

Applied Deep Learning

DAT310

Lecture:
Deep Learning Best
Practices



UNIVERSITY OF
CALGARY

Start
something.

Agenda



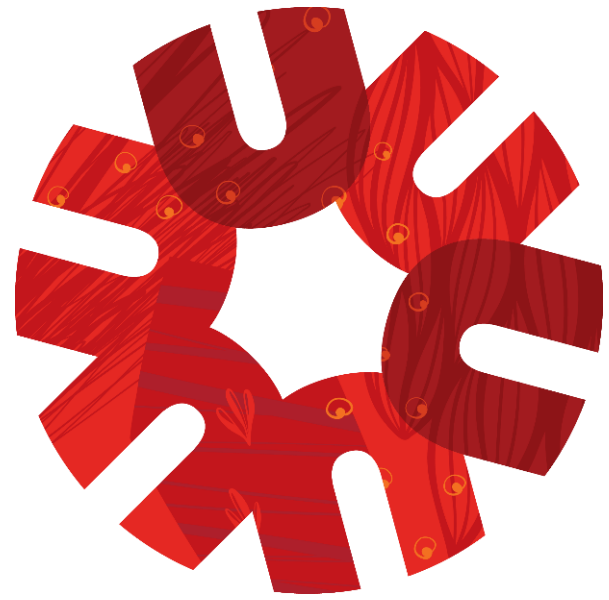
UNIVERSITY OF
CALGARY

**Start
something.**

- General ML Pipeline
- Rules of Thumb
- Debugging Networks
- Improving Results
- Faster Networks



General ML Pipeline



UNIVERSITY OF CALGARY
CONTINUING EDUCATION

**Start
something.**

General ML Pipeline



UNIVERSITY OF
CALGARY

**Start
something.**

1. Get a dataset
2. Build a dataset and dataloader
3. Choose some models for the task
4. Train your models
 - a. Need optimizer and loss function
 - b. Compute loss w/ output of network, compute gradients using backwards
 - c. Apply gradients with optimizer.step()
5. Evaluate your models



Rules of Thumb

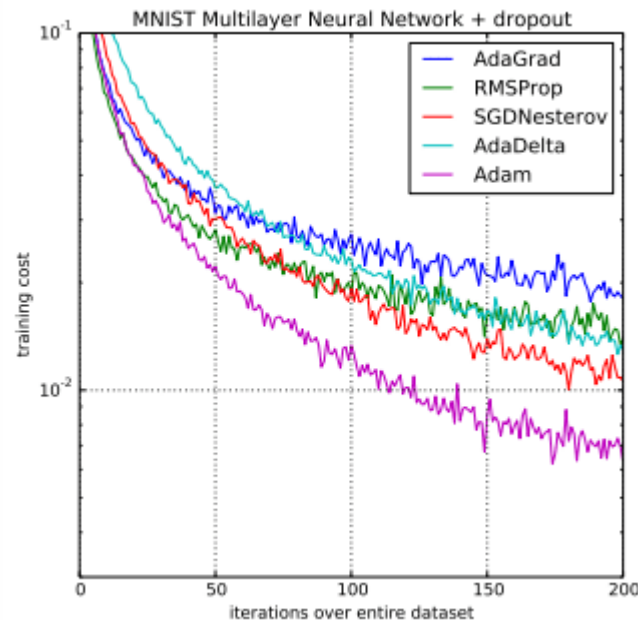


UNIVERSITY OF CALGARY
CONTINUING EDUCATION

**Start
something.**

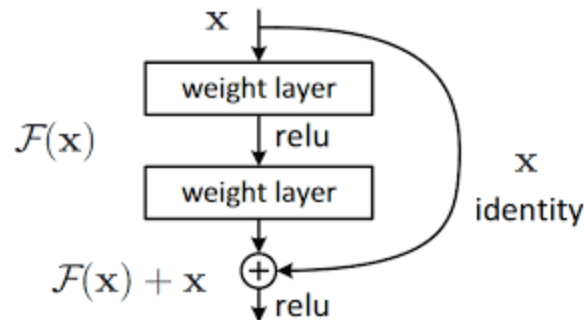
Optimizer

- Pick Adam as your optimizer
 - Combo of Momentum and RMS Prop
 - Automatically figures out a learning rate for you!



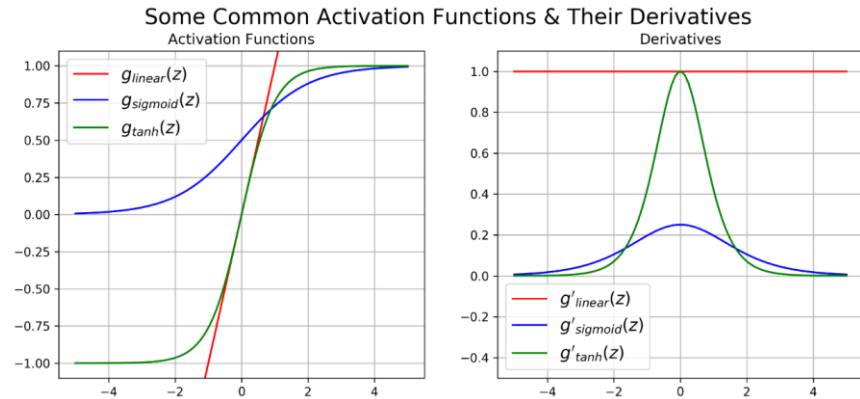
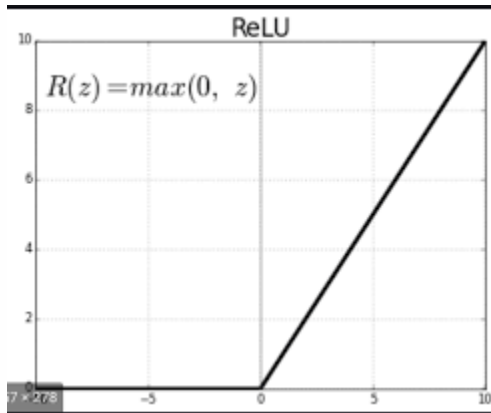
Non Linearities

- Deeper VS Wider -> Deeper is better
 - Why? More non linearities
- How to make deeper?
 - Skip Connections



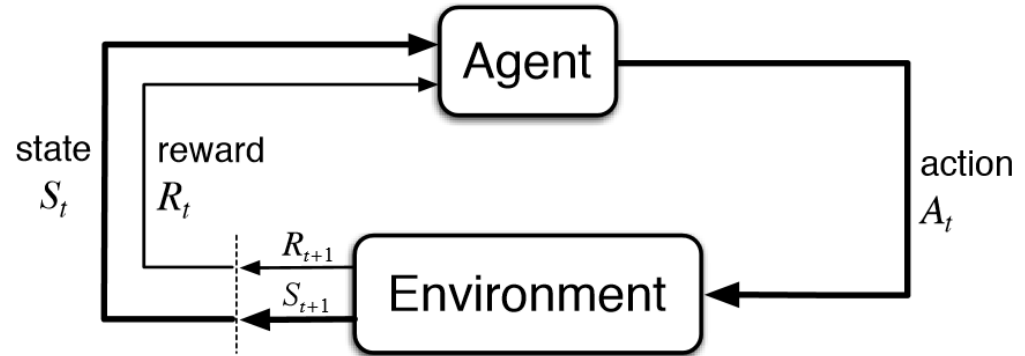
Non Linearities

- Pick Relu
 - No exploding or vanishing gradients like tanh or sigmoid



Reinforcement Learning

- When do you TYPICALLY use reinforcement learning?
 - No loss function with a gradient available
- You have an environment that produces a state
 - I.e a video game



Rules of Thumb



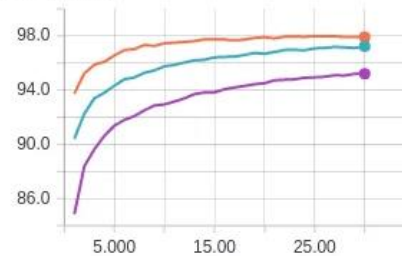
UNIVERSITY OF
CALGARY

Start
something.

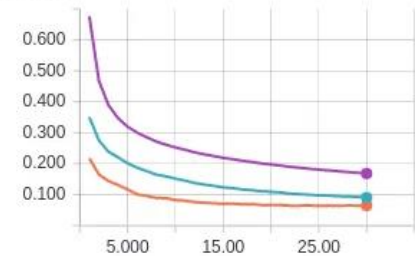
Batch Size

- Bigger = Faster training, worse generalization
- Smaller = Slower training, better generalization

test accuracy



test loss



Testing loss and accuracy when the model is trained using different batch sizes.

- Purple = Batch Size 1024
- Blue = 256
- Orange = 64



Rules of Thumb

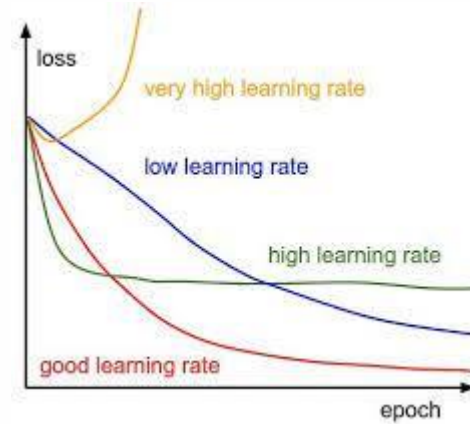


UNIVERSITY OF
CALGARY

Start
something.

Learning Rate

- Plot your loss per LR



Rules of Thumb

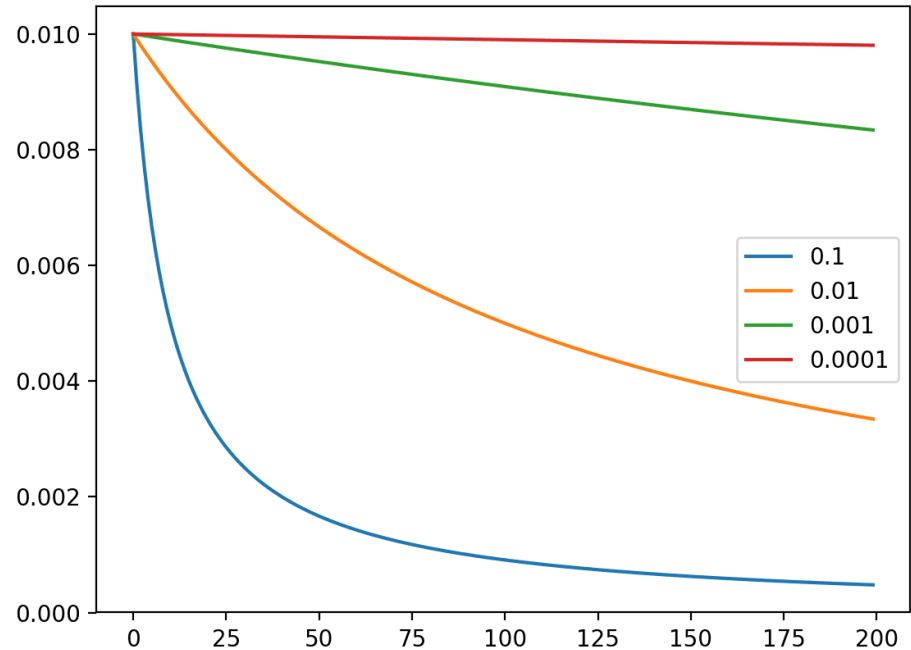


UNIVERSITY OF
CALGARY

Start
something.

Learning Rate

- Plot your loss per LR



Rules of Thumb



UNIVERSITY OF
CALGARY

**Start
something.**

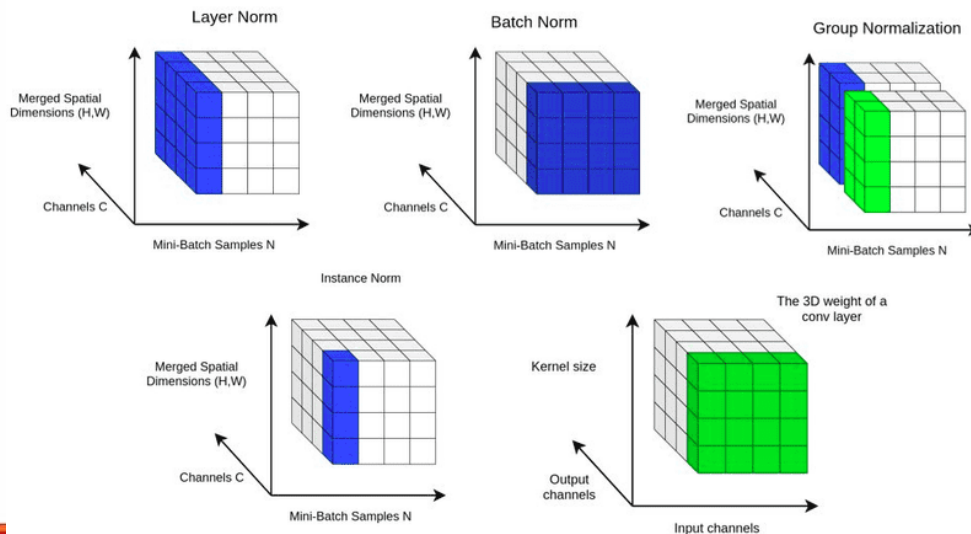
Transfer Learning

- Always use pretrained models!
 - ImageNet or COCO for CV
 - Language Models for NLP



BatchNorm

- Always use BatchNorm (or some variation)!
 - Trains 6X faster
 - Less sensitive to LR
 - Helps w/ vanishing and exploding gradients



Model Size VS Dataset Size

- Model Size
 - If you are positive your implementation of the network and training is correct, but you still cannot overfit to just a small subset of training data.... Make your model deeper + wider
- Dataset Size
 - As long as your training error is low, you can always decrease generalization error by collecting more training data



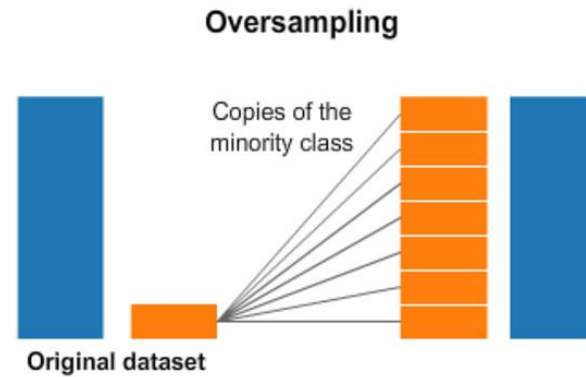
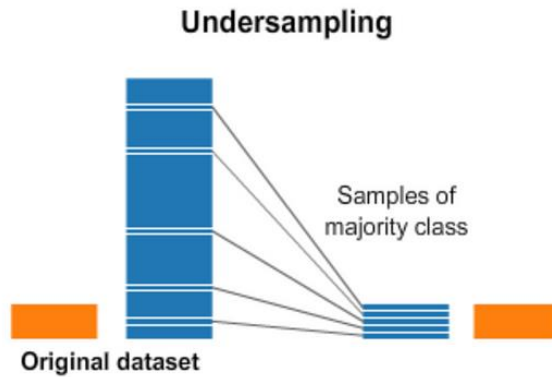
Output Activation

- Binary Classification:
 - Sigmoid
- Multi Class:
 - Softmax
- Regression:
 - Sigmoid
 - None



Imbalanced Data

- Oversample
- Undersample
- Weight the loss
- Use balanced accuracy and a confusion matrix to evaluate with



Rules of Thumb



UNIVERSITY OF
CALGARY

**Start
something.**

Regularization

- L2 is less forgiving, set lambda small i.e 0.0001
- L2 is used by default in pytorch



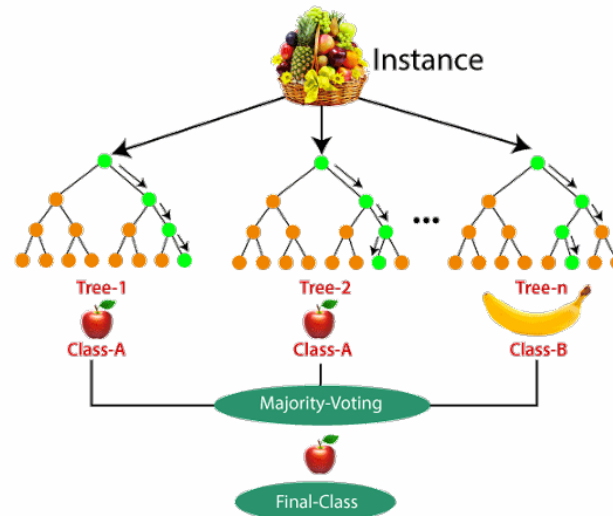
Gradient Clipping

- I loathe gradient clipping
- Prevents exploding grads
- Should fix exploding grads before using this
- Prevents gradients from being over a certain size
 - I find this cripples your learning



Don't forget about Classic ML Methods

- Linear Regression, Logistic Regression, Random Forest, Gradient Boosting... all excellent models that are fast to train
- IME they do very well on tabular data, awful on text and images



Debugging Networks



UNIVERSITY OF CALGARY
CONTINUING EDUCATION

**Start
something.**

Debugging Networks



UNIVERSITY OF
CALGARY

**Start
something.**

1. Check dataloader
2. Overfit on a small subset of data
3. Play with LR
4. Are your parameters registered?
5. Check gradients



Debugging Networks



UNIVERSITY OF
CALGARY

**Start
something.**

Check Dataloader

- For CV insure your images are loaded properly
- Insure you're labels are aligned
- If you apply augmentations to your images, you MIGHT need to apply them to the labels
 - During segmentation for example



Debugging Networks



UNIVERSITY OF
CALGARY

Start
something.

Overfit on a Small Subset

- Your network should be able to memorize a small subset of data
- This is MUCH faster than using the entire dataset
- Allows you to check that the network works, LR is correct, not too much regularization, etc....
QUICKLY



Debugging Networks



UNIVERSITY OF
CALGARY

**Start
something.**

Play With LR

- Covered the plotting a lot already
- Make the plots!



Check Parameters are Registered

- Print out the parameters of your model
- Make sure all layers are there
- For parameter in model.parameters():
 print(parameter)



Debugging Networks

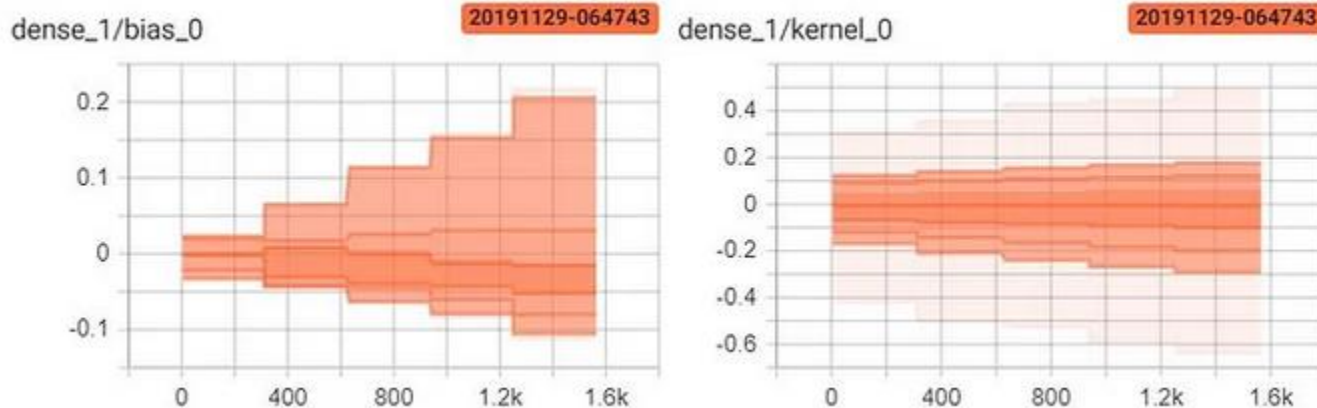


UNIVERSITY OF
CALGARY

Start
something.

Check Gradients

- Can use tensorboardX to viz gradients
- Can check if exploding or vanishing
- Can check if proper size throughout the network



Improving Performance



UNIVERSITY OF CALGARY
CONTINUING EDUCATION

**Start
something.**

Hyperparameter Search

The learning rate is the most important hyperparameter. If bound by time, focus on tuning it.

- The learning rate can be picked by monitoring learning curves that plot the objective function over time.
- The optimal learning rate is typically higher than the learning rate that yields the best performance after the first ~100 iterations, but not so high that it causes instability..



Improving Performance

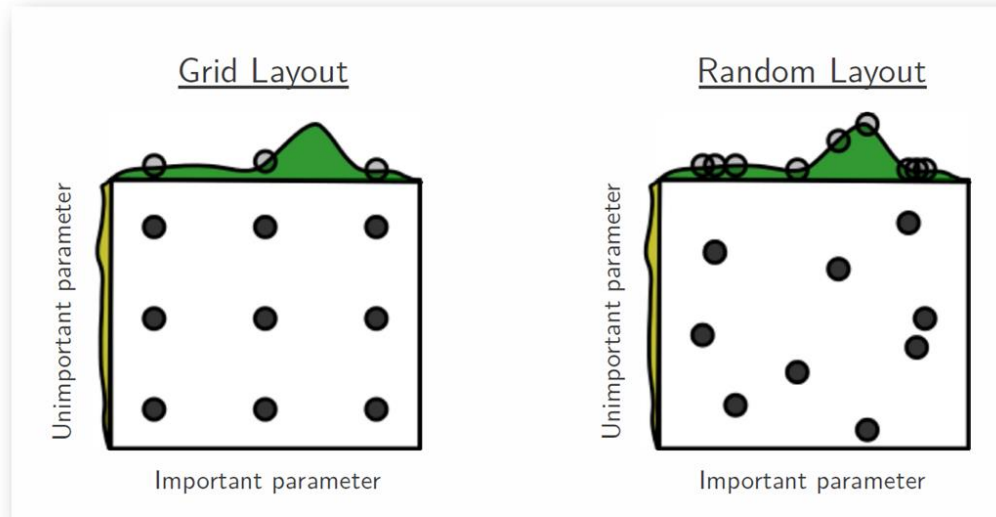


UNIVERSITY OF
CALGARY

Start
something.

Hyperparameter Search

Random search typically converges to good hyperparameters faster than grid search.



Improving Performance



UNIVERSITY OF
CALGARY

**Start
something.**

Check Examples w/ Largest Loss Values

- Save losses on validation or training
- Check the examples with the largest loss values
- Is there a pattern?
- If so, look for papers that address this, or get more data labelled!



Improving Performance



UNIVERSITY OF
CALGARY

**Start
something.**

Early Stopping

- Compute validation loss + accuracy every epoch
- Save the best epoch, even if it's not the last one
- That's your model!



Improving Performance

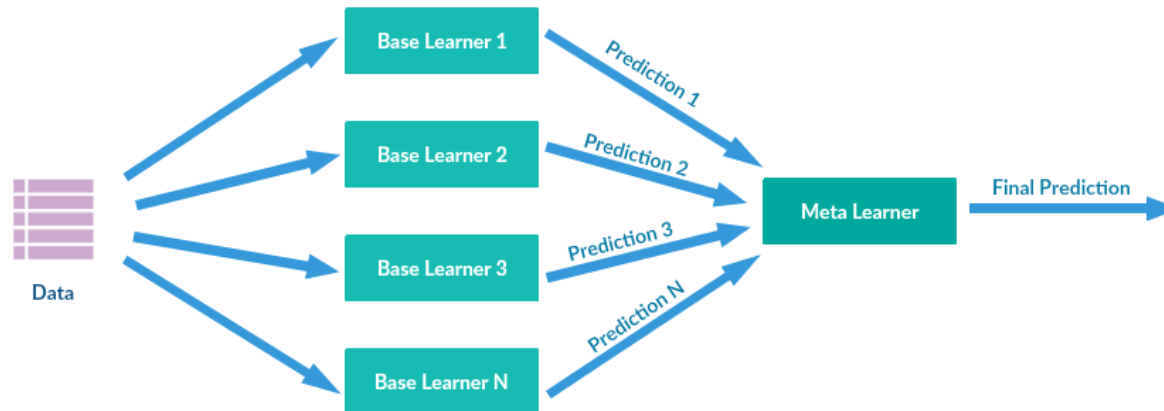


UNIVERSITY OF
CALGARY

Start
something.

Ensembling

- Train multiple models and combine their predictions
 - Usually through averaging the output probabilities
- Great for increasing performance, bad in actual production
 - Too expensive to run!



Improving Performance



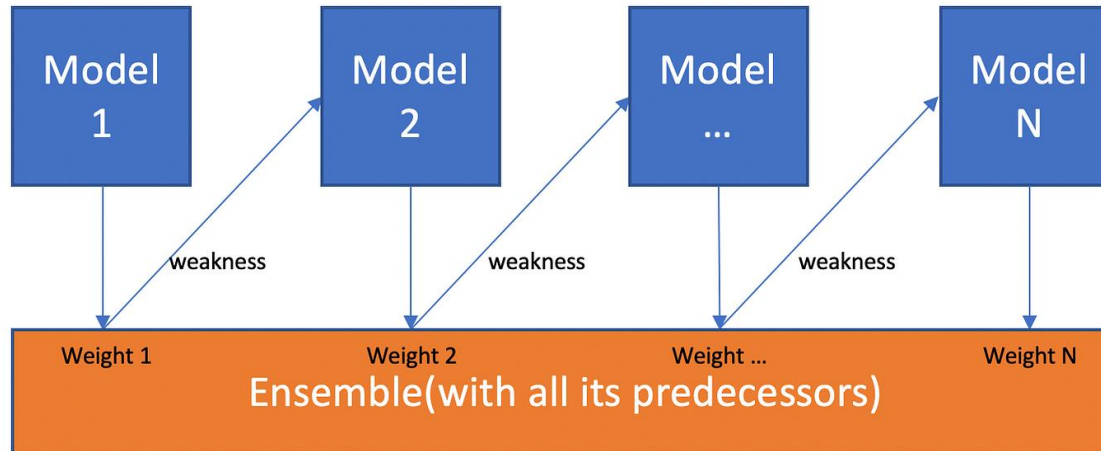
UNIVERSITY OF
CALGARY

Start
something.

Boosting

- Boosting = Training a model on the errors of your previous model
- Then ensembling the two models
- Fancier version of ensembling

Model 1,2,..., N are individual models (e.g. decision tree)



Improving Performance



UNIVERSITY OF
CALGARY

Start
something.

Increase Resolution

- Higher input resolution almost ALWAYS gives better results
- Can test / validate with a higher input resolution than trained with
 - Form of regularization



Improving Performance

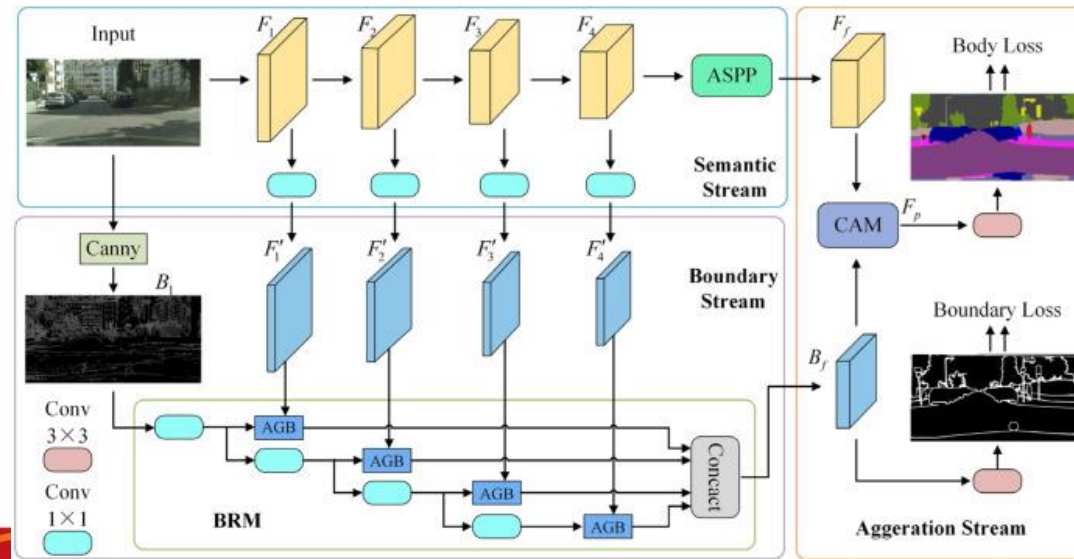


UNIVERSITY OF
CALGARY

Start
something.

Auxiliary Losses

- Use a loss that addresses the weakness of your current model
 - Boundary loss + Segmentation loss for example



Improving Performance

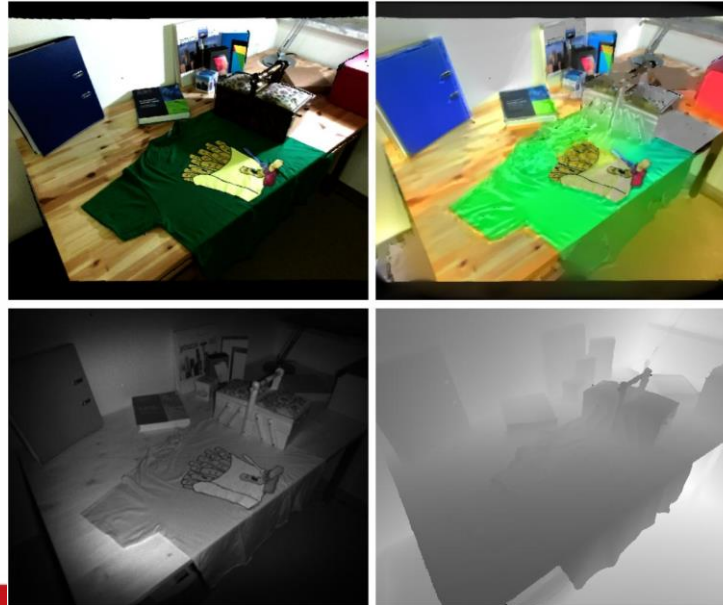


UNIVERSITY OF
CALGARY

Start
something.

Don't just use RGB Input

- Can use RGB + Depth + Edges + a Segmentation etc....
- Sometimes adding this extra channel improves performance



Improving Performance

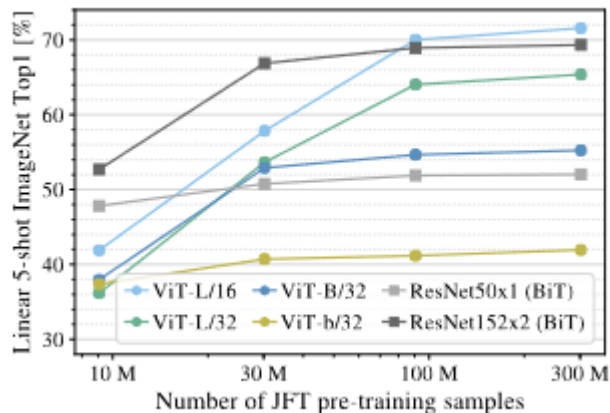


UNIVERSITY OF
CALGARY

Start
something.

When Transformers Excel

- Need AT LEAST 10 MILLION training examples
- Not much better than CNNs until over 100 MILLION training examples



Faster Networks



UNIVERSITY OF CALGARY
CONTINUING EDUCATION

**Start
something.**

Pin Memory

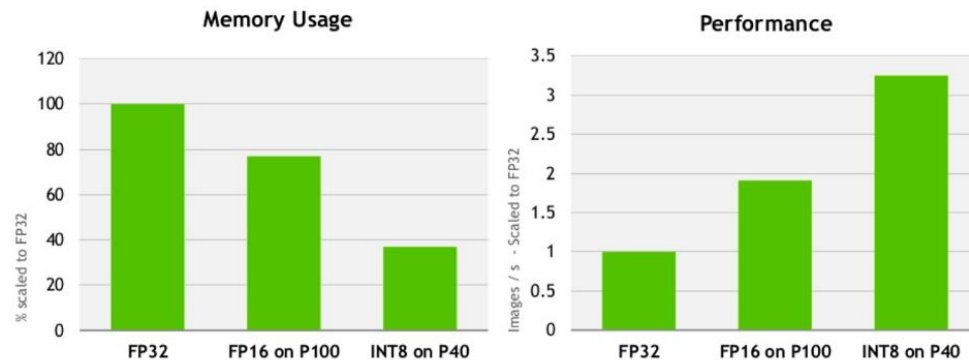
- `Dataloader(dataset, pin_memory=True)`
- Pin_memory transfers to the GPU faster

Setting for the training DataLoader	Time for one training epoch
<code>{'num_workers': 0, 'pin_memory': False}</code>	8.2 s
<code>{'num_workers': 1, 'pin_memory': False}</code>	6.75 s
<code>{'num_workers': 1, 'pin_memory': True}</code>	6.7 s
<code>{'num_workers': 2, 'pin_memory': True}</code>	4.2 s
<code>{'num_workers': 4, 'pin_memory': False}</code>	4.5 s
<code>{'num_workers': 4, 'pin_memory': True}</code>	4.1 s
<code>{'num_workers': 8, 'pin_memory': True}</code>	4.5 s

Half Precision

- By default Pytorch uses float 32
- Half precision is faster, and uses half as much memory
 - `model.half()`
- Does take a minor performance (accuracy) hit

SMALLER AND FASTER



developer.nvidia.com/tensorrt

ResNet50 Model, Batch Size = 128, TensorRT 2.1 RC prerelease

19  NVIDIA

Zero Grads Properly

```
model.zero_grad()
```

```
# or
```

```
optimizer.zero_grad()
```

- executes **memset** for every parameter in the model
- backward pass updates gradients with **"+="** operator (read + write)

```
for param in model.parameters():  
    param.grad = None
```

```
# or (in PyT >= 1.7)
```

```
model.zero_grad(set_to_none=True)
```

- doesn't execute **memset** for every parameter
- memory is zeroed-out by the allocator in a more efficient way
- backward pass updates gradients with **"="** operator (write)

Faster Networks

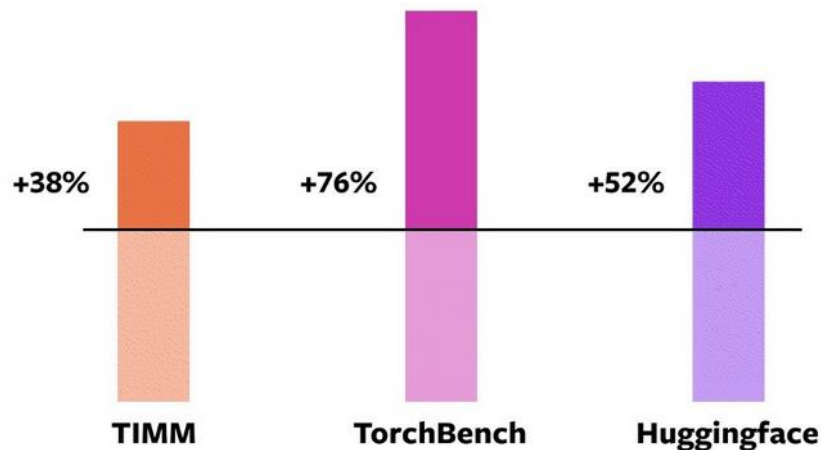


UNIVERSITY OF
CALGARY

Start
something.

Pytorch 2.0

- Pytorch 2.0 allows you to compile your model
- `Model = torch.compile(model)`
- Provides a massive speed up



Speedups for torch.compile against eager mode on an NVIDIA A100 GPU