CS11-747 Neural Networks for NLP

# Debugging Neural Networks for NLP

Graham Neubig

**Carnegie Mellon University**

**Language Technologies Institute**

Site
https://phontron.com/class/nn4nlp2021/

# In Neural Networks, Debugging is Paramount!

- Models are often **complicated and opaque**

- **Everything is a hyperparameter** (network size, model variations, batch size/strategy, optimizer/ learning rate)

- Non-convex, stochastic optimization has **no guarantee of decreasing/converging loss**

# Understanding Your Problem

# A Typical Situation

- You've implemented a nice model

- You've looked at the code, and it looks OK

- Your accuracy on the test set is bad

- **What do I do?**

# Possible Causes

- **Training time problems**
  - Lack of model capacity
  - Inability to train model properly
  - Training time bug
- **Decoding time bugs**
  - Disconnect between test and decoding
  - Failure of search algorithm
- **Overfitting**
- **Mismatch between optimized function and eval**

Don't debug all at once! Start top and work down.

# Debugging at Training Time
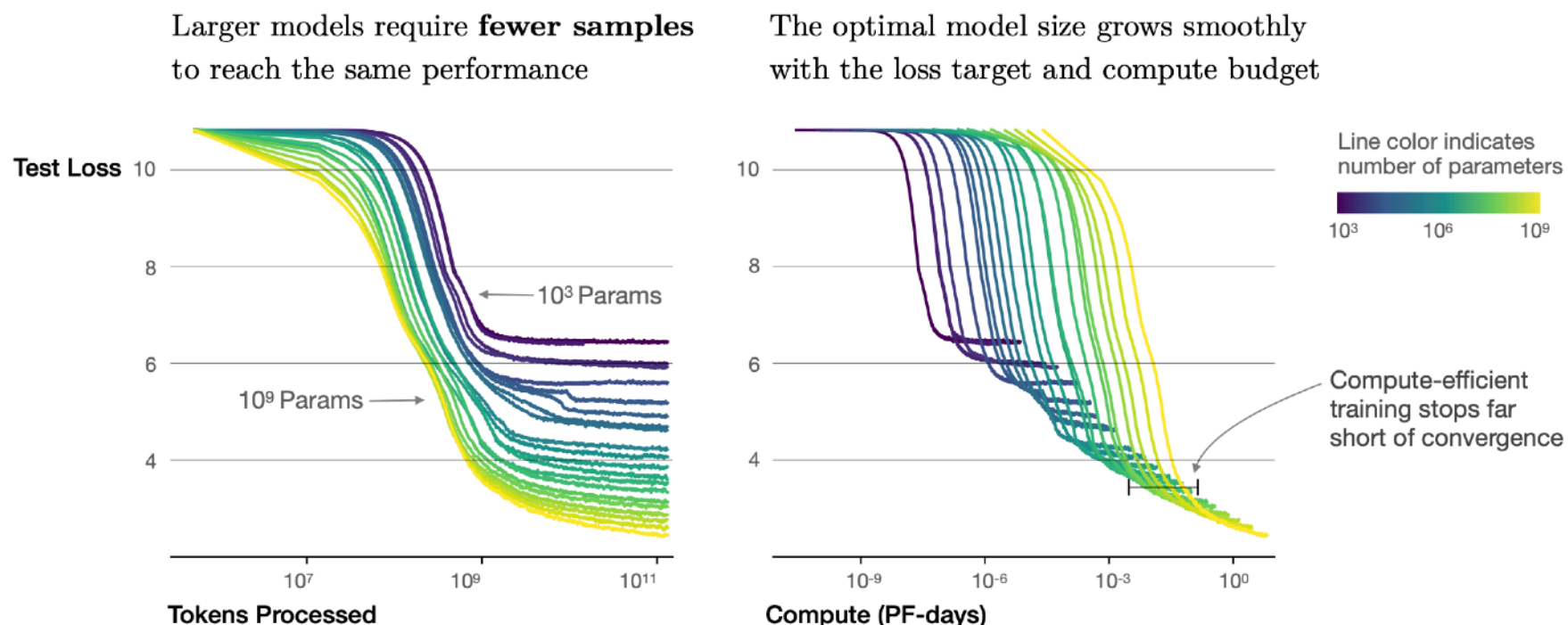
# Identifying Training Time Problems

- Look at the **loss function** calculated on the **training set**

  - Is the loss function going down?

  - Is it going down basically to zero if you run training long enough (e.g. 20-30 epochs)?

  - If not, does it go down to zero if you use very small datasets?

# Is My Model Too Weak?

- Larger models tend to perform better, esp. when pre-trained (e.g. Raffel et al. 2020)

| Model | GLUE Average | CoLA Matthew's | SST-2 Accuracy | MRPC F1 | MRPC Accuracy | STS-B Pearson | STS-B Spearman |
|---|---|---|---|---|---|---|---|
| Previous best | $89.4^a$ | $69.2^b$ | $97.1^a$ | $\mathbf{93.6}^b$ | $\mathbf{91.5}^b$ | $92.7^b$ | $92.3^b$ |
| T5-Small | 77.4 | 41.0 | 91.8 | 89.7 | 86.6 | 85.6 | 85.0 |
| T5-Base | 82.7 | 51.1 | 95.2 | 90.7 | 87.5 | 89.4 | 88.6 |
| T5-Large | 86.4 | 61.2 | 96.3 | 92.4 | 89.9 | 89.9 | 89.2 |
| T5-3B | 88.5 | 67.1 | 97.4 | 92.5 | 90.0 | 90.6 | 89.8 |
| T5-11B | **90.3** | **71.6** | **97.5** | 92.8 | 90.4 | **93.1** | **92.8** |

- Larger models can learn with fewer steps (Kaplan et al. 2020, Li et al. 2020)

# Trouble w/ Optimization

- If increasing model size doesn't help, you may have an optimization problem

- **Possible causes:**

  - Bad optimizer

  - Bad learning rate

  - Bad initialization

  - Bad minibatching strategy

# Reminder: Optimizers

- **SGD:** take a step in the direction of the gradient

- **SGD with Momentum:** Remember gradients from past time steps to prevent sudden changes

- **Adagrad:** Adapt the learning rate to reduce learning rate for frequently updated parameters (as measured by the variance of the gradient)

- **Adam:** Like Adagrad, but keeps a running average of momentum and gradient variance

- **Many others:** RMSProp, Adadelta, etc.
(See Ruder 2016 reference for more details)

# Learning Rate

- Learning rate is an important parameter

  - Too low: will not learn or learn vey slowly

  - Too high: will learn for a while, then fluctuate and diverge

- **Common strategy:** start from an initial learning rate then gradually decrease

- **Note:** need a different learning rate for each optimizer! (SGD default is 0.1, Adam 0.001)
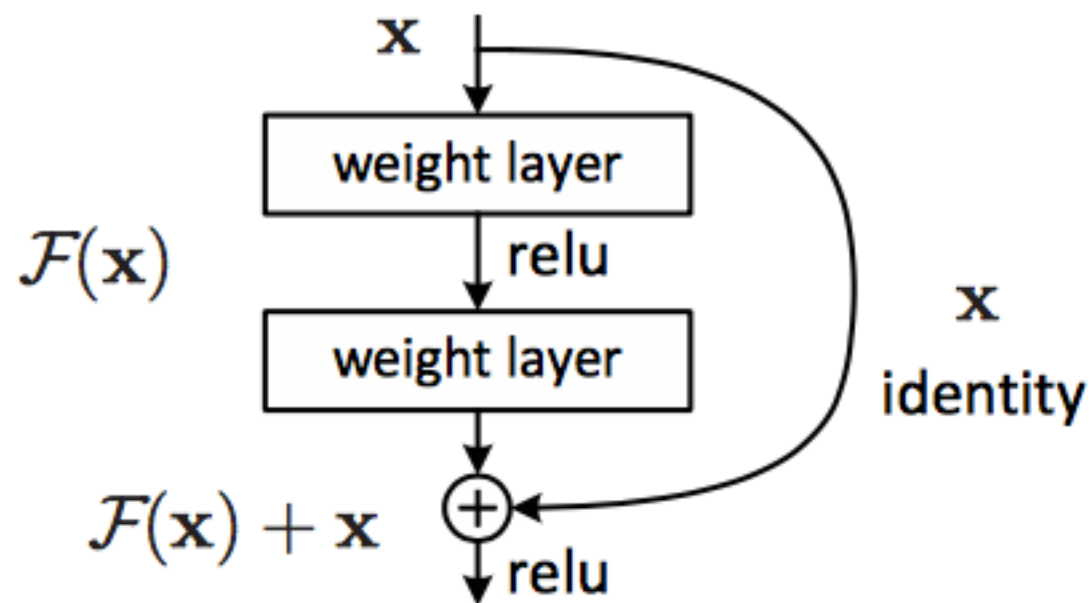
# Initialization

- Neural nets are sensitive to initialization, which results in different sized gradients

- Standard initialization methods:

  - **Gaussian initialization:** initialize with a zero-mean Gaussian distribution

  - **Uniform range initialization:** simply initialize uniformly within a range

  - **Glorot initialization, He initialization:** initialize in a uniform manner, where the range is specified according to net size

- Latter is common/default, but read prior work carefully

# Be Careful of Multi-layer Models

- Extra layers can help, but can also hurt if you're not careful due to vanishing gradients

- Solutions:

**Residual Connections (He et al. 2015)**     **Highway Networks (Srivastava et al. 2015)**



$$\mathbf{y} = H(\mathbf{x}, \mathbf{W_H}) \cdot T(\mathbf{x}, \mathbf{W_T}) + \mathbf{x} \cdot (1 - T(\mathbf{x}, \mathbf{W_T}))$$

# Debugging at Test Time

# Training/Decoding Disconnects

- Usually your loss calculation and prediction will be implemented in different functions

- Especially true for structured prediction models (e.g. encoder-decoders)

  - See `enc_dec.py` example from this class, which has `calc_loss()` and `generate()` functions

- Like all software engineering: **duplicated code is a source of bugs**!

- Also, usually loss calculation is minibatched, generation not.

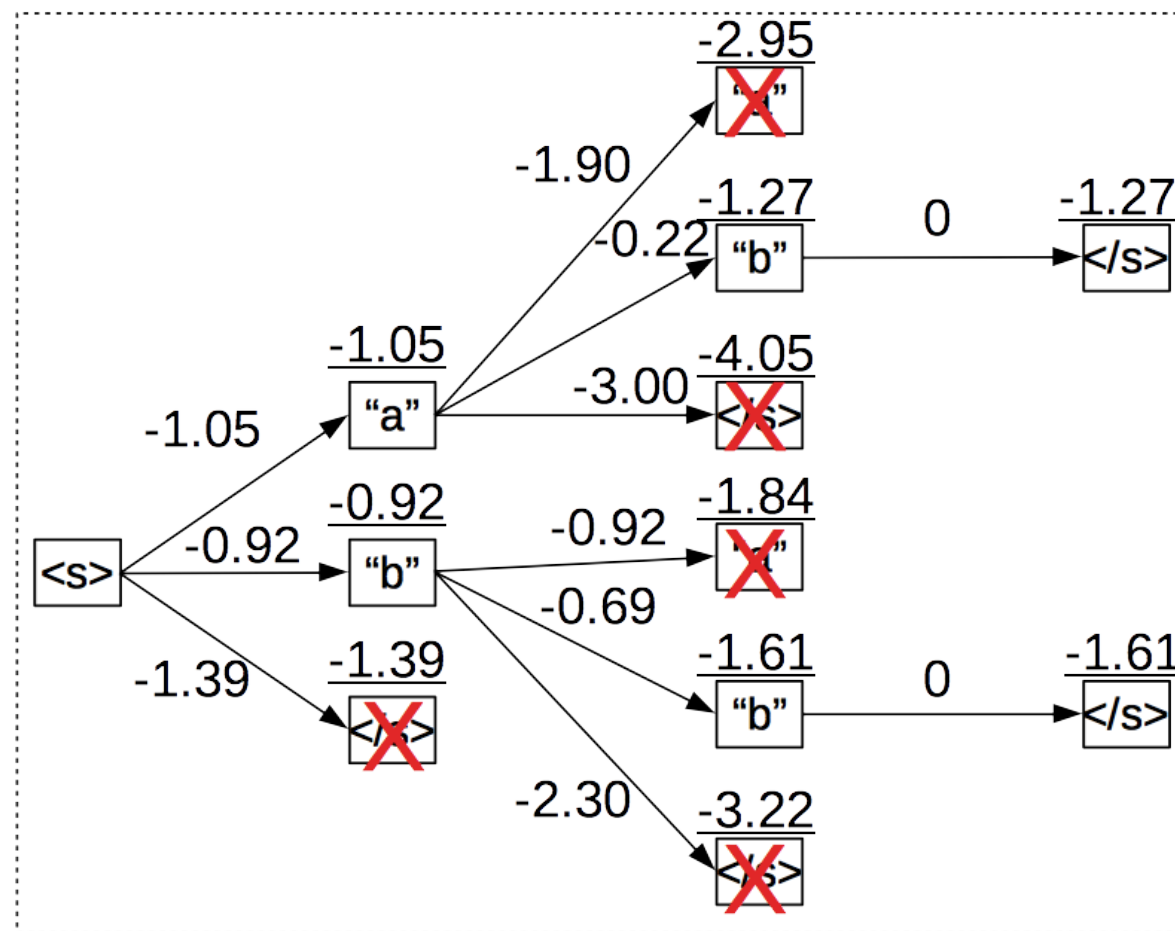# Debugging Minibatching

- Debugging mini-batched loss calculation

  - Calculate loss with **large batch size** (e.g. 32)

  - Calculate loss for **each sentence individually and sum**

  - The values should be the same (modulo numerical precision)

- Create a unit test that tests this!

# Debugging Structured Generation

- Your decoding code should get the same score as loss calculation

- Test this:

    - Call **decoding function**, to generate an output, and keep track of its score

    - Call **loss function** on the generated output

    - The score of the two functions should be the same

- Create a unit test doing this!

# Beam Search

- Instead of picking one high-probability word, maintain several paths



- More in a later class

# Debugging Search

- As you make search better, the model score should get better (almost all the time)

- Run search with varying beam sizes and make sure you get a better overall model score with larger sizes

- Create a unit test testing this!

# Look At Your Data!

- Decoding problems can often be detected by looking at outputs and realizing something is wrong

- e.g. The first word of the sentence is dropped every time
  > went to the store yesterday
  > bought a dog

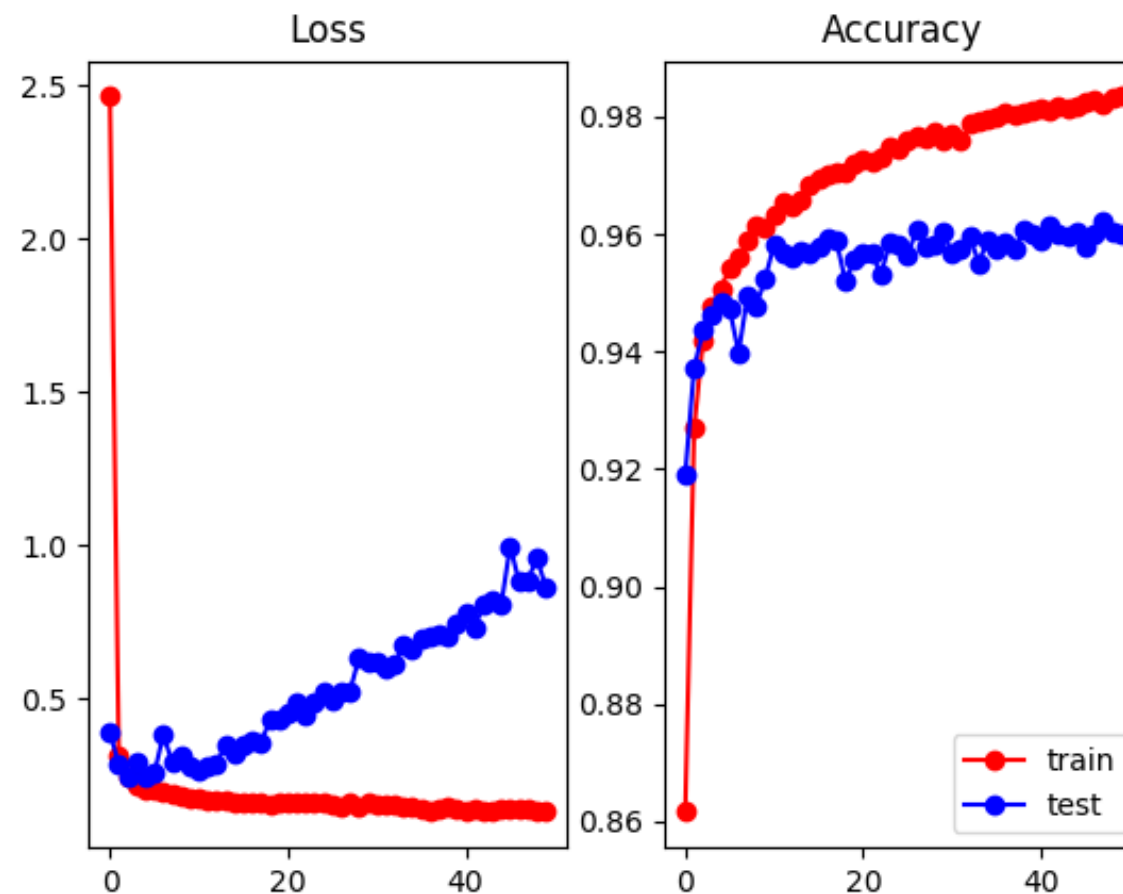- e.g. our system was <unk>ing University of Nebraska at Kearney

# Mismatch b/t Optimized Function and Evaluation Metric

# Loss Function, Evaluation Metric

- It is very common to optimize for maximum likelihood for training

- But even though likelihood is getting better, accuracy can get worse

# Example w/ Classification
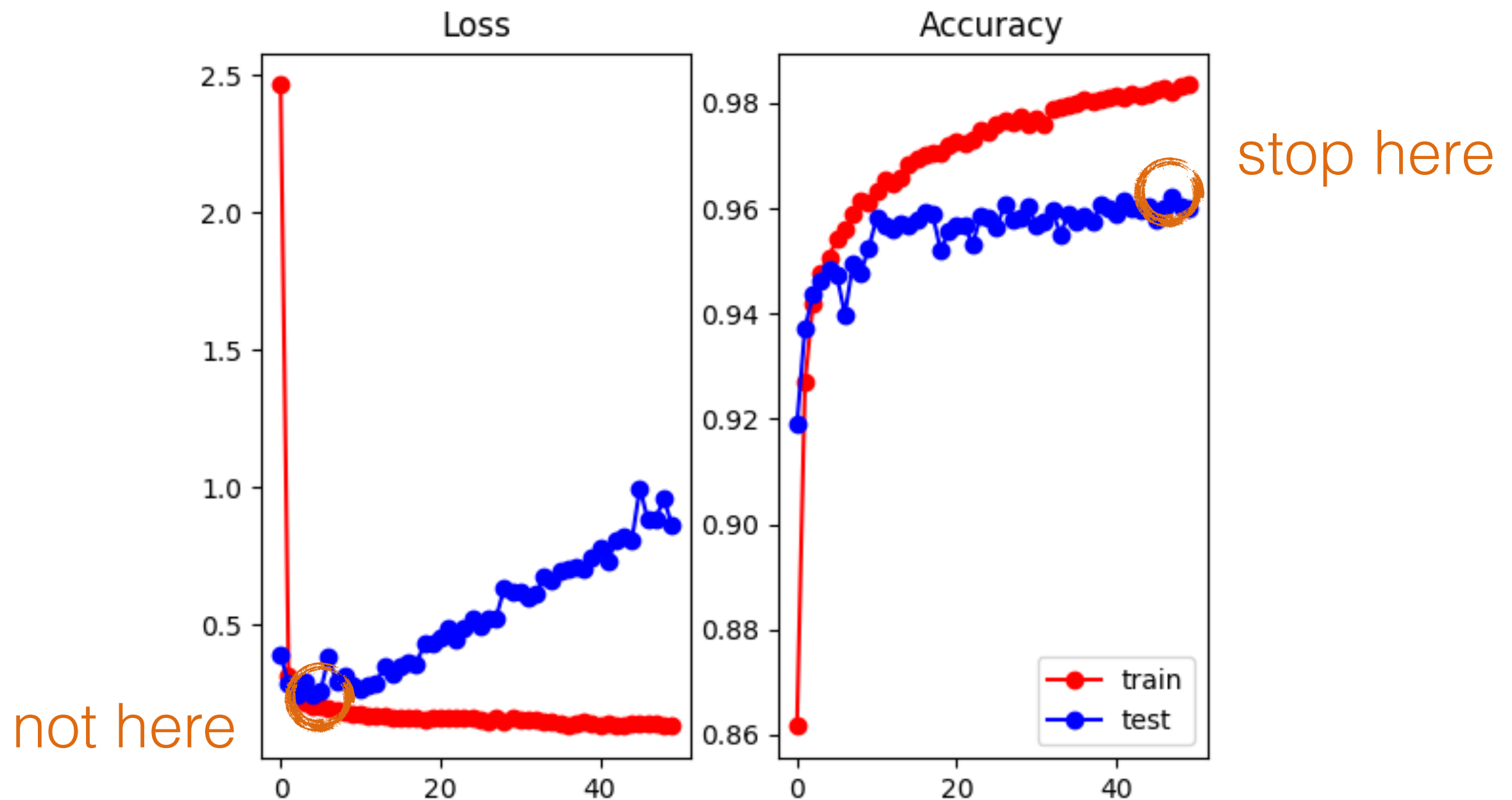
- Loss and accuracy are de-correlated (see dev)



- Why? Model gets more confident about its mistakes.

# Managing Loss Function/ Eval Metric Differences

- Most principled way: use structured prediction techniques to be discussed in future classes

  - Structured max-margin training

  - Minimum risk training

  - Reinforcement learning

  - Reward augmented maximum likelihood

# A Simple Method:
# Early Stopping w/ Eval Metric

# Final Words

# Reproducing Previous Work

- Reproducing previous work is hard because everything is a hyper-parameter

- If code is released, find and reduce the differences one by one

- If code is not released, try your best

- Feel free to contact authors about details, they will usually respond!

# Questions?