# The Effect of Batch Normalization on Deep Convolutional Neural Networks
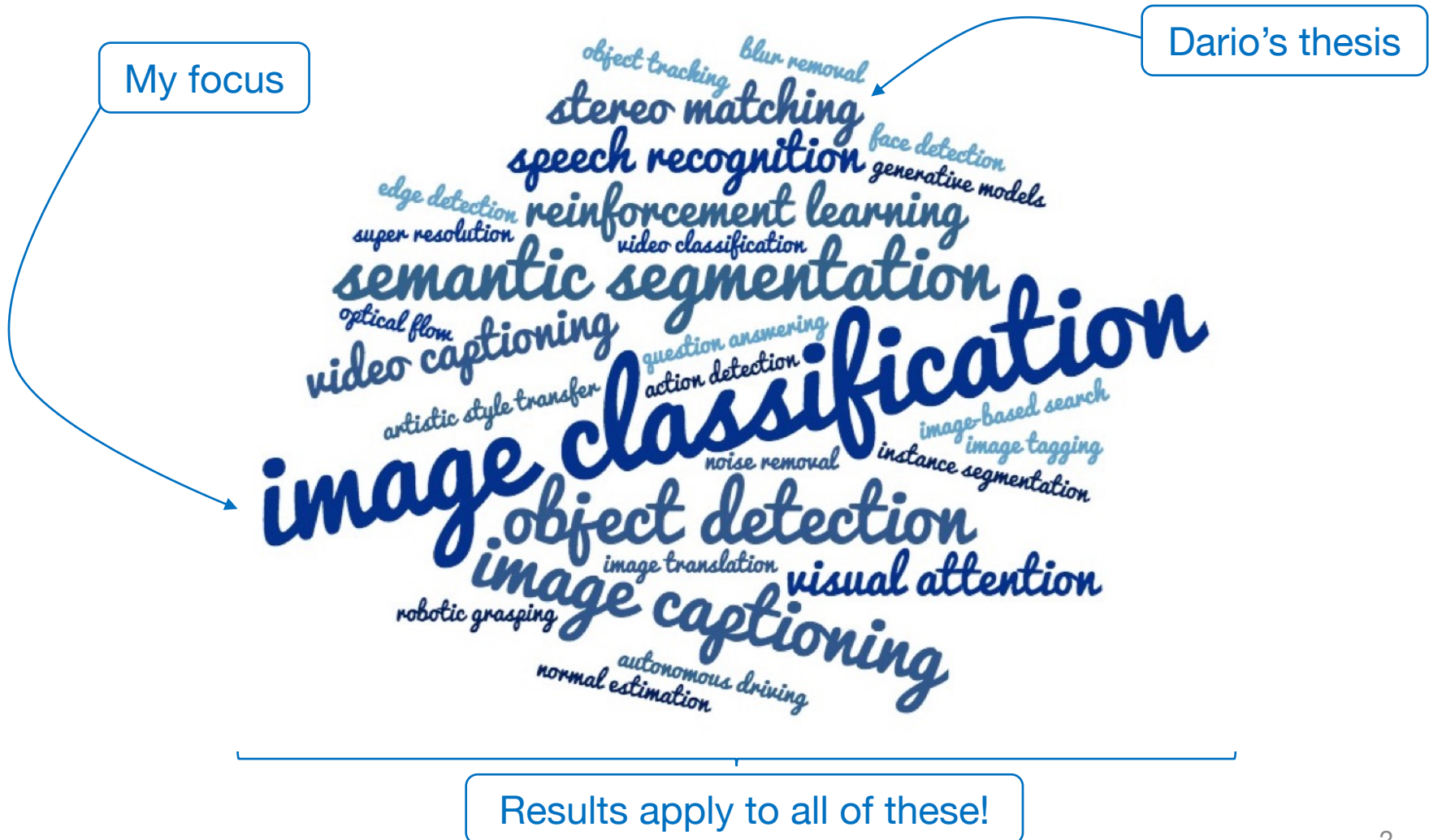
Fabian Schilling
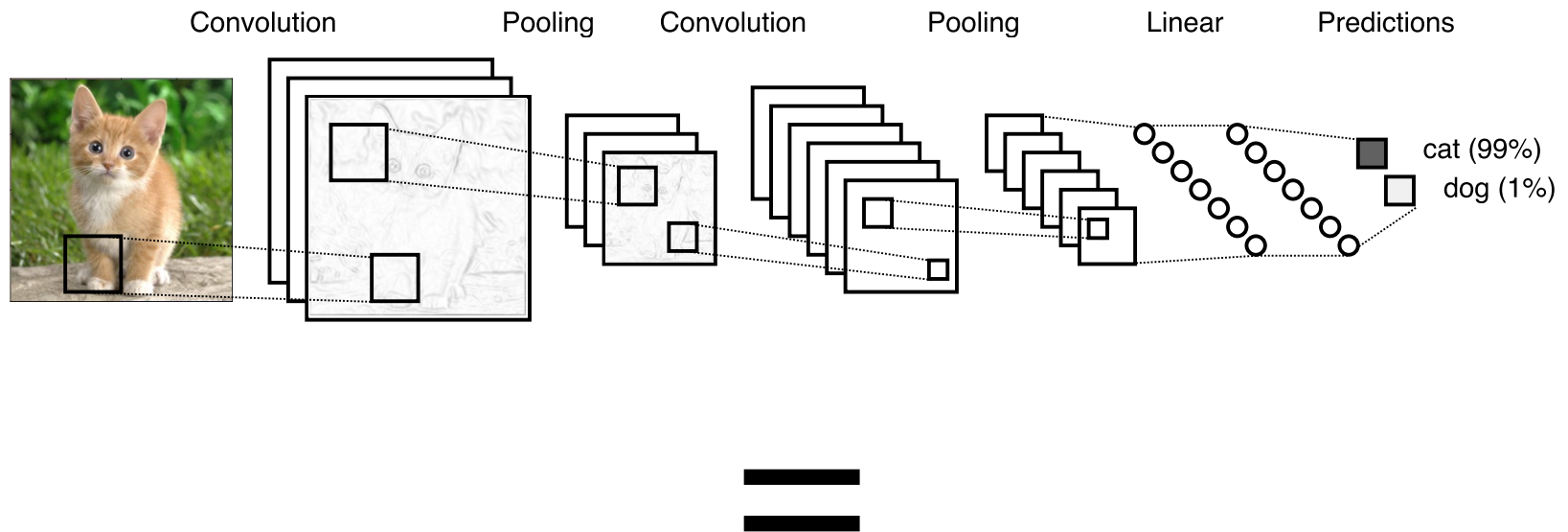
# Convolutional networks are versatile



My focus

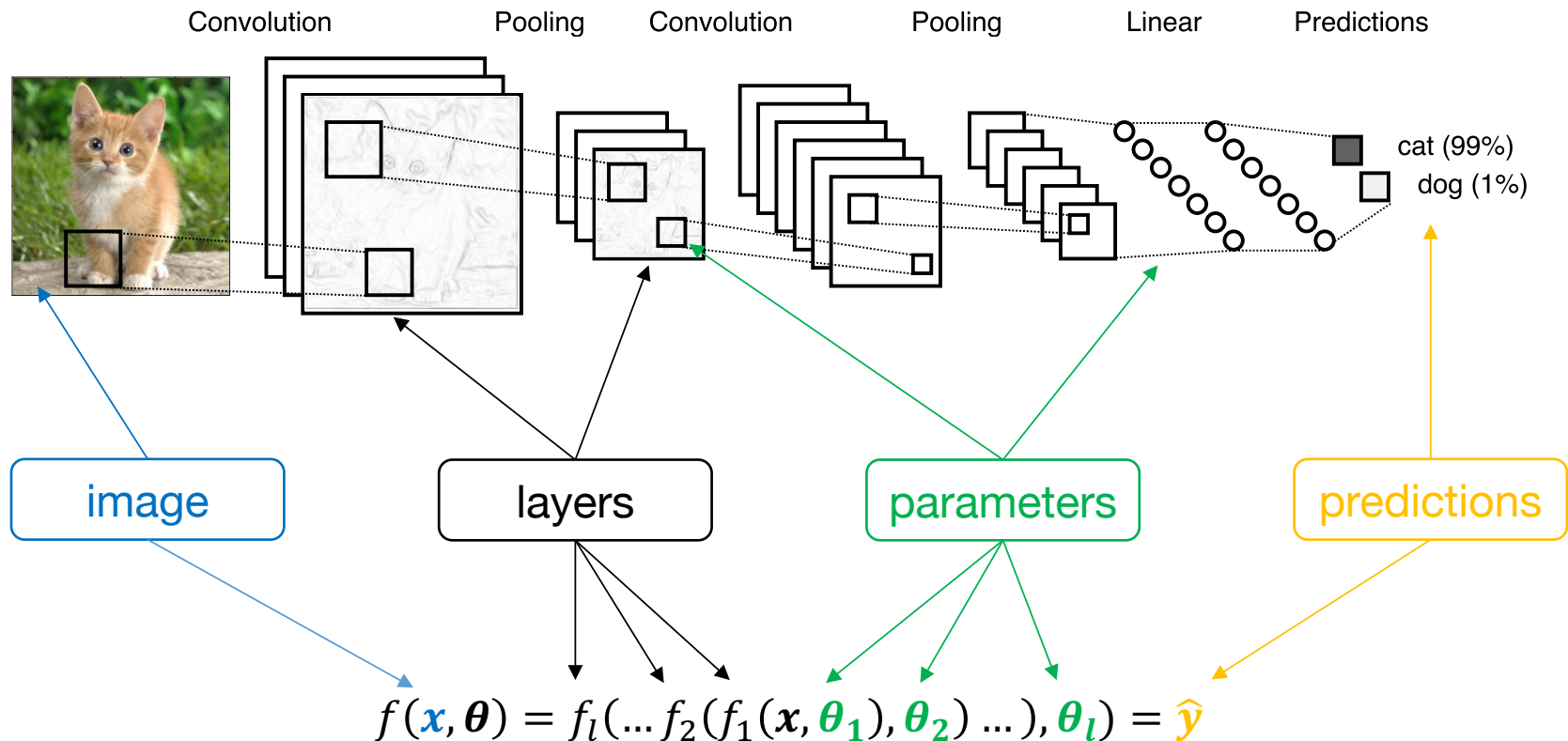Dario's thesis

object tracking • blur removal
stereo matching • face detection
speech recognition • generative models
edge detection • reinforcement learning
super resolution • video classification
semantic segmentation
optical flow
video captioning • question answering
artistic style transfer • action detection • image-based search
image classification • image tagging
noise removal • instance segmentation
object detection
image translation • visual attention
image captioning
robotic grasping
autonomous driving
normal estimation

Results apply to all of these!

# Background & theory

What are the prerequisites for batchnorm?

# Overview of convnet architecture



Convolution  Pooling  Convolution  Pooling  Linear  Predictions
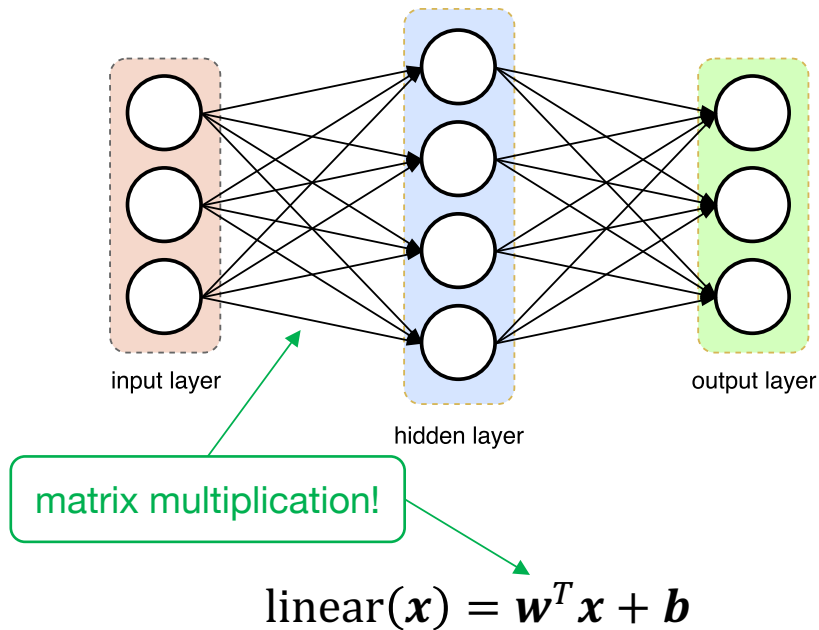
cat (99%)
dog (1%)

=

$$f(\boldsymbol{x}, \boldsymbol{\theta}) = f_l(\dots f_2(f_1(\boldsymbol{x}, \boldsymbol{\theta_1}), \boldsymbol{\theta_2}) \dots), \boldsymbol{\theta_l}) = \widehat{\boldsymbol{y}}$$
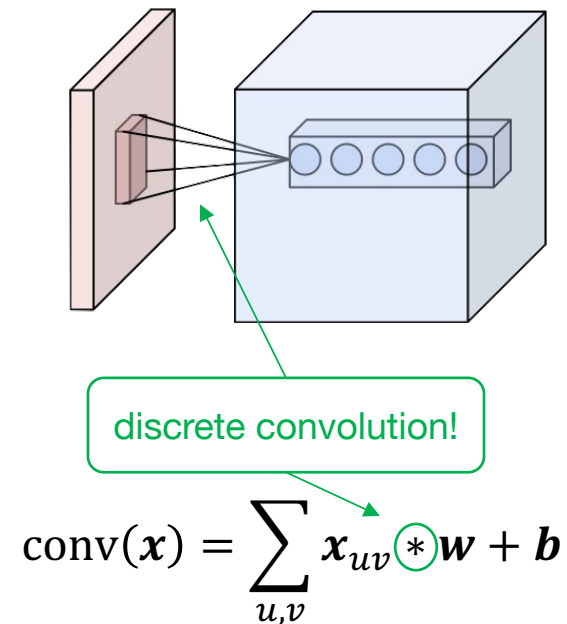
# Overview of convnet architecture



$$f(\boldsymbol{x}, \boldsymbol{\theta}) = f_l(\dots f_2(f_1(\boldsymbol{x}, \boldsymbol{\theta_1}), \boldsymbol{\theta_2}) \dots), \boldsymbol{\theta_l}) = \widehat{\boldsymbol{y}}$$

# Linear and convolutional layer

**Linear layer**



input layer

hidden layer

output layer

matrix multiplication!

$$\text{linear}(\boldsymbol{x}) = \boldsymbol{w}^T \boldsymbol{x} + \boldsymbol{b}$$

**Convolutional layer**



discrete convolution!

$$\text{conv}(\boldsymbol{x}) = \sum_{u,v} \boldsymbol{x}_{uv} * \boldsymbol{w} + \boldsymbol{b}$$

# Pooling and activation layer

**Pooling layer**

| 4 | 5 | 9 | 6 |
|---|---|---|---|
| 8 | 7 | 2 | 0 |
| 7 | 5 | 9 | 0 |
| 3 | 5 | 1 | 1 |

| 8 | 9 |
|---|---|
| 7 | 9 |

just downsampling!

$$\text{maxpool}(\boldsymbol{x}) = \max_{u,v} \ \boldsymbol{x}_{uv}$$

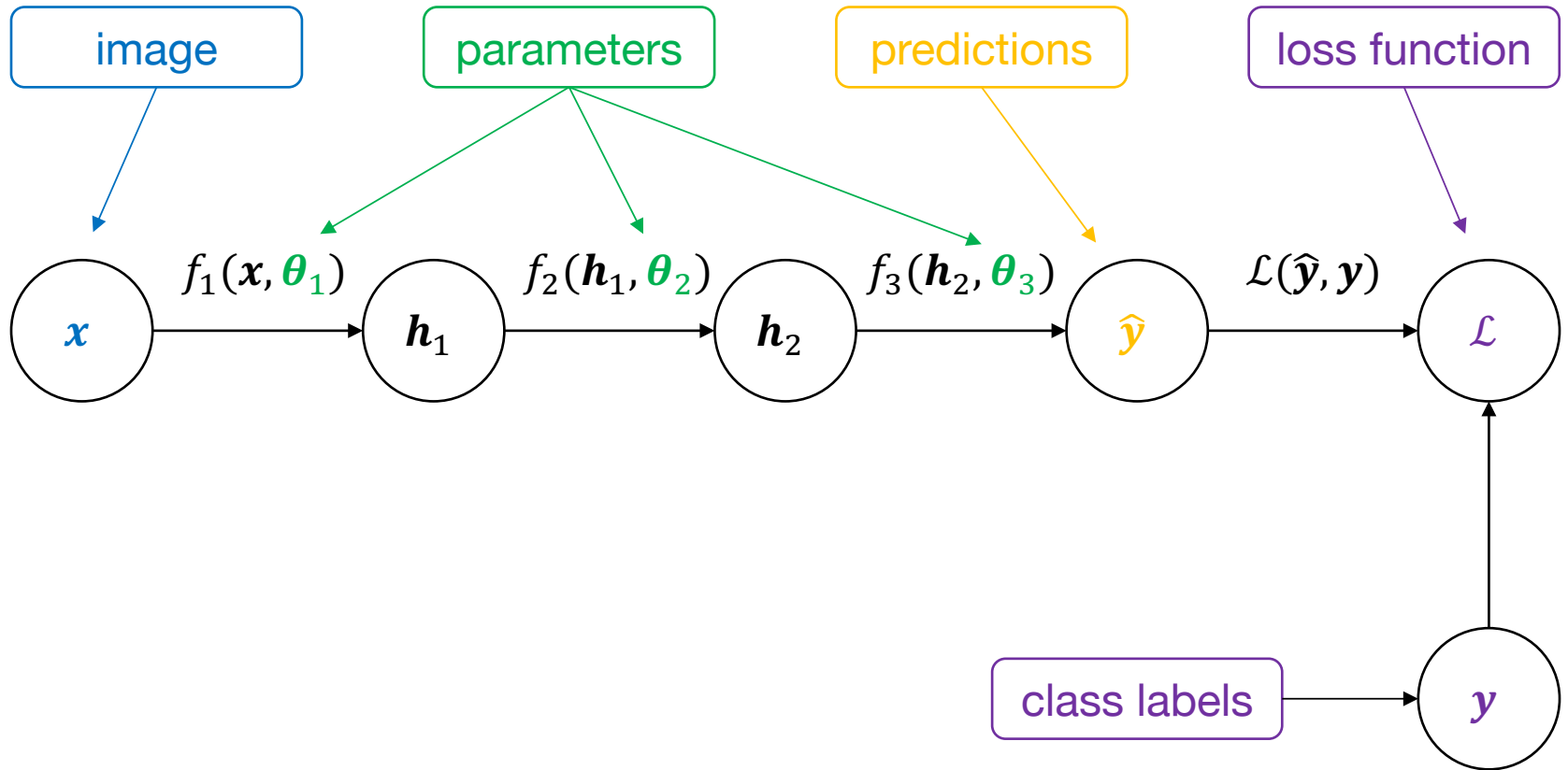**Activation layer**

these are cruicial!

— Sigmoid
-- ReLU
······ ELU

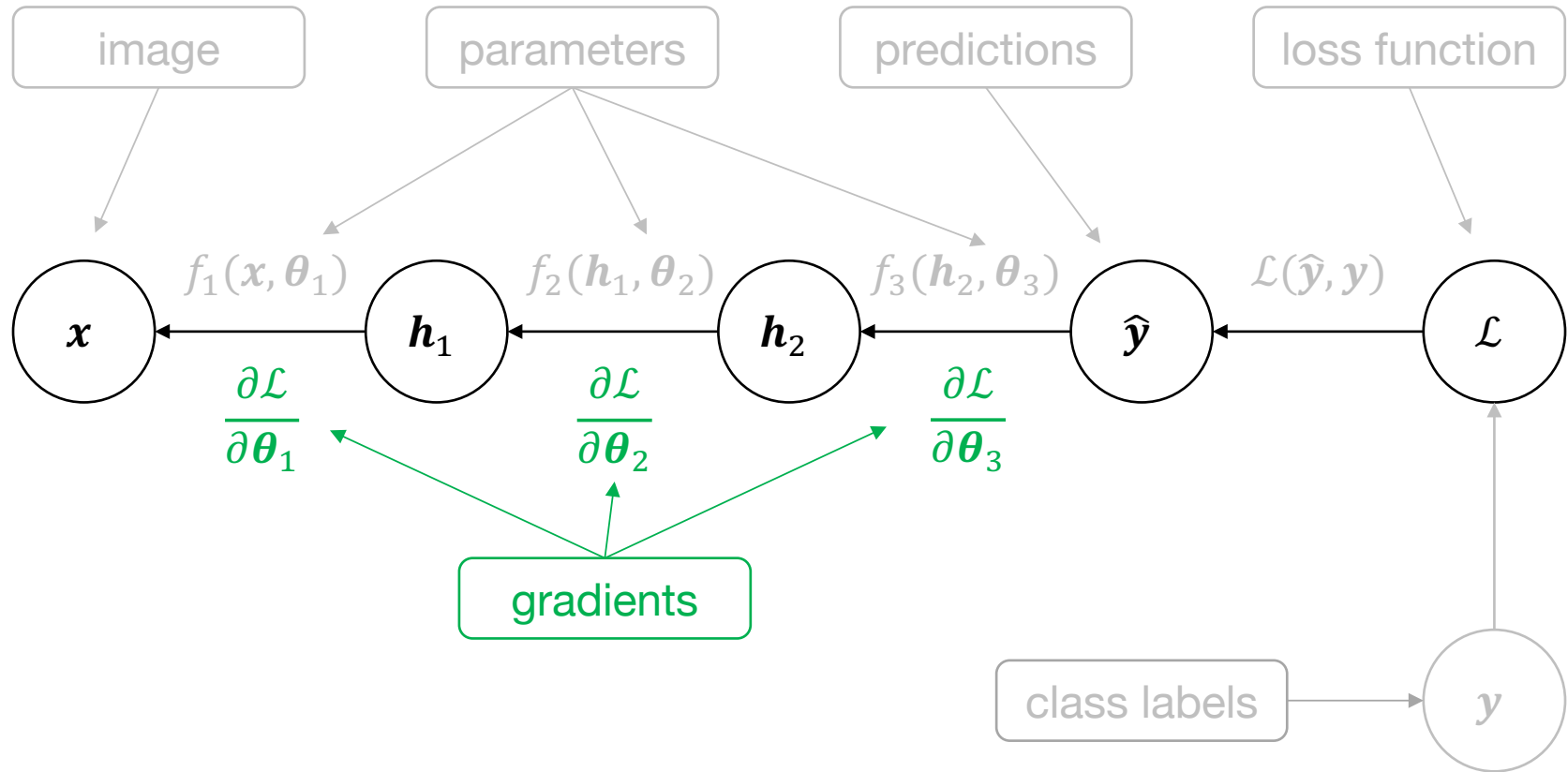$$\text{Sigmoid}(x) = \frac{1}{1 + e^{-x}}$$

$$\text{ReLU}(x) = \max(x, 0)$$

$$\text{ELU}(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha(e^x - 1) & \text{if } x \le 0 \end{cases}$$

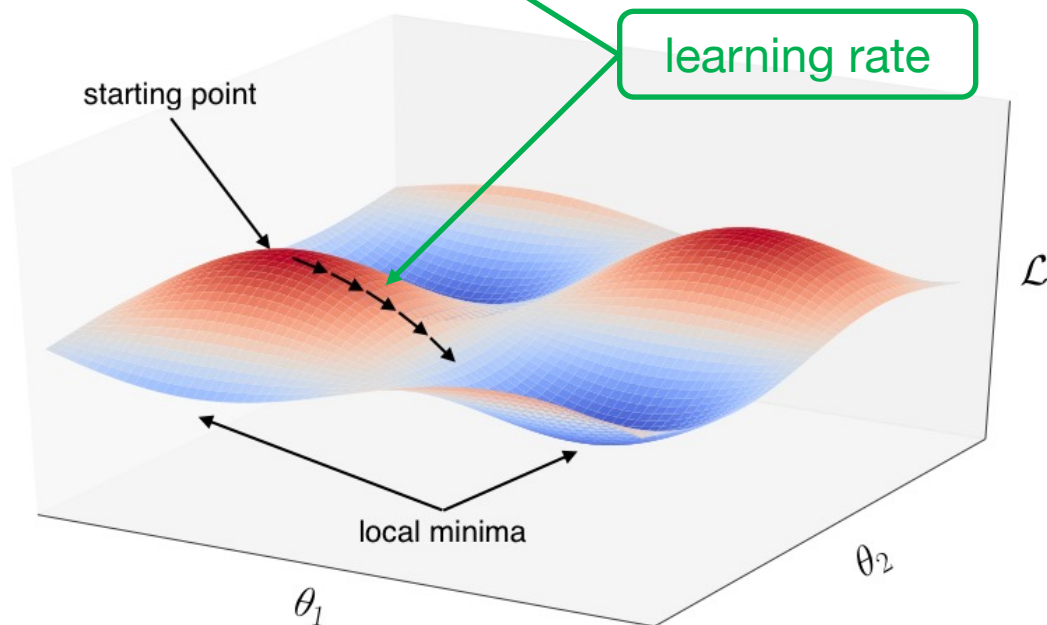# Backpropagation: forward pass

# Backpropagation: backward pass



key takeaway: backprop works for any (differentiable) computational DAG!

# Mini-batch stochastic gradient descent

batch size

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \cdot \frac{1}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \frac{\partial \mathcal{L}(\hat{\boldsymbol{y}}_i, \boldsymbol{y}_i)}{\partial \boldsymbol{\theta}}$$

learning rate

starting point

local minima

$\theta_1$
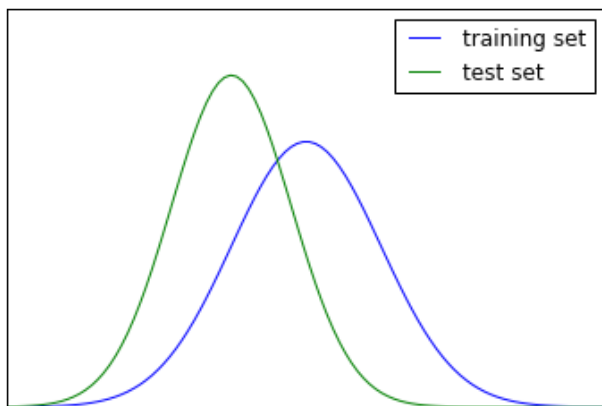
$\theta_2$

$\mathcal{L}$

Why use mini-batches?

- **Memory constraints**, entire dataset might not fit

- **Time constraints**, can't parallelize one example

- Noisy updates help escape local minima!
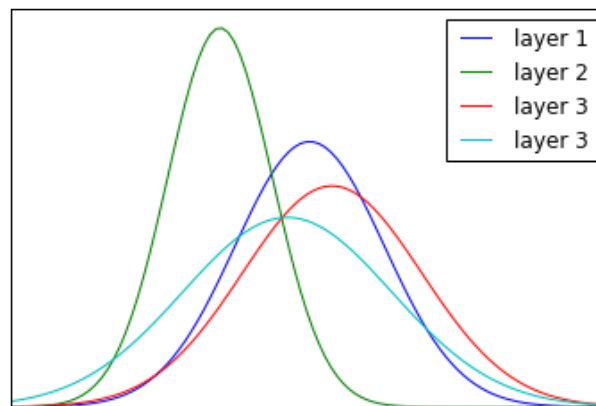
# The problem: internal covariate shift

**Covariate shift**



Distribution of training and test set different

solution: data preprocessing

**Internal covariate shift \***



Distributions of **individual network layers** change during training

solution: batchnorm!

* Ioffe et. al "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift" (2015)

# The solution: batch normalization*

- Subtract batch mean $\mu_{\mathcal{B}}$
- Divide by batch standard deviation $\sigma_{\mathcal{B}}$
- Add $\epsilon$ for numerical stability

$\mu_{\mathcal{B}}$ and $\sigma_{\mathcal{B}}$ are replaced with population statistics at test time!

$$\text{batchnorm}(x) = \gamma \cdot \frac{x - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} + \beta$$

- Scale and shift parameters $\gamma$ and $\beta$ can be learned with backprop
- Batch normalization is differentiable!

# Objective: the effect of batchnorm

- What **general** effects does batchnorm have on…
  - …classification accuracy?
  - …convergence speed?

- What **specific** effects does batchnorm have on...
  - …activation functions?
  - …learning rates?
  - …weight initialization?
  - …regularization methods?
  - …batch sizes?
  - …wall time?

# Experimental setup

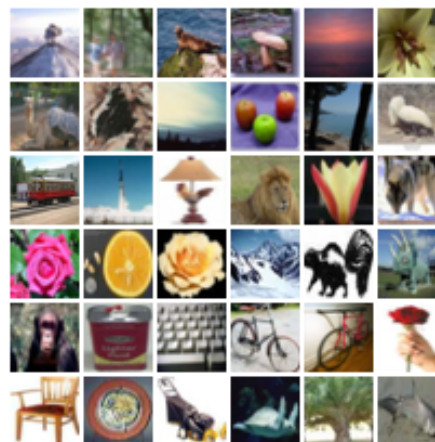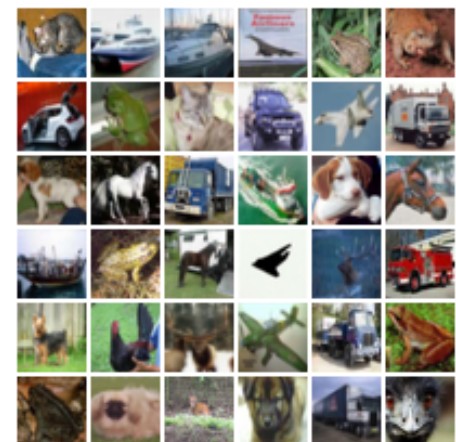How did we conduct our experiments?

# Datasets used in experiments

| MNIST | SVHN | CIFAR10 | CIFAR100 |
|-------|------|---------|----------|



- Handwritten digits
- **32x32 grayscale!**
- 10 classes
- Approx. balanced
- 60k /10k split
- 99.8% best acc.

- Housing numbers
- 32x32 color
- 10 classes
- **Unbalanced!**
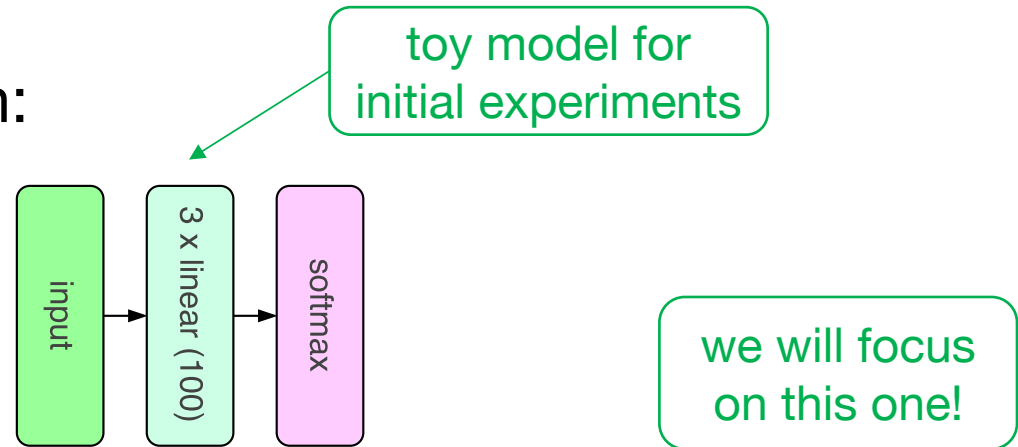- 73k / 26k split
- 98.3% best acc.

- Objects & animals
- 32x32 color
- 10 classes
- Perfectly balanced
- 50k / 10k split
- 96.5% best acc.

- Objects & animals
- 32x32 color
- **100 classes!**
- Perfectly balanced
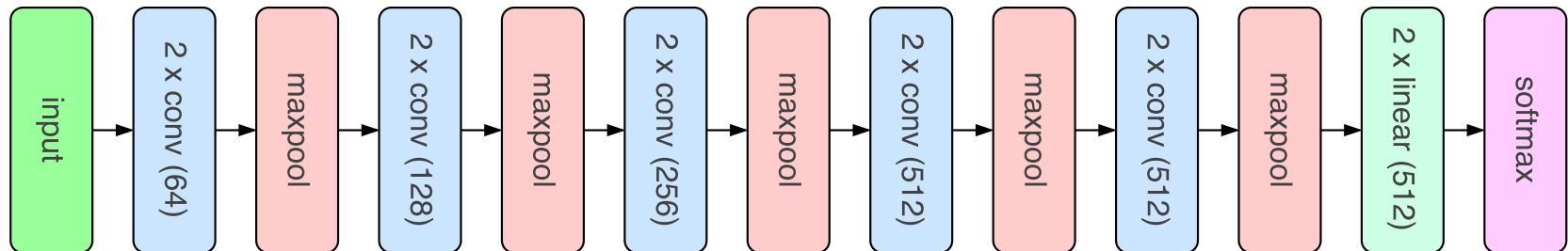- 50k / 10k split
- 75.7% best acc.

increasing difficulty!

# Models used in the experiments
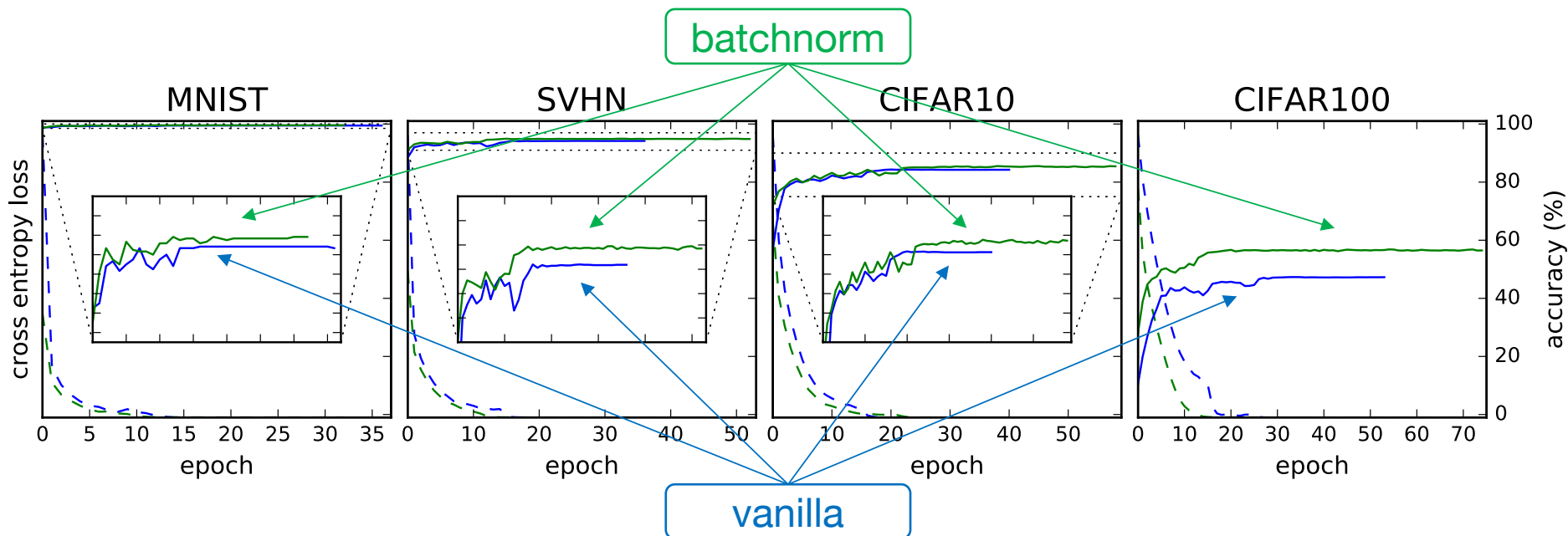
**Multi-layer perceptron:**

toy model for
initial experiments

input → 3 x linear (100) → softmax

we will focus
on this one!

**Deep convolutional neural network:**

input → 2 x conv (64) → maxpool → 2 x conv (128) → maxpool → 2 x conv (256) → maxpool → 2 x conv (512) → maxpool → 2 x conv (512) → maxpool → 2 x linear (512) → softmax

Activation function after each convolution & linear layer, 3x3 kernels and zero padding for convolutions, 2x2 max-pooling
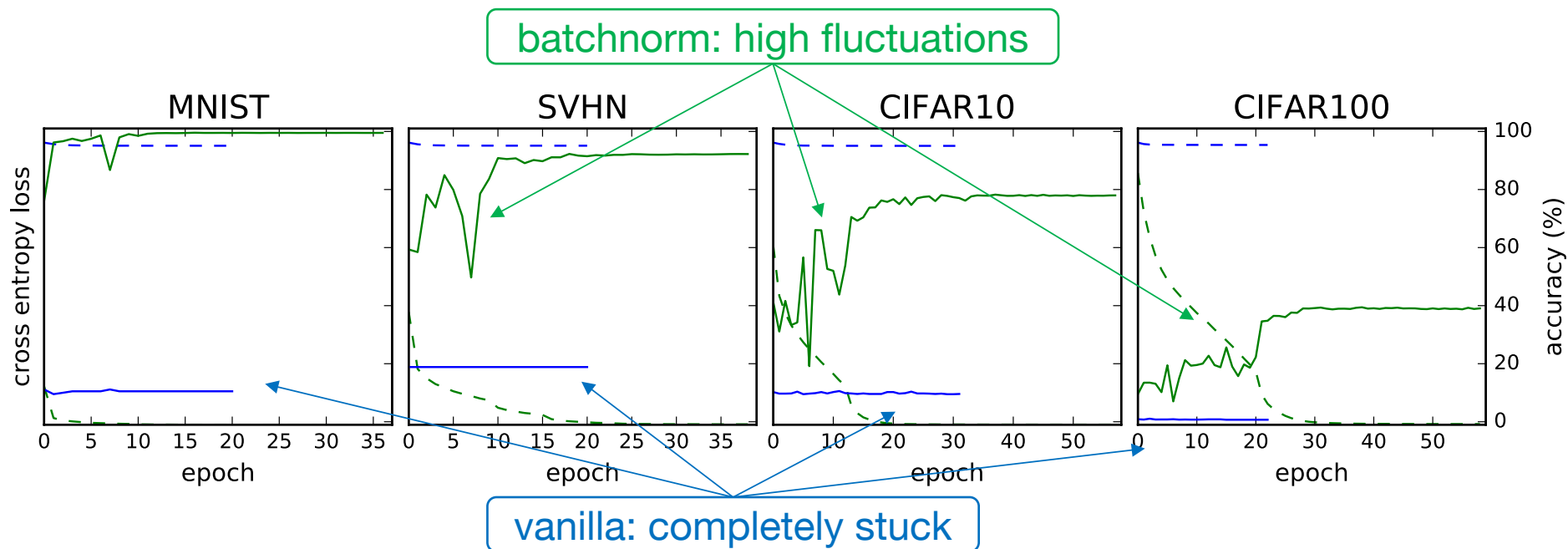
# Experimental results

What did we find out?

# Batchnorm achieves higher accuracies



- **Baseline model** with ReLU activations and Kaiming* initialization
- Batchnorm layer added before all activation functions
- Near-zero losses on all datasets (overfitting)
- Batchnorm beats the vanilla network on all datasets!
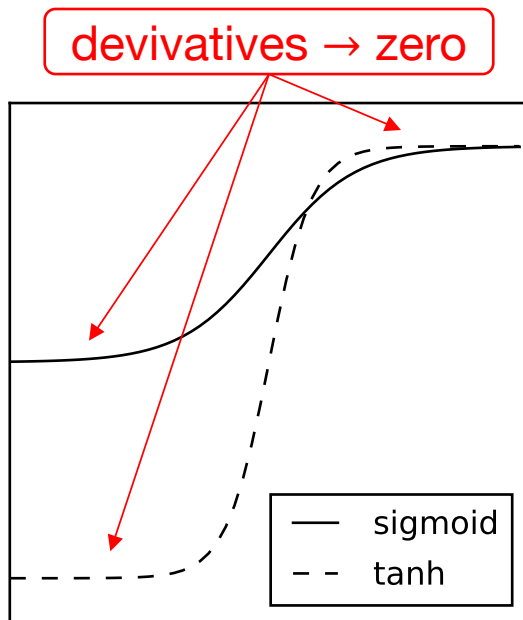- Score: Batchnorm 1 – vanilla 0

* Kaiming He et. al: "Delving Deep into Rectifiers: Surpassing Human-Level Performance on ImageNet Classification" (2015)

# Batchnorm can handle Sigmoids



batchnorm: high fluctuations

MNIST    SVHN    CIFAR10    CIFAR100

cross entropy loss

accuracy (%)

vanilla: completely stuck

- Exchange ReLU nonlinearities with Sigmoids (and Xavier initialization)
- Vanilla network is stuck at random guessing and cannot escape
- Batchnorm model converges but with major fluctuations
- Score: Batchnorm 2 – vanilla 0
- Why does this happen to the vanilla model?

* Xavier Glorot et. al: "Understanding the Difficulty of Training Deep Feedforward Neural Networks" (2010)

# A sidenote on saturating nonlinearities
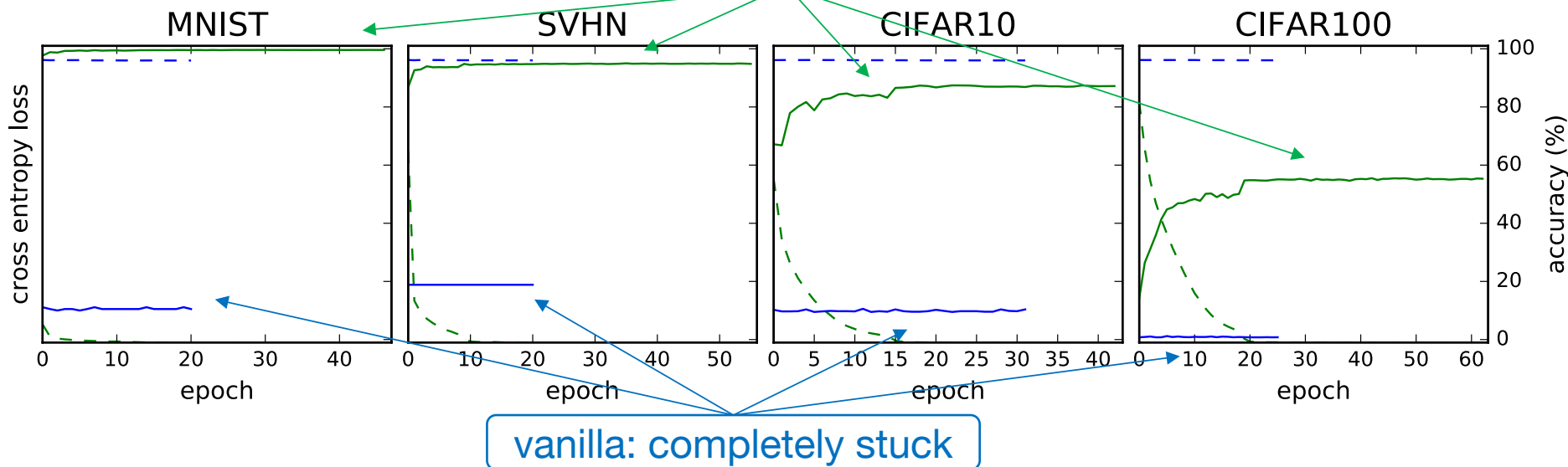
- Sigmoids (and `tanh` also) activation functions *saturate*
- This causes *vanishing gradient problems*\*

devivatives → zero



- Activations stuck in saturated regime
- More severe as network depth increases (multiplicative effect!)
- Batchnorm's scale and shift parameters push activations to non-saturated regime!
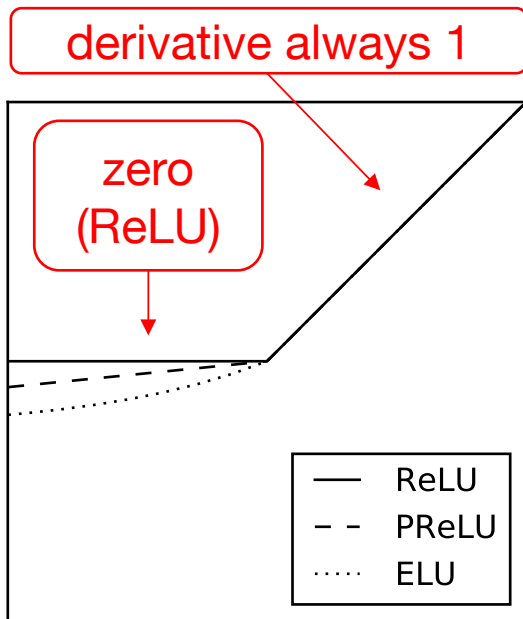
# Batchnorm favors higher learning rates

batchnorm: better accuracy and convergence speed



vanilla: completely stuck

- Back to ReLU baseline, increase (10x) initial learning rate to $\eta = 0.1$
- Batchnorm: **higher accuracy & convergence speed** than baseline
- Vanilla model is stuck at random guessing…
- Score: batchnorm 3 – vanilla 0
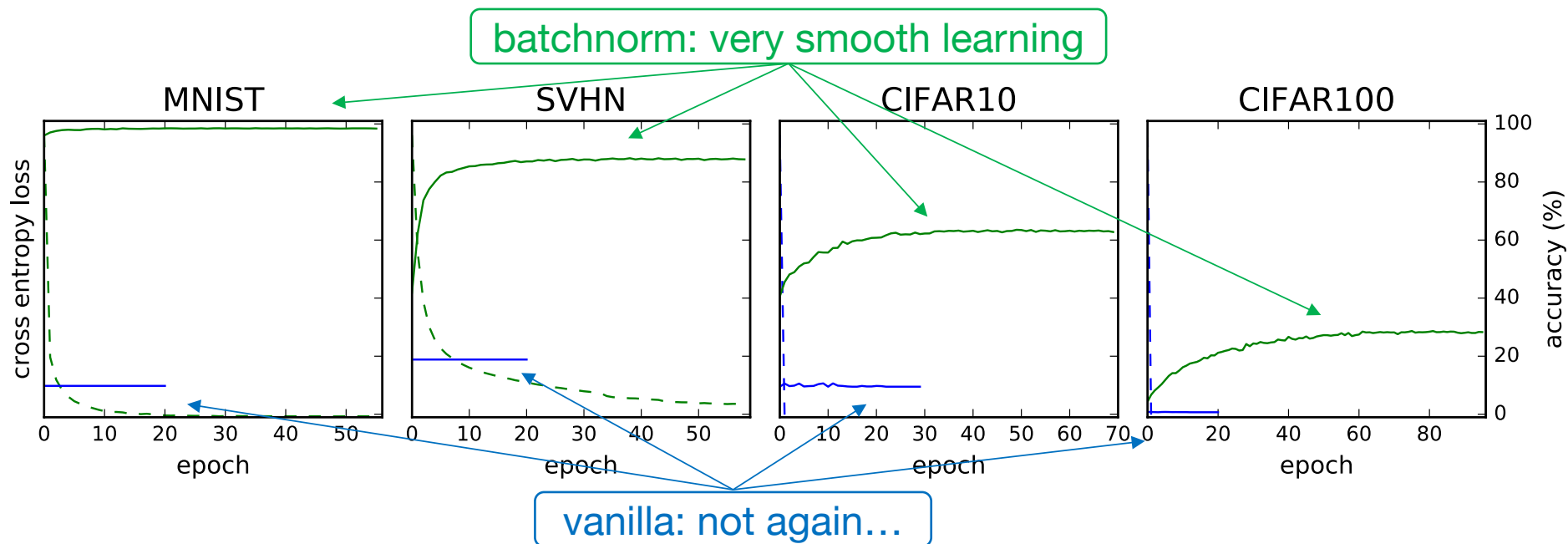- What causes this behavior in the vanilla model?

# A sidenote on rectified nonlinearities

- ReLUs (also ELUs) are not bounded by saturation
- This causes *exploding gradient problems**



derivative always 1

zero (ReLU)

—— ReLU
– – PReLU
..... ELU

- Opposite of vanishing gradients!
- ReLUs cause so-called *dead neurons*
- Batchnorm scales and shifts activations and thus helps to prevents dead neurons
- Another approach: *gradient clipping*

# Batchnorm handles arbitrary inits

batchnorm: very smooth learning



vanilla: not again…

- Back to ReLU baseline, initialize weights with $\mathcal{N}(\mu = 0, \sigma = 1)$
- Batchnorm: smooth learning but **worse accuracy** than baseline init
- Vanilla model is stuck again…
- Score: batchnorm 4 – vanilla 0
- What causes this behavior in the vanilla model?

23

# A sidenote on weight initialization

*Symmetry-breaking* and *variance-preserving* inits crucial!

Batchnorm forward pass:

$$\text{batchnorm}\big((\alpha\boldsymbol{w})^T\boldsymbol{x} + \boldsymbol{b}\big) = \text{batchnorm}(\boldsymbol{w}^T\boldsymbol{x})$$

$\alpha$ has no effect

Batchnorm backward pass:

$$\frac{\partial\,\text{batchnorm}\big((\alpha\boldsymbol{w}^T\boldsymbol{x})\big)}{\partial\,\boldsymbol{x}} = \frac{\partial\,\text{batchnorm}(\boldsymbol{w}^T\boldsymbol{x})}{\partial\,\boldsymbol{x}}$$

$\alpha$ has no effect

$$\frac{\partial\,\text{batchnorm}\big((\alpha\boldsymbol{w}^T\boldsymbol{x})\big)}{\partial\,\alpha\boldsymbol{w}} = \frac{1}{\alpha}\cdot\frac{\partial\,\text{batchnorm}(\boldsymbol{w}^T\boldsymbol{x})}{\partial\,\boldsymbol{w}}$$

larger weights
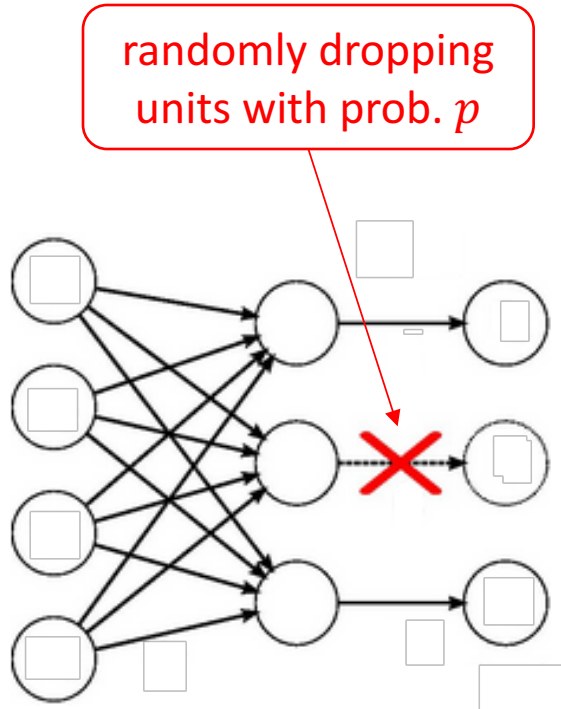→ smaller gradients
and vice versa!

# Batchnorm likes some weight decay

batchnorm: always ahead



- Back to baseline, add $L_2$ weight decay with strength $\lambda = 0.0005$
- Both models achieve **higher accuracies** than baseline models
- Batchnorm model consistently ahead of vanilla counterpart
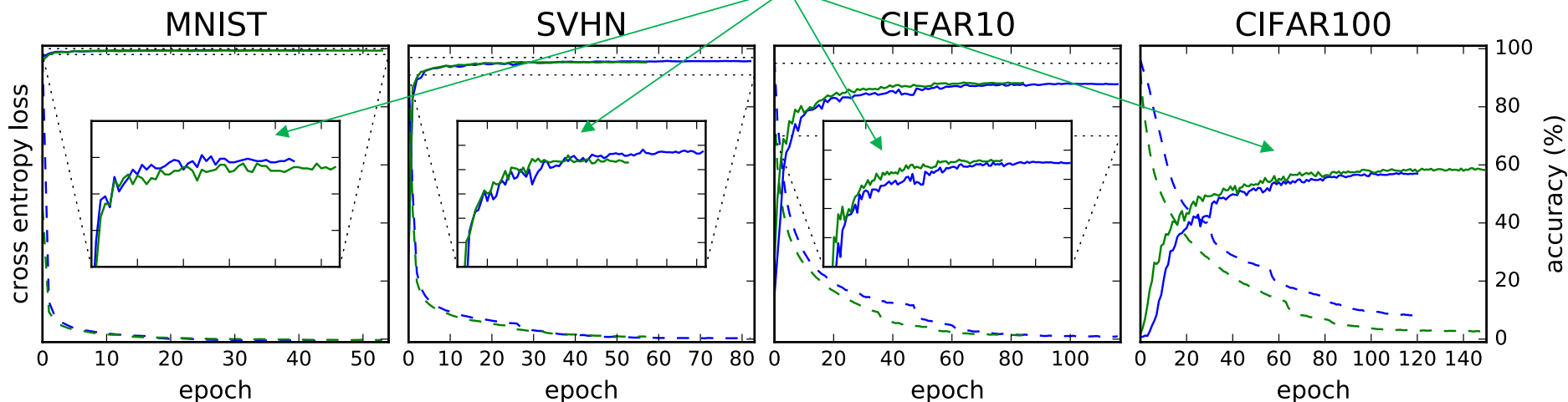- Score: batchnorm 5 – vanilla 0

25

# A quick intro to dropout *

randomly dropping
units with prob. $p$



- Prevents co-adaptation of features

- Can be seen as ensemble learning

- We do **not** drop units at test time!

* Nitish Srivastava et. al: "Dropout: A Simple Way to Prevent Neural Networks from Overfitting" (2014)
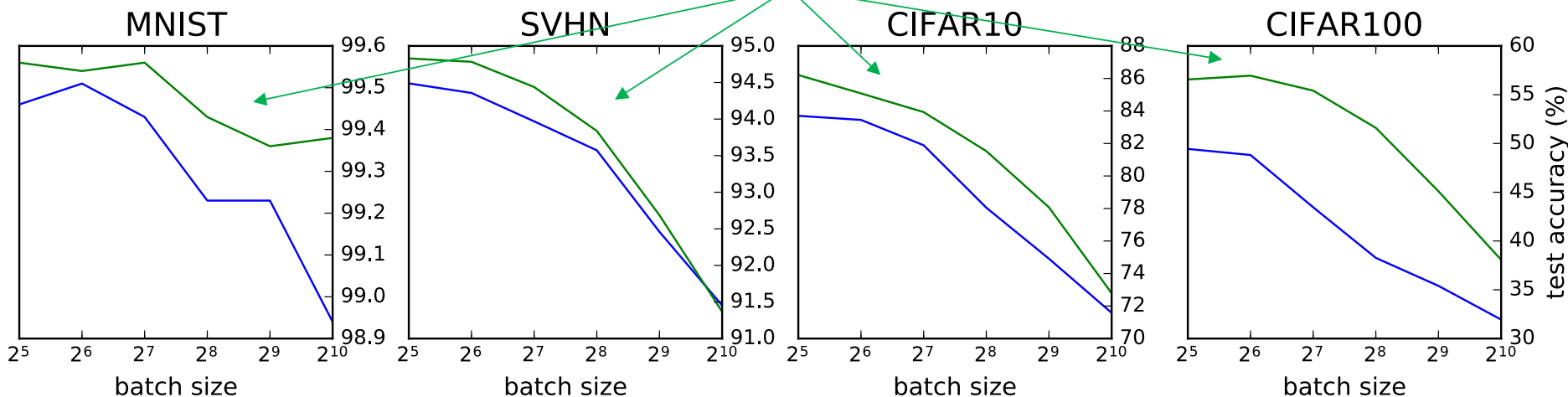
# Batchnorm suffers under dropout



- Back to baseline, add dropout with keep probability $p = 0.5$
- Both models converge **much slower** than before, **best accuracies!**
- Batchnorm not ahead of vanilla counterpart
- Score: batchnorm 5 – vanilla 0 (we'll call this a draw!)
- Why? Batch statistics might be of lower value with dropped activations

# Batchnorm can handle all batch sizes



- Back to baseline, we select batch sizes as in: $\mathcal{B} = \{2^5, 2^6, 2^7, 2^8, 2^9, 2^{10}\}$
- Batchnorm ahead independent of batch size
- General trend: **smaller batch** sizes work better (for our models)
- Ultimately, it is a time/memory tradeoff!
- Score: batchnorm 6 – vanilla 0

# Batchnorm time overhead noticeable

vanilla → MNIST (V)
batchnorm → MNIST (B)

| Batch size | Mean overhead | MNIST | | SVHN | | CIFAR10 | | CIFAR100 | |
|---|---|---|---|---|---|---|---|---|---|
| | | V | B | V | B | V | B | V | B |
| 32 | 19.7% | 141 | 162 | 166 | 201 | 112 | 137 | 113 | 137 |
| 64 | 19.0% | 103 | 122 | 128 | 153 | 87 | 104 | 87 | 104 |
| 128 | 16.4% | 86 | 101 | 107 | 124 | 73 | 85 | 73 | 85 |
| 256 | 15.5% | 77 | 89 | 95 | 109 | 65 | 75 | 65 | 75 |
| 512 | 15.3% | 72 | 83 | 89 | 102 | 60 | 69 | 60 | 70 |
| 1024 | 15.2% | 69 | 80 | 86 | 99 | 58 | 66 | 58 | 66 |

Table B.4: CNN average time per epoch for different batch sizes in seconds.

- Batchnorm scales nicely with increasing batch sizes
- Score: batchnorm 6 – vanilla 1 (we have to admit defeat here)
- **Your mileage may vary! ***

* Using CUDA 7.5 & cuDNN v4. Logging introduced overhead (but same for vanilla & batchnorm on absolute scale)

# Conclusion: the effect of batchnorm

- What general effects does batchnorm have on…
  - …classification accuracy? Higher accuracies!
  - …convergence speed? Slightly faster! *

- What specific effects does batchnorm have on…
  - …activation functions? Can handle them all!
  - …learning rates? Higher learning rates favorable.
  - …weight initialization? Variance-preserving still best.
  - …regularization methods? Weight decay helps, not dropout. *
  - …batch sizes? Can handle all of them, smaller better. *
  - …wall time? Noticeable overhead, but worth it. *

* Very much in general. There are a lot of "but if…"'s

# Conclusion: why does batchnorm work?

- Examples are normalized over different batches which leads to a regularizing effect

- Vanishing and exploding gradients can be avoided due to scale and shift parameters

- More stable gradient propagation in general *

- Layers can "focus" on the predictive relationship, rather than adapting to shift in layer distributions

# Key takeaway

Use batchnorm if your model has convergence problems!
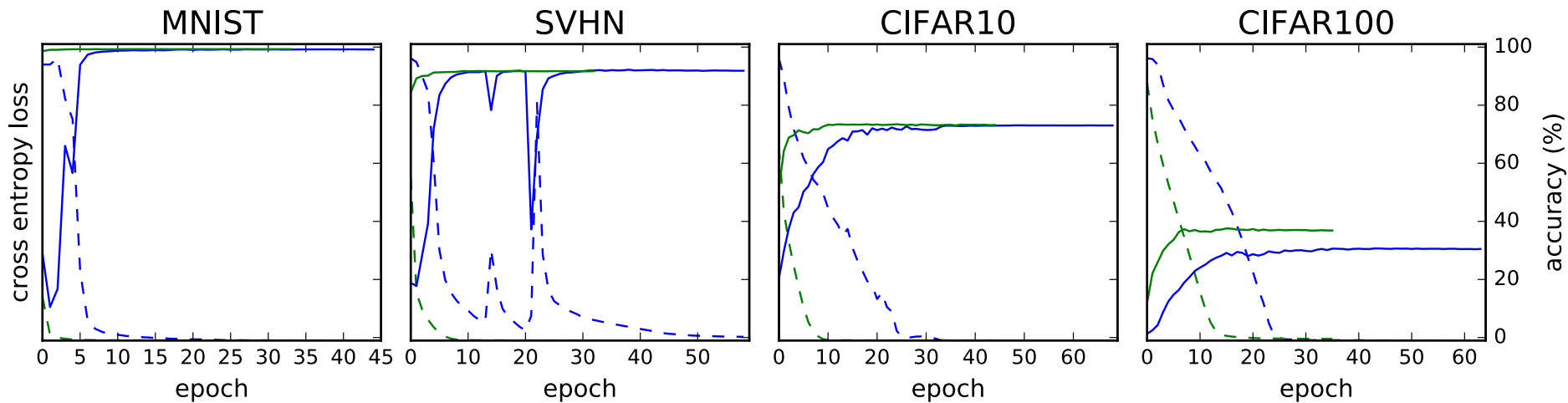
## Questions?

# Backup slides

Did not have time to cover this

# Training strategy for the experiments

- Global feature standardization

- Shuffle training set, 10% validation split

- Cross entropy loss

- Mini-batch SGD with batch size $\mathcal{B} = 64$
  - **Nesterov momentum** * with $\mu = 0.9$

- Initial learning rate of $\eta = 0.01$

- Adaptive learning rate decay with $k = 0.5$
  - After 5 epochs without improvement

- "Early" stopping after 20 epochs without improvement

# Funky convergence problems

# An intro to $L_2$ weight decay

add $L_2$ norm to loss

$$\mathcal{L}(\hat{\boldsymbol{y}}, \boldsymbol{y}, \boldsymbol{w}) = \frac{1}{\mathcal{B}} \sum_{i=1}^{\mathcal{B}} \mathcal{L}_i(\hat{y}_i, y_i) + \frac{\lambda}{2} ||w||_2^2$$

- - Data loss

data loss minimal

total loss

regularization loss minimal

- Gaussian prior on weights
- Perform MAP instead of MLE
- Independent of batchnorm?

$w_2$

$w_1$

—— $L_1$
—— $L_2$
—— Elastic net

# SGD: what could possibly go wrong?

We have the following simple model:

note: effects up to order $l$!

$$f(x, \boldsymbol{w}) = x \cdot w_1 \cdot w_2 \cdot ... \cdot w_l = \hat{y}$$

Assume $\frac{\partial \mathcal{L}}{\partial \hat{y}} = 1$ so we wish to decrease $\hat{y}$ slighly.

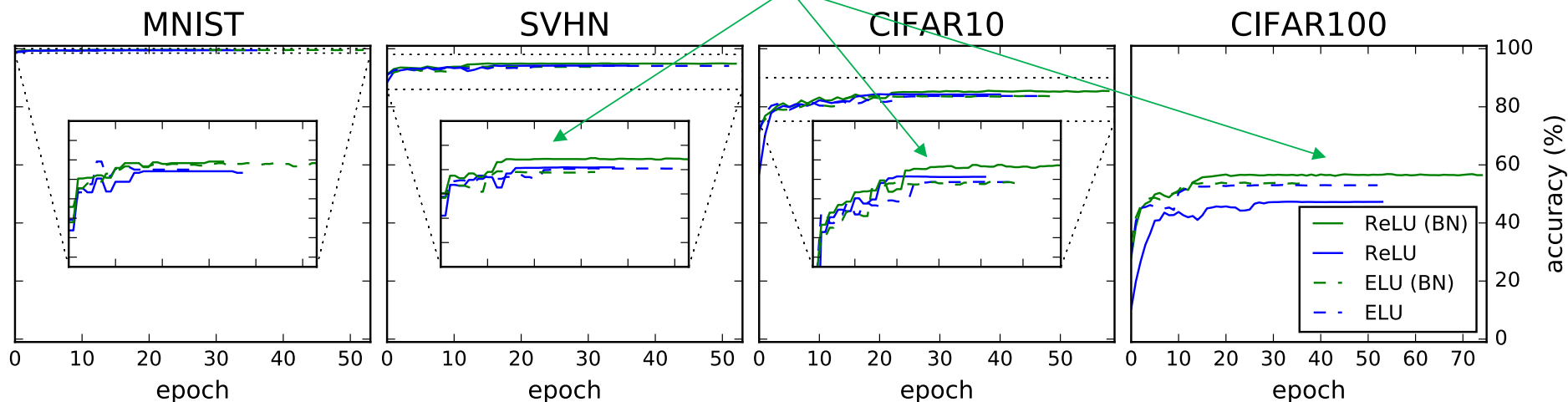After gradient descent update the network computes:

$$f(x, \boldsymbol{w}) = x \cdot \left( w_1 - \eta \frac{\partial \mathcal{L}}{\partial w_1} \right) \cdot \left( w_2 - \eta \frac{\partial \mathcal{L}}{\partial w_2} \right) \cdot ... \cdot \left( w_l - \eta \frac{\partial \mathcal{L}}{\partial w_l} \right) = \hat{y}$$

adaptive optimization can help!

how can wechoose the learning rate?

37

# Batchnorm + ReLU outperforms ELU*



- Train both models with ELU nonlinearities
- Paper claims that vanilla ELUs outperform ReLU models with batchnorm
- We found that this is **not** the case (at least for our setup)
- Theoretical properties of ELUs are questionable

* D. Clevert et. al: "Fast and Accurate Deep Network Learning by Exponential Linear Units (ELUs)" (2015)

# Conclusion: practical recommendations

- Initial learning rate & its decay are crucial
  - If you optimize one thing, optimize those!
- Use ReLU (give parametric & leaky ReLU a try)
- Use dropout!
- Try residual connections and inception
- Use batchnorm if convergence seems problematic