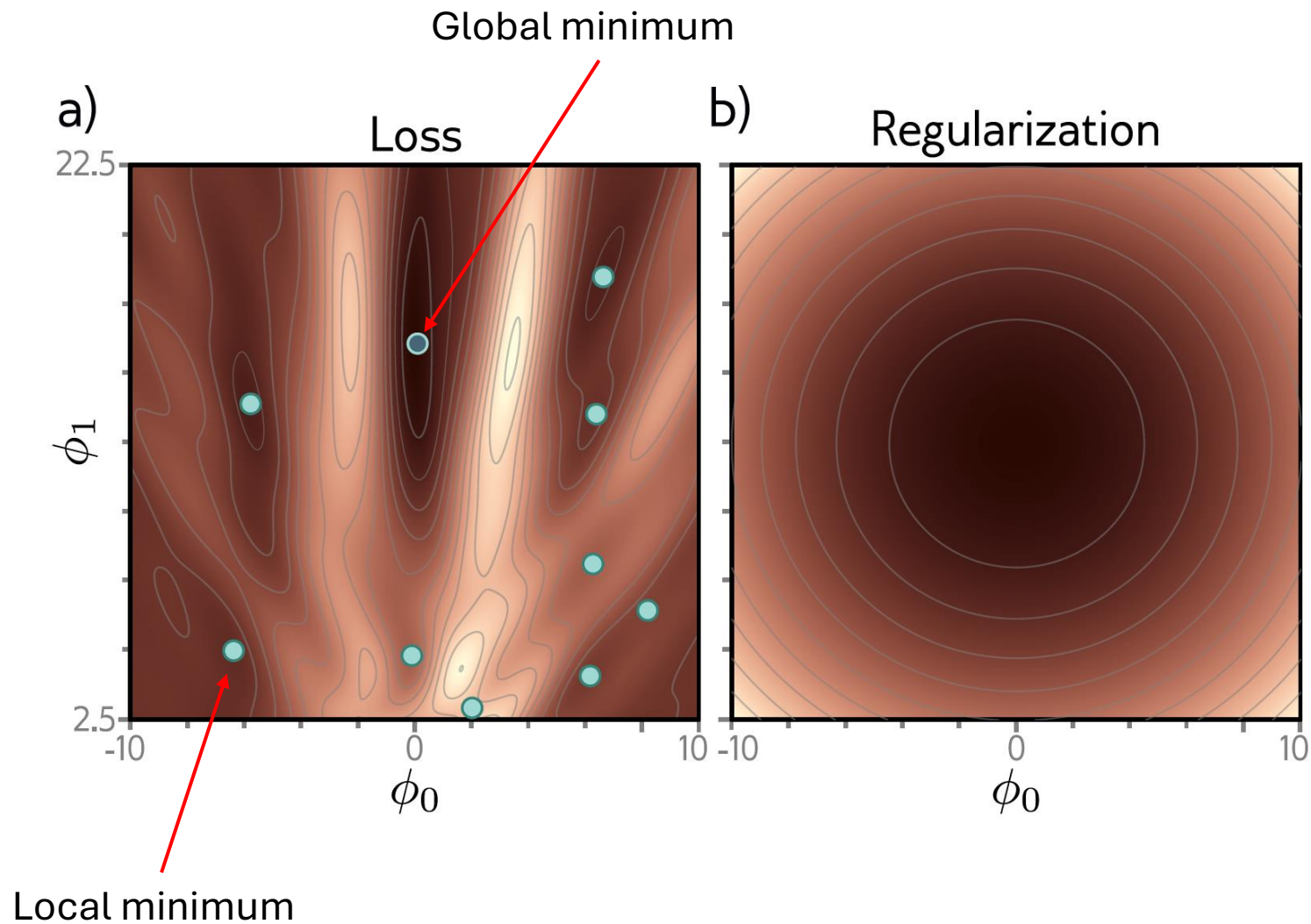
- Regularization adds an extra term

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}}\left[\sum_{i=1}^{I} \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot \mathrm{g}[\phi]\right]$$

# Explicit regularization

- Standard loss function:
$$\hat{\phi} = \underset{\phi}{\mathrm{argmin}}\left[L[\phi]\right]$$

$$= \underset{\phi}{\mathrm{argmin}}\left[\sum_{i=1}^{I}\ell_i[\mathbf{x}_i,\mathbf{y}_i]\right]$$

- Regularization adds an extra term
$$\hat{\phi} = \underset{\phi}{\mathrm{argmin}}\left[\sum_{i=1}^{I}\ell_i[\mathbf{x}_i,\mathbf{y}_i] + \lambda \cdot g[\phi]\right]$$

- Favors some parameters, disfavors others.
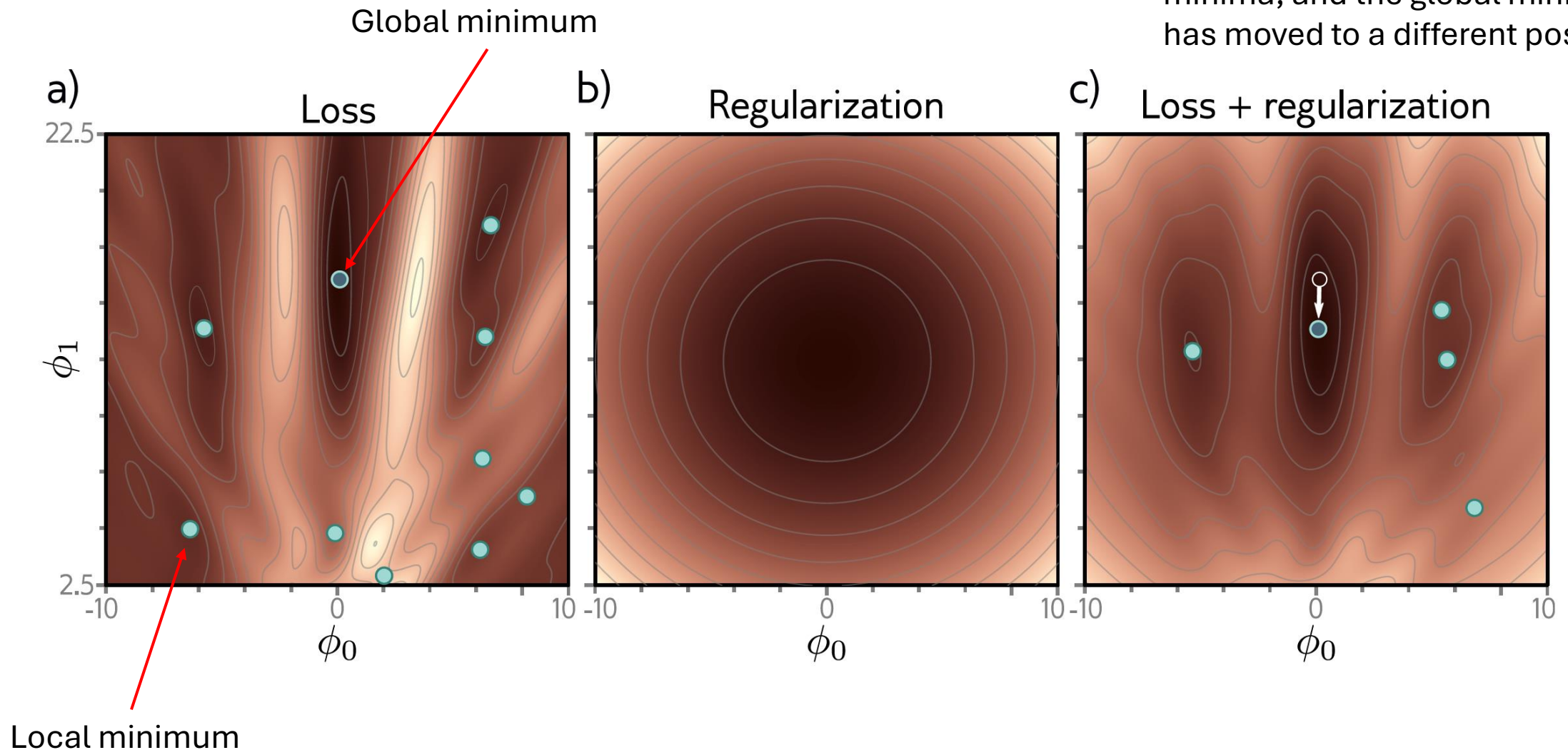- $\lambda > 0$ controls the strength

# Explicit regularization

# Explicit regularization

# Explicit regularization

# Probabilistic interpretation

- Maximum likelihood:

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmax}} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\phi}) \right]$$

# Probabilistic interpretation

- Maximum likelihood:

$$\hat{\boldsymbol{\phi}} = \operatorname*{argmax}_{\boldsymbol{\phi}} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\phi}) \right]$$

- Regularization is equivalent to a adding a prior over parameters

$$\hat{\boldsymbol{\phi}} = \operatorname*{argmax}_{\boldsymbol{\phi}} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\phi}) Pr(\boldsymbol{\phi}) \right]$$

Maximum a posteriori (MAP) criterion

... what you know about parameters *before* seeing the data

# Equivalence

- Explicit regularization:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmin}} \left[ \sum_{i=1}^{I} \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot \mathrm{g}[\phi] \right]$$

- Probabilistic interpretation:

$$\hat{\phi} = \underset{\phi}{\operatorname{argmax}} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{x}_i, \phi) Pr(\phi) \right]$$

# Equivalence

- Explicit regularization:

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ \sum_{i=1}^{I} \ell_i[\mathbf{x}_i, \mathbf{y}_i] + \lambda \cdot \mathrm{g}[\boldsymbol{\phi}] \right]$$

- Probabilistic interpretation:

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmax}} \left[ \prod_{i=1}^{I} Pr(\mathbf{y}_i | \mathbf{x}_i, \boldsymbol{\phi}) Pr(\boldsymbol{\phi}) \right]$$

- Mapping: $\quad \lambda \cdot \mathrm{g}[\boldsymbol{\phi}] = -\log[Pr(\boldsymbol{\phi})]$

# L2 Regularization

- Can only use very general terms
- Most common is L2 regularization
- Favors smaller parameters

$$\hat{\boldsymbol{\phi}} = \underset{\boldsymbol{\phi}}{\operatorname{argmin}} \left[ \mathrm{L}[\boldsymbol{\phi}, \{\mathbf{x}_i, \mathbf{y}_i\}] + \lambda \sum_j \phi_j^2 \right]$$

- Also called Tikhonov regularization, ridge regression, or Frobenius norm regularization
- In neural networks, usually just for weights (not for the biases) and called weight decay

# Why does L2 regularization help?

- Discourages memorizing the data (overfitting)
- Encourages smoothness between datapoints

# L2 regularization

# Implicit regularization

Gradient descent approximates a differential equation (infinitesimal step size)

Finite step size equivalent to regularization

Add in that regularization and differential equation converges to same place



a) Loss

b) Regularization

c) Loss + regularization

Continuous gradient descent path. The finite step size causes reaching a different final position

# Implicit regularization

- Gradient descent disfavors areas where gradients are steep

$$\tilde{\mathrm{L}}_{GD}[\phi] = \mathrm{L}[\phi] + \frac{\alpha}{4}\left\|\frac{\partial L}{\partial \phi}\right\|^2$$

# Implicit regularization

- Gradient descent disfavors areas where gradients are steep

$$\tilde{\text{L}}_{GD}[\phi] = \text{L}[\phi] + \frac{\alpha}{4}\left\|\frac{\partial L}{\partial \phi}\right\|^2$$

- SGD likes all batches to have similar gradients

$$\tilde{\text{L}}_{SGD}[\phi] = \tilde{\text{L}}_{GD}[\phi] + \frac{\alpha}{4B}\sum_{b=1}^{B}\left\|\frac{\partial L_b}{\partial \phi} - \frac{\partial L}{\partial \phi}\right\|^2$$

$$= \text{L}[\phi] + \frac{\alpha}{4}\left\|\frac{\partial L}{\partial \phi}\right\|^2 + \frac{\alpha}{4B}\sum_{b=1}^{B}\left\|\frac{\partial L_b}{\partial \phi} - \frac{\partial L}{\partial \phi}\right\|^2$$

- Depends on learning rate – perhaps why larger learning rates generalize better.

a) Loss, $L[\phi]$

b) GD modification

Global minimum

Implicit regularization term from gradient descent penalizes the squared gradient magnitude.

c) SGD modification

d) Modified loss, $\tilde{L}_{SGD}[\phi]$

Additional implicit regularization from stochastic gradient descent penalizes the variance of the batch gradients

Global minimum can be different from that of the original loss function $L[\varphi]$

**a)** MNIST1D no label noise

Legend:
- Full batch, LR = 0.5
- Full batch, LR = 0.1
- Full batch, LR = 0.05

Test

Train

Error (%) vs Hidden layer size

Effect of learning rate (LR) and batch size for 4000 training and 4000 test examples from MNIST-1D for a NN with two hidden layers

**b)** MNIST1D no label noise

Legend:
- Batch size 10, LR = 0.1
- Batch Size 100, LR = 0.1
- Full batch (4000), LR = 0.1

Test

Train

Error (%) vs Hidden layer size

Generally, performance is:
- best for larger learning rates
- best with smaller batches

# Early stopping

- If we stop training early,  weights don't have time to overfit to noise
- Weights start small, don't have time to get large
- Reduces effective model complexity
- Known as early stopping
- Don't have to re-train

a) Iter = 0, Loss = 32.24
b) Iter = 1000, Loss = 1.64
c) Iter = 5000, Loss = 1.10
d) Iter = 10000, Loss = 0.80
e) Iter = 50000, Loss = 0.36
f) Iter = 200000, Loss = 0.16

# Ensembling

- Average together the predictions of several models – an ensemble
- Can take mean or median
- Different initializations / different models
- Different subsets of the data resampled with replacements – bagging
- The assumption is that model errors are independent and will cancel out

a) Original  b) Model 1  c) Model 2  d) Model 3  e) Model 4  f) Ensemble

# Dropout



a)

b)

c)

d)

Training:
- Drop units with probability p

Testing:
- Do not drop any unit
- Scale the output of each unit with the value p

# Dropout

a) Original

b) Turn off hidden unit 8

c) 2000 iters dropout (7/8/9)

Output, $y$

Input, $x$

Can eliminate sudden changes in function that are far from data and don't contribute to training loss

Dropout in iterative optimization: A probabilistic process to 'augment' the training set during iterative training and increase invariance:

Standard neural network training

Dropout-based training
At each iteration, each neuron is active with probability p (using Bernoulli distribution and cut-off value of e.g. p = 0.5)



All iterations

Training iteration 1 ••• Training iteration t

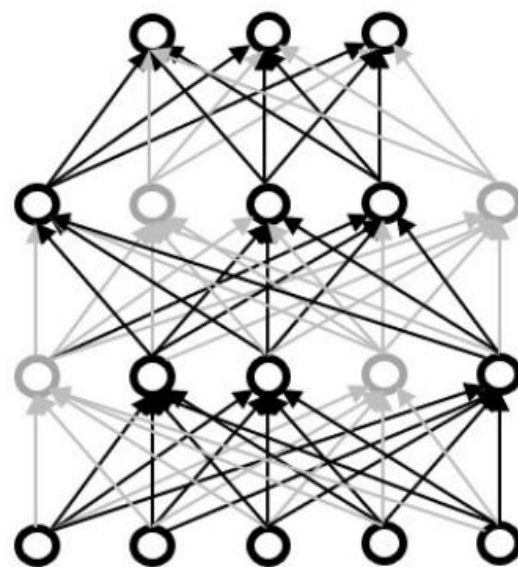# Continuous Dropout-based iterative optimization:

Standard neural network training

Continuous dropout-based training
At each iteration, each neuron is 'suppressed' (multiplied) with masks sampled from $\mu \sim U(0, 1)$ or $g \sim N(0.5, \sigma^2)$
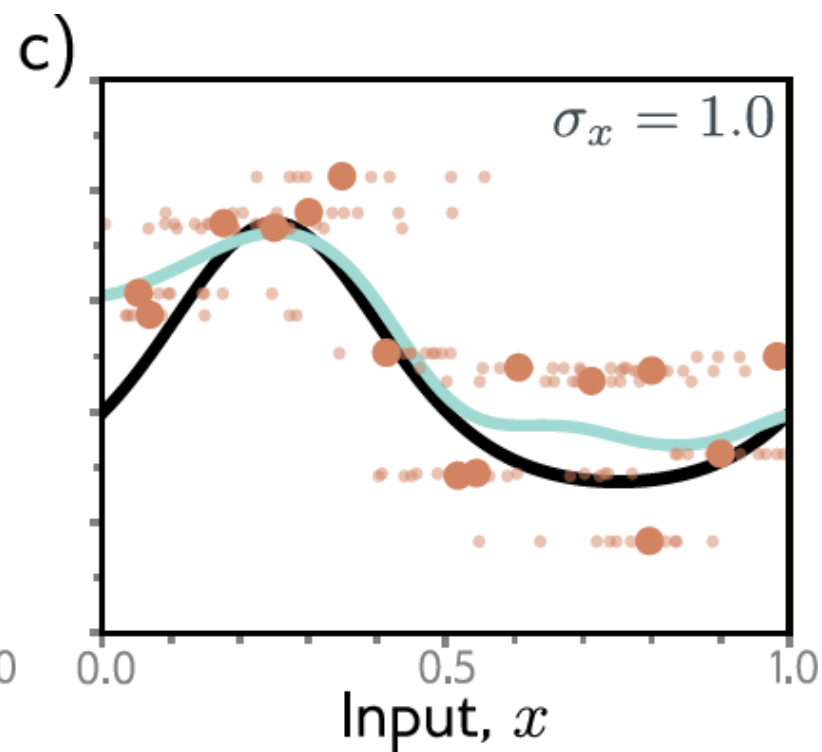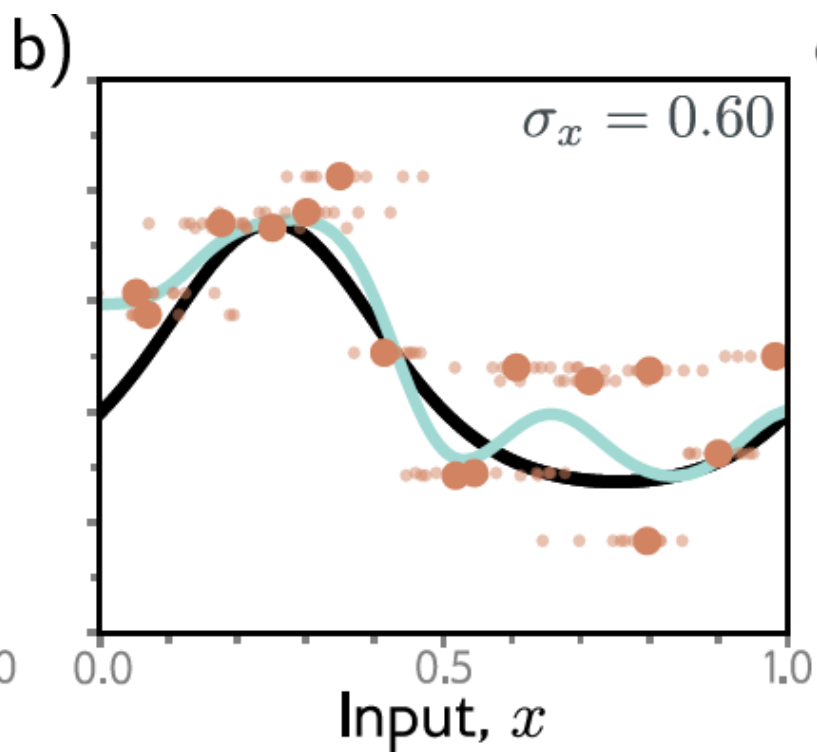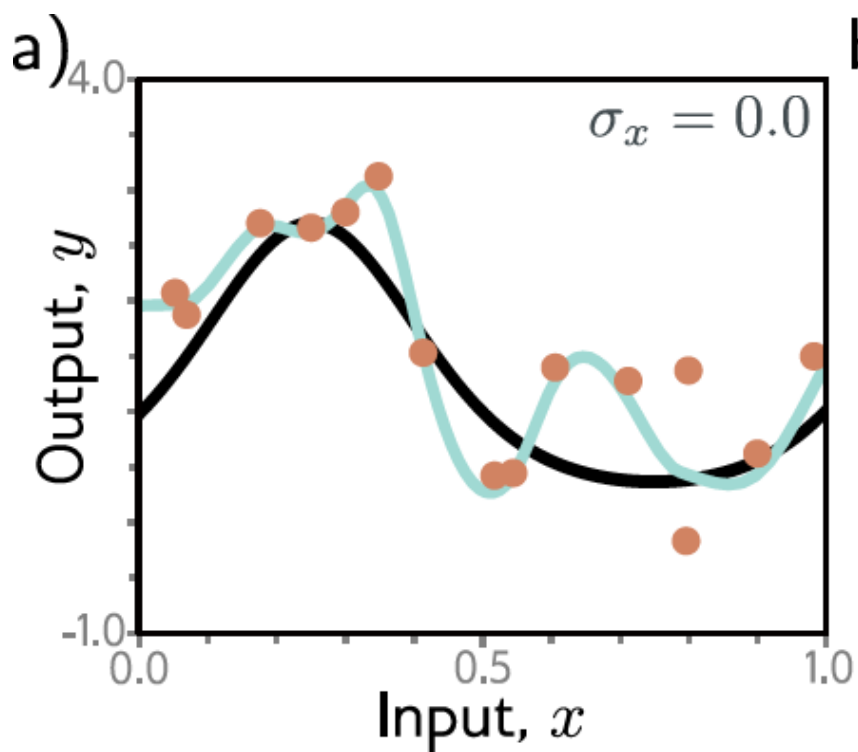


All iterations

Training iteration 1 $\cdots$ Training iteration t

# Adding noise



- to inputs
- to weights
- to outputs (labels)

# Bayesian approaches

- There are many parameters compatible with the data

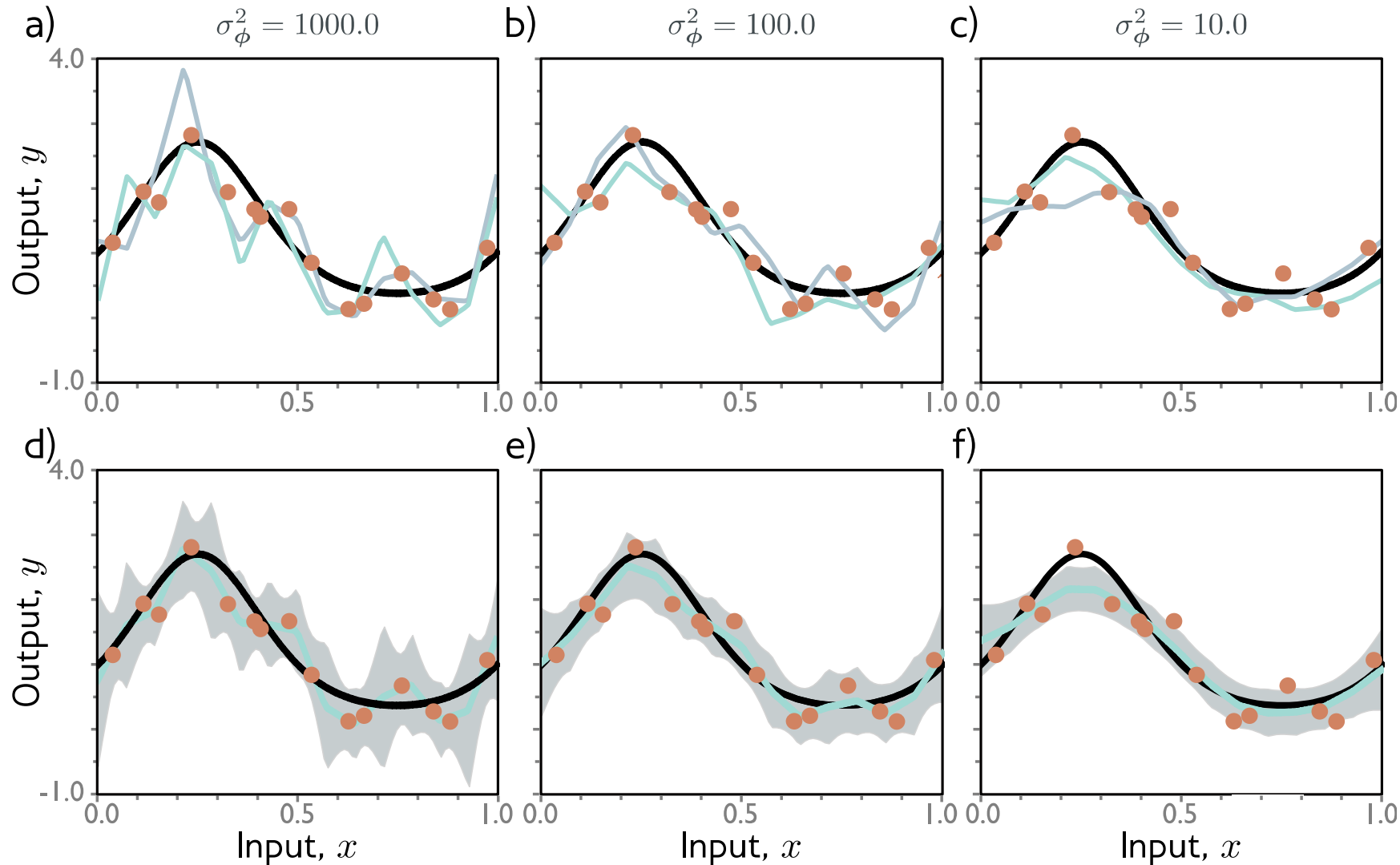- Can find a probability distribution over them

Prior info about parameters

$$Pr(\phi|\{\mathbf{x}_i, \mathbf{y}_i\}) = \frac{\prod_{i=1}^{I} Pr(\mathbf{y}_i|\mathbf{x}_i, \phi)Pr(\phi)}{\int \prod_{i=1}^{I} Pr(\mathbf{y}_i|\mathbf{x}_i, \phi)Pr(\phi)d\phi}$$

# Bayesian approaches

- There are many parameters compatible with the data

- Can find a probability distribution over them

$$Pr(\boldsymbol{\phi}|\{\mathbf{x}_i, \mathbf{y}_i\}) = \frac{\prod_{i=1}^{I} Pr(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\phi})Pr(\boldsymbol{\phi})}{\int \prod_{i=1}^{I} Pr(\mathbf{y}_i|\mathbf{x}_i, \boldsymbol{\phi})Pr(\boldsymbol{\phi})d\boldsymbol{\phi}}$$

- Take all possible parameters into account when make prediction

$$Pr(\mathbf{y}|\mathbf{x}, \{\mathbf{x}_i, \mathbf{y}_i\}) = \int Pr(\mathbf{y}|\mathbf{x}, \boldsymbol{\phi})Pr(\boldsymbol{\phi}|\{\mathbf{x}_i, \mathbf{y}_i\})d\boldsymbol{\phi}$$

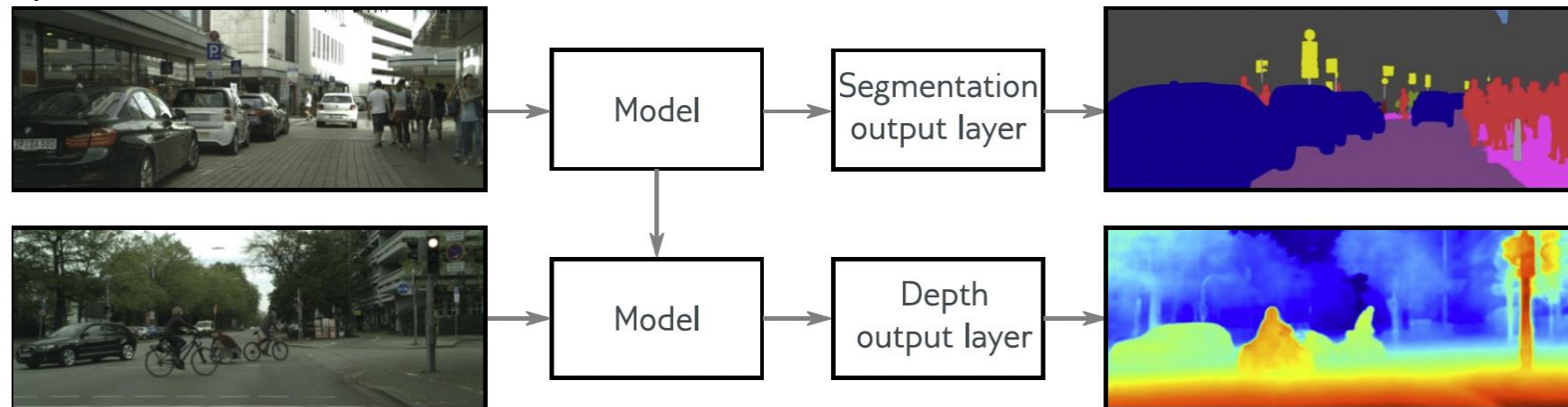Infinite weighted ensemble

# Bayesian approaches



a–c) Two sets of parameters (cyan and gray curves) sampled from the posterior using normally distributed priors with mean zero and three variances.
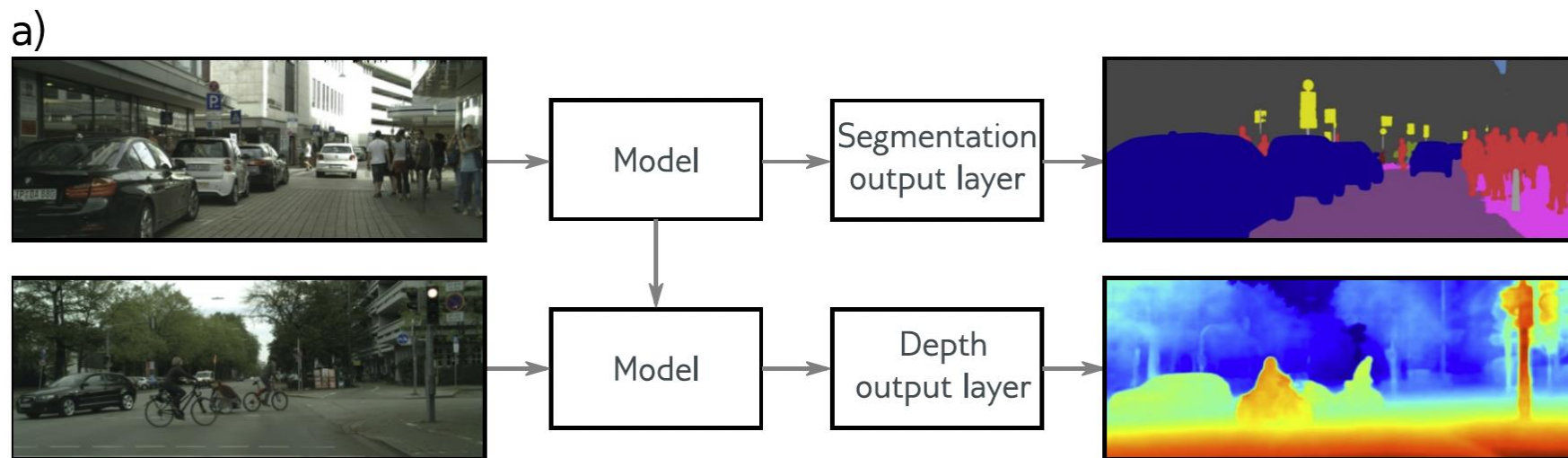
d–f) Inference proceeds by taking a weighted sum over all possible parameter values where the weights are the posterior probabilities. This produces both a prediction of the mean (cyan curves) and the associated uncertainty (gray region is two standard deviations).
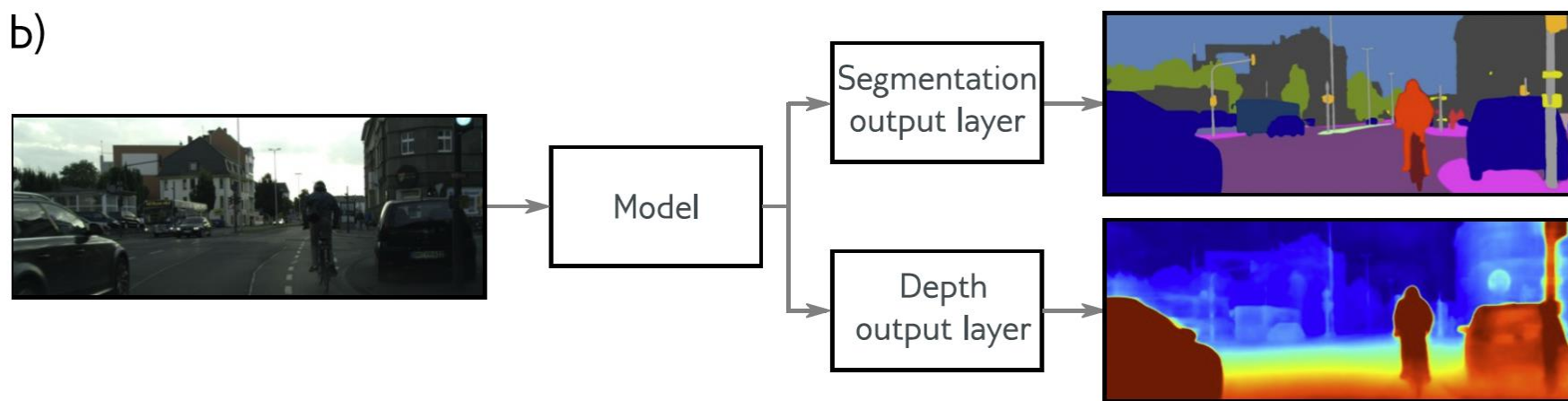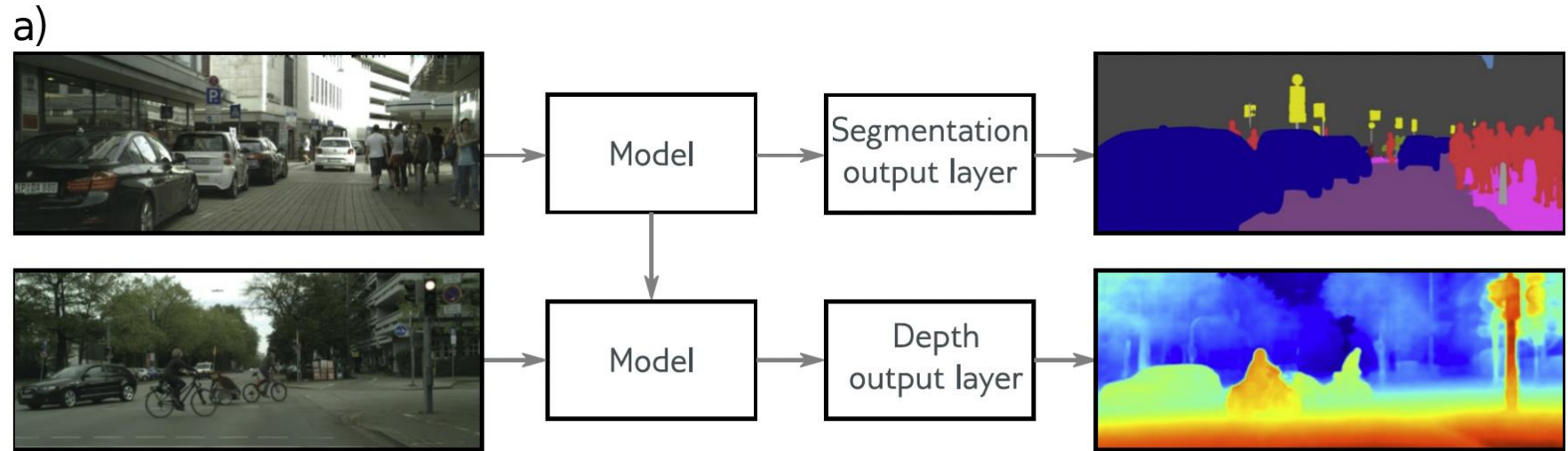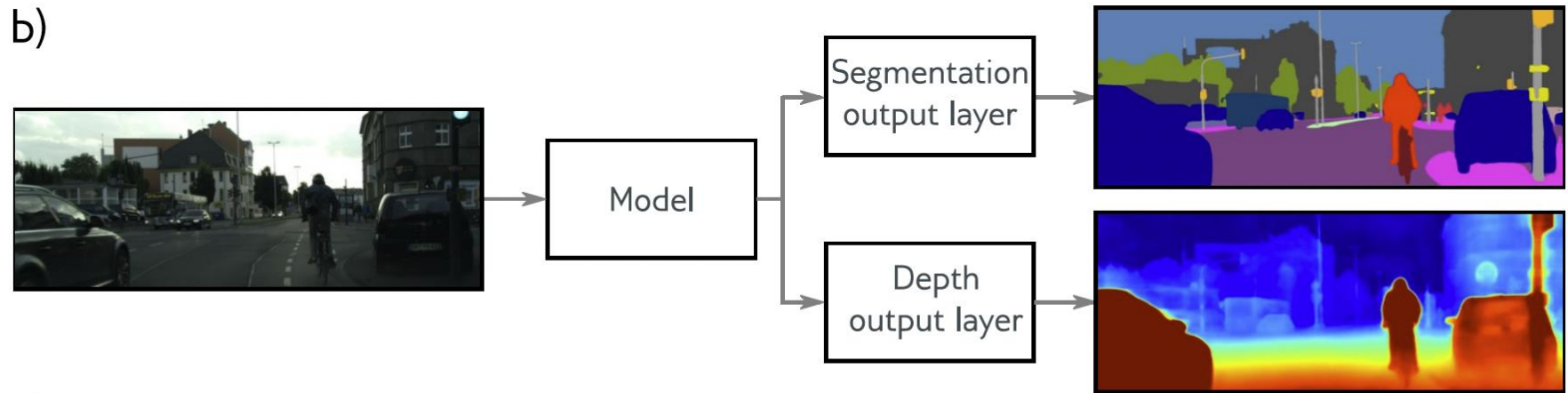
- Transfer learning
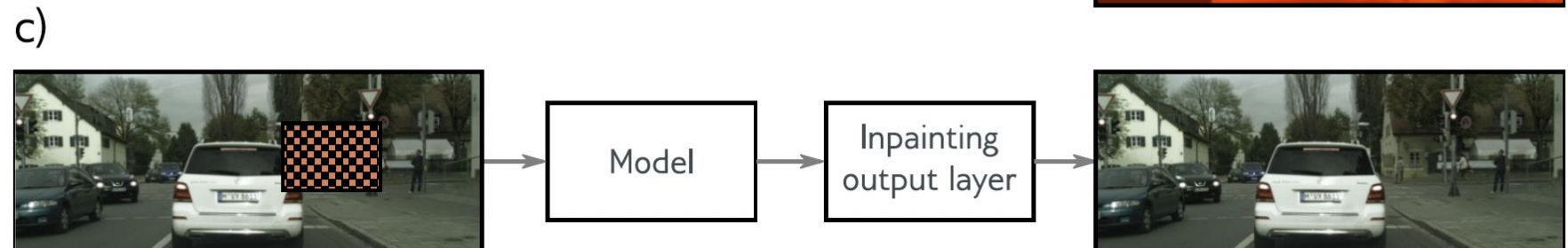
a)

- Transfer learning

a)



- Multi-task learning

b)

- Transfer learning



a)

- Multi-task learning



b)

- Self-supervised learning



c)

# Data augmentation



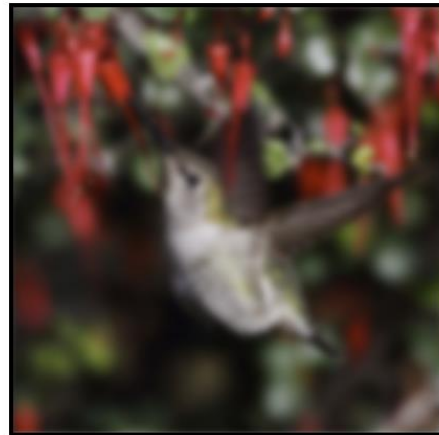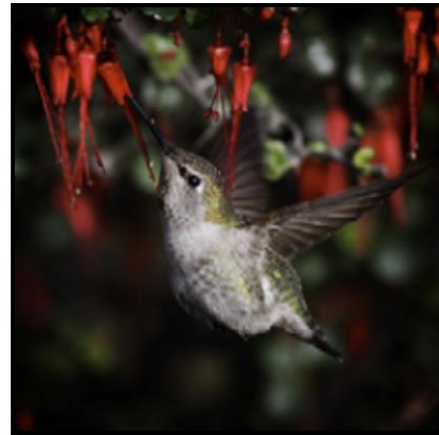a) Original    b) Flip    c) Rotate and crop    d) Vertical stretch

e) Color balance    f) Blur    g) Vignette    h) Pincushion

# Regularization overview



Make function smoother

Increase data

Explicit L2 regularization

Apply noise to inputs

Data augmentation

Multi-task learning

Early stopping

Apply noise to outputs (label smoothing)

Transfer learning

Ensembling

Implicit regularization

Bayesian approach

Dropout

Apply noise to weights

Combine multiple models

Find wider minima