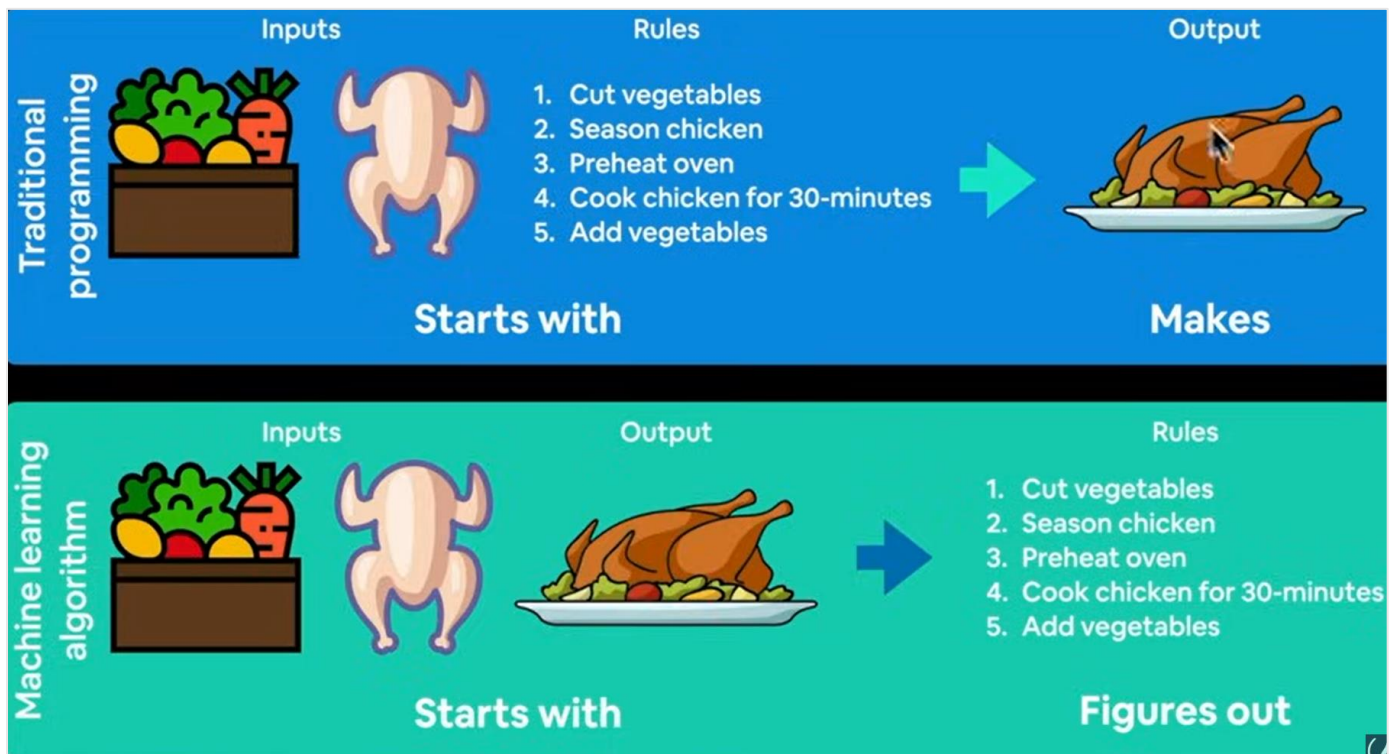


# PyTorch for Deep Learning & Machine Learning

## Machine Learning Algorithms vs Traditional Programming:



### Why use machine learning and deep learning?

**Reason:** For a complex problem, can you think of all the rules? (Probably not)

“If you can build a simple rule-based system that doesn't require machine learning, do that.”

(Google machine learning handbook)

### What deep learning is good for?

- **Problems with long lists of rules:** When the traditional approach fails, machine learning/deep learning may help.
- **Continually changing environments:** Deep learning can adapt (learn) to new scenarios.
- **Discovering insights within large collections of data:** Can you imagine trying to hand-craft rules for what 101 different kinds of food look like?

### What deep learning is not good for?

- **When you need explainability:** The patterns learned by a deep learning model are typically uninterpretable by a human.
- **When the traditional approach is a better option:** If you can accomplish what you need with a simple rule-based system.

- **When errors are unacceptable:** Since the outputs of deep learning model aren't always predictable.
  - **When you don't have much data:** Deep learning models usually require a fairly large amount of data to produce great results.
- 

## Machine learning vs Deep learning:

# Machine Learning vs. Deep Learning

(common algorithms)

- Random forest
- Gradient boosted models
- Naive Bayes
- Nearest neighbour
- Support vector machine
- ...many more

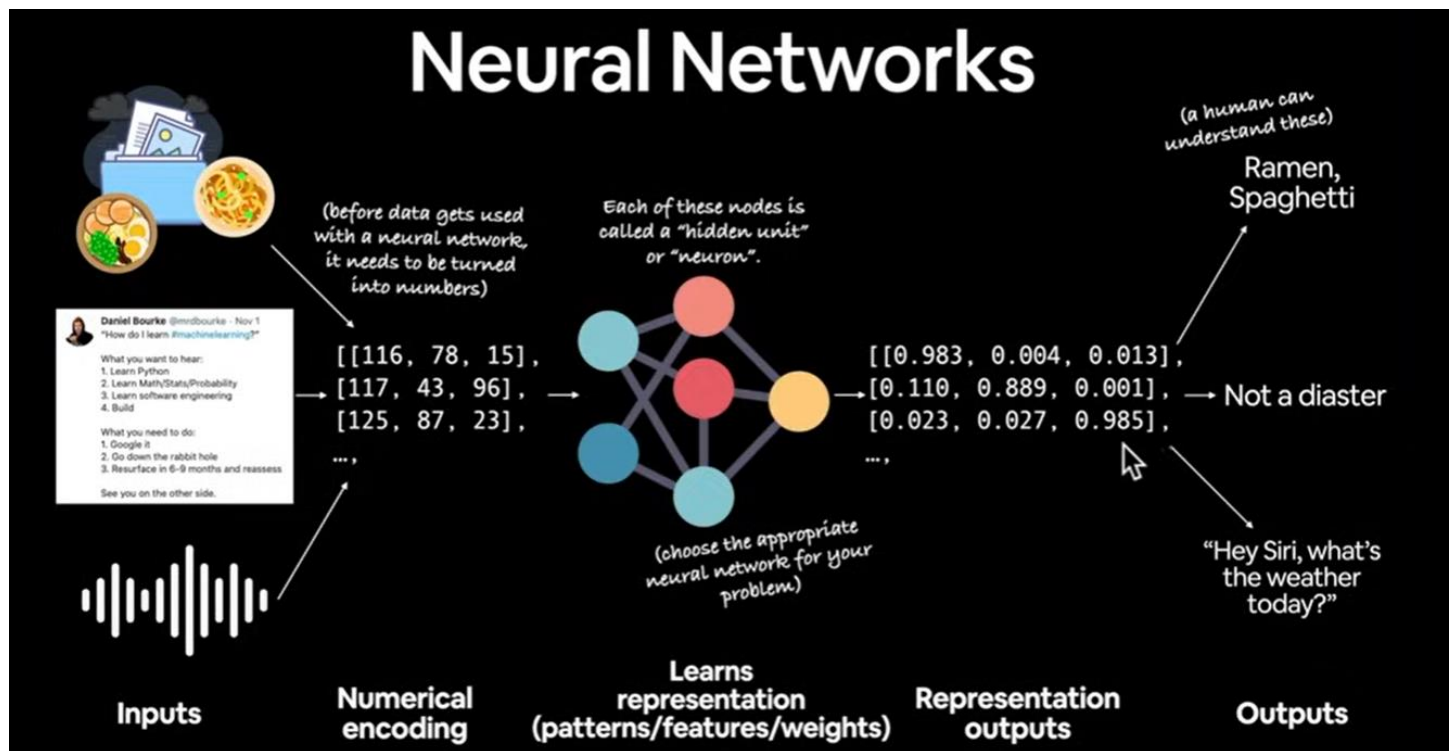
*(since the advent of deep learning these are often referred to as "shallow algorithms")*

- Neural networks
- Fully connected neural network
- Convolutional neural network
- Recurrent neural network
- Transformer
- ...many more

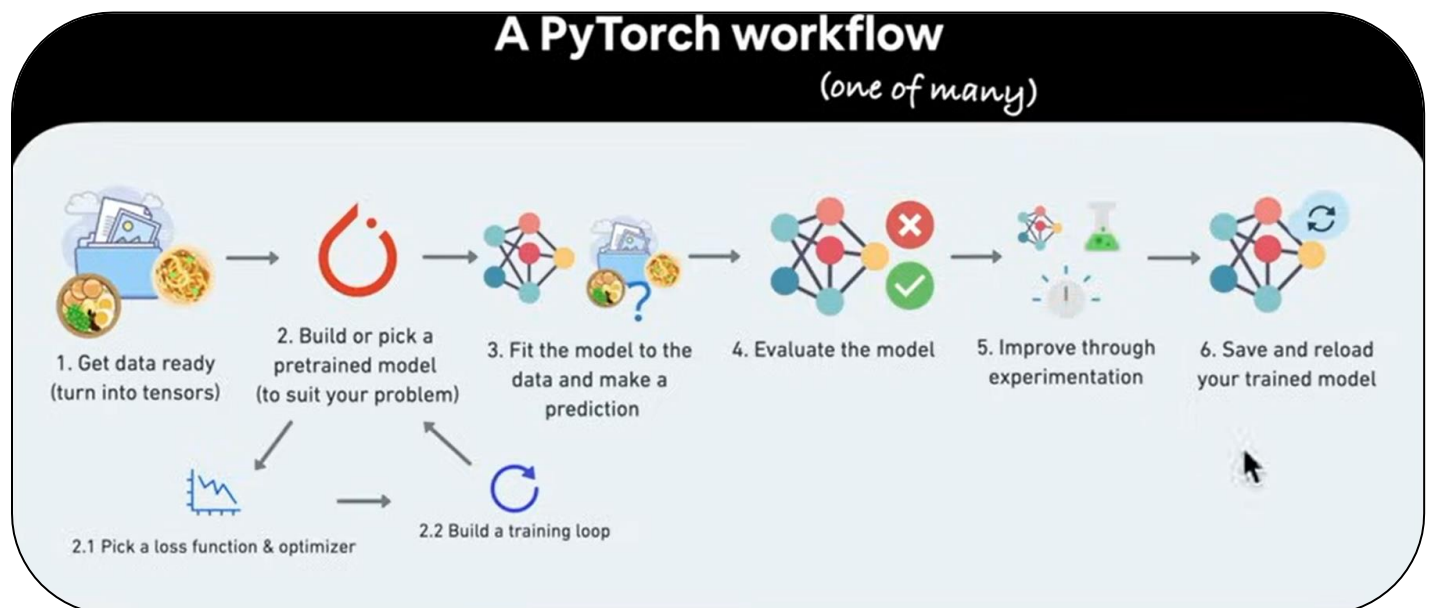
Structured data

Unstructured data

## Neural Networks:



## PyTorch



## What is PyTorch?

- Most popular research deep learning framework
- Write fast deep learning code in Python (able to run on a GPU/many GPUs)
- Able to access many pre-built deep learning models (Torch Hub / torchvision.models)
- Whole stack:** Preprocess data, model data, deploy model in your application/cloud.

- Originally designed and used in-house by Facebook/Meta (now open-source and used by companies such as Tesla, Microsoft, OpenAI)

## Why PyTorch?

- ❖ It is research favorite.

## What is a GPU/TPU?

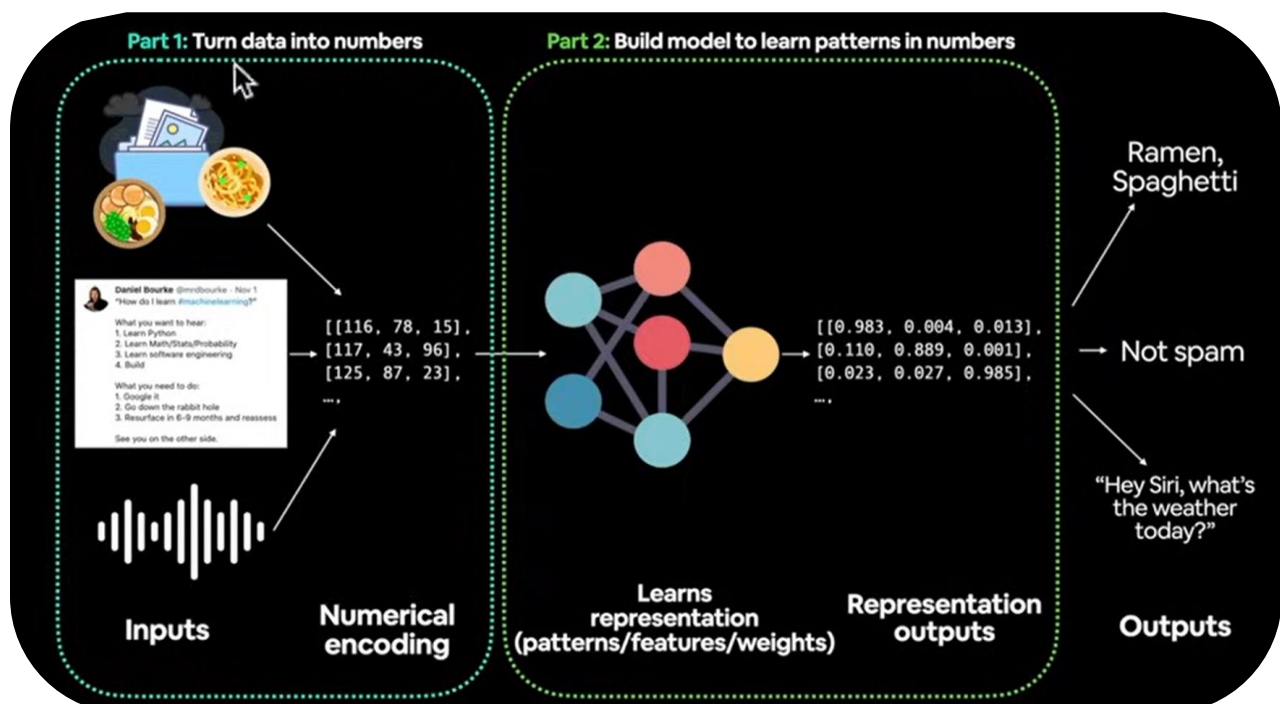
- Graphic Processing Unit (GPU):** A specialized processor designed to accelerate graphics rendering and complex calculations in parallel.
- Tensor Processing Unit (TPU):** A custom-designed chip by Google optimized for machine learning and artificial intelligence applications, particularly for handling tensor computations.
- CUDA (Compute Unified Device Architecture):** A parallel computing platform and programming model developed by NVIDIA that enables general-purpose processing on GPUs, allowing software to perform complex computations efficiently.

Machine Learning is a game of two parts:

- Get data into a numerical representation.
- Build a model to learn pattern in that numerical representation.

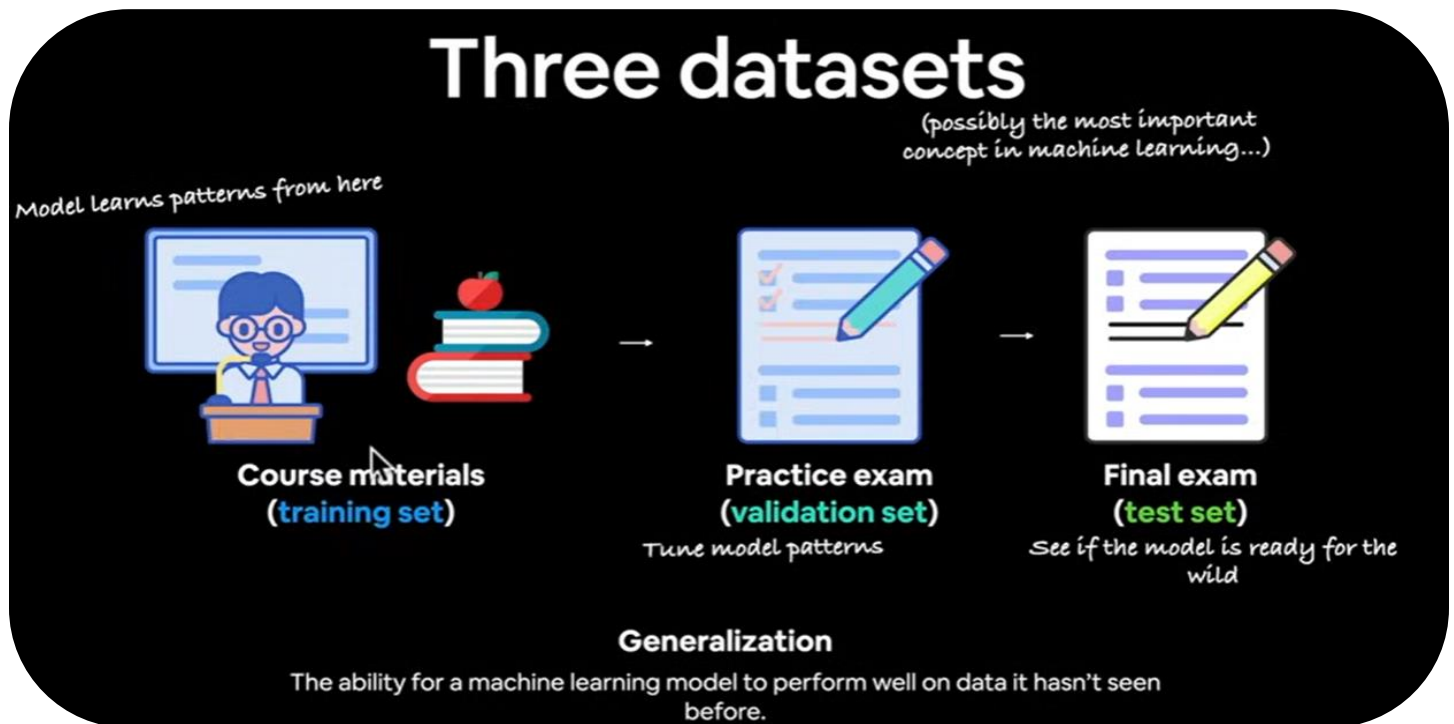
## PyTorch `permute()` function:

- ❖ **Purpose:** Reorders dimensions of a tensor.
- ❖ **Syntax:** `tensor.permute(dims)`
- ❖ **Key:** Does not modify original tensor, returns a new one.
- ❖ **Common use:** Adjusting tensor shape for operations (e.g., from [batch, channels, height, width] to [batch, height, width, channels])

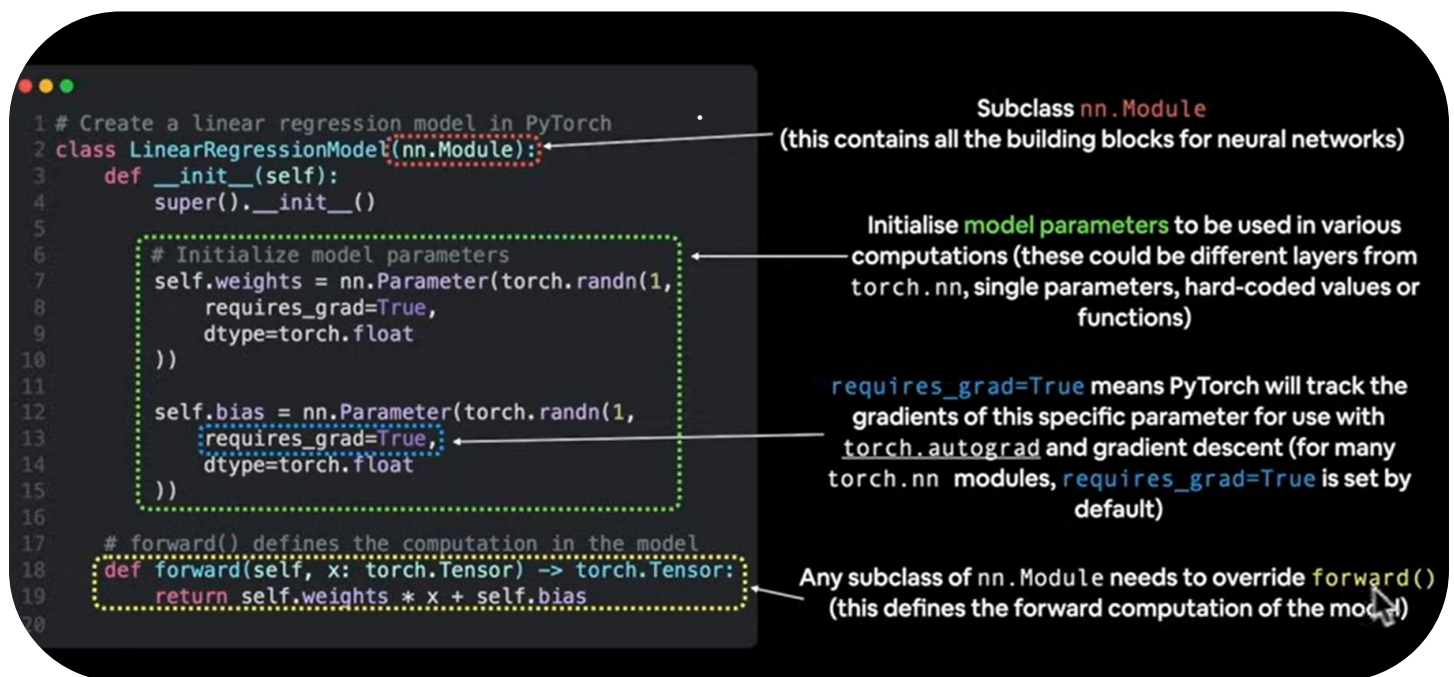




## Three Datasets



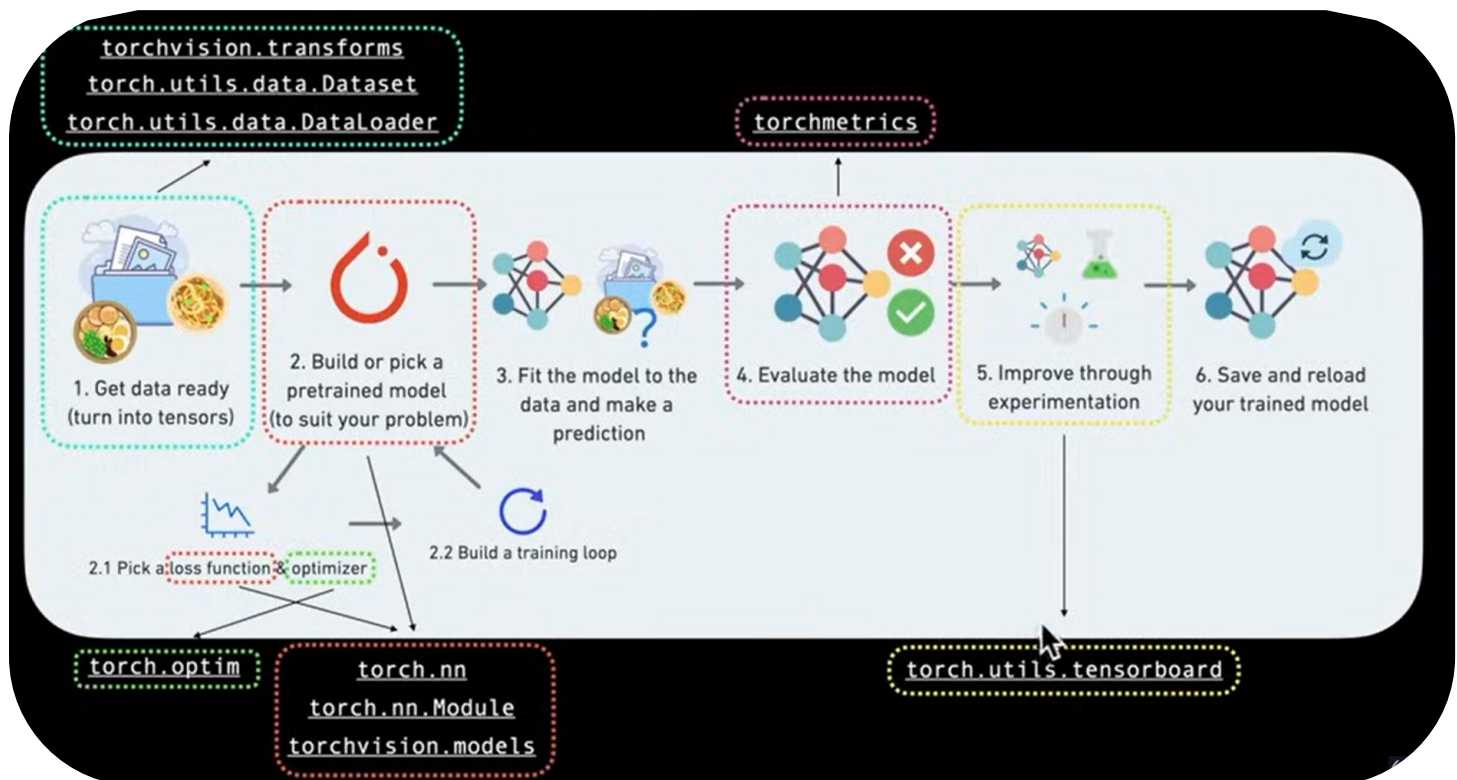
## Build Model:



## PyTorch Essential Neural Networks Building Modules:

## PyTorch essential neural network building modules

PyTorch module	What does it do?
<code>torch.nn</code>	Contains all of the building blocks for computational graphs (essentially a series of computations executed in a particular way).
<code>torch.nn.Module</code>	The base class for all neural network modules, all the building blocks for neural networks are subclasses. If you're building a neural network in PyTorch, your models should subclass <code>nn.Module</code> . Requires a <code>forward()</code> method be implemented.
<code>torch.optim</code>	Contains various optimization algorithms (these tell the model parameters stored in <code>nn.Parameter</code> how to best change to improve gradient descent and in turn reduce the loss).
<code>torch.utils.data.Dataset</code>	Represents a map between key (label) and sample (features) pairs of your data. Such as images and their associated labels.
<code>torch.utils.data.DataLoader</code>	Creates a Python iterable over a torch Dataset (allows you to iterate over your data).



## PyTorch Training Loop:

## PyTorch training loop

```

1 # Pass the data through the model for a number of epochs (e.g. 100):
2 for epoch in range(epochs):
3     # Put model in training mode (this is the default state of a model)
4     model.train()
5
6     # 1. Forward pass on train data using the forward() method inside
7     y_pred = model(X_train)
8
9     # 2. Calculate the loss (how different are the model's predictions to the true values):
10    loss = loss_fn(y_pred, y_true)
11
12    # 3. Zero the gradients of the optimizer (they accumulate by default)
13    optimizer.zero_grad()
14
15    # 4. Perform backpropagation on the loss
16    loss.backward()
17
18    # 5. Progress/step the optimizer (gradient descent)
19    optimizer.step()
20

```

Note: all of this can be turned into a function

Pass the data through the model for a number of **epochs** (e.g. 100 for 100 passes of the data)

Pass the data through the model, this will perform the **forward()** method located within the model object

Calculate the **loss value** (how wrong the model's predictions are)

Zero the **optimizer gradients** (they accumulate every epoch, zero them to start fresh each forward pass)

Perform **backpropagation** on the loss function (compute the gradient of every parameter with `requires_grad=True`)

Step the **optimizer** to update the model's parameters with respect to the gradients calculated by `loss.backward()`

## PyTorch testing loop

```

1 # Setup empty lists to keep track of model progress
2 epoch_count = []
3 train_loss_values = []
4 test_loss_values = []
5
6 # Pass the data through the model for a number of epochs (e.g. 100) epochs:
7 for epoch in range(epochs):
8
9     ### Training loop code here ###
10
11     ### Testing starts ###
12
13     # Put the model in evaluation mode
14     model.eval()
15
16     # Turn on inference mode context manager
17     with torch.inference_mode():
18         # 1. Forward pass on test data
19         test_pred = model(X_test)
20
21         # 2. Calculate loss on test data
22         test_loss = loss_fn(test_pred, y_test)
23
24     # Print out what's happening every 10 epochs
25     if epoch % 10 == 0:
26         epoch_count.append(epoch)
27         train_loss_values.append(loss)
28         test_loss_values.append(test_loss)
29         print(f"Epoch: {epoch} | MAE Train Loss: {loss} | MAE Test Loss: {test_loss}")
30

```

Create empty lists for storing useful values (helpful for tracking model progress)

Tell the model we want to **evaluate** rather than train (this turns off functionality used for training but not evaluation)

Turn on **`torch.inference_mode()`** context manager to disable functionality such as gradient tracking for inference (gradient tracking not needed for inference) *(faster performance!)*

Pass the test data through the model (this will call the model's implemented **forward()** method)

Calculate the **test loss value** (how wrong the model's predictions are on the test dataset, lower is better)

Display **information outputs** for how the model is doing during training/testing every ~10 epochs (note: what gets printed out here can be adjusted for specific problems)



## Improving a Model:

## Improving a model

(from a model's perspective)

```

1 # Create a model
2 model = nn.Sequential(
3     nn.Linear(in_features=3, out_features=100),
4     nn.Linear(in_features=100, out_features=100),
5     nn.ReLU(),
6     nn.Linear(in_features=100, out_features=3)
7 )
8
9 # Setup a loss function and optimizer
10 loss_fn = nn.BCEWithLogitsLoss()
11 optimizer = torch.optim.SGD(params=model.parameters(),
12                             lr=0.001)
13
14 # Training code...
15 epochs = 10
16
17 # Testing code...

```

Smaller model

```

1 # Create a larger model
2 model = nn.Sequential(
3     nn.Linear(in_features=3, out_features=128),
4     nn.ReLU(),
5     nn.Linear(in_features=128, out_features=256),
6     nn.ReLU(),
7     nn.Linear(in_features=256, out_features=128),
8     nn.ReLU(),
9     nn.Linear(in_features=128, out_features=3)
10 )
11
12 # Setup a loss function and optimizer
13 loss_fn = nn.BCEWithLogitsLoss()
14 optimizer = torch.optim.Adam(params=model.parameters(),
15                              lr=0.0001)
16
17 # Training code...
18 epochs = 100
19
20 # Testing code...

```

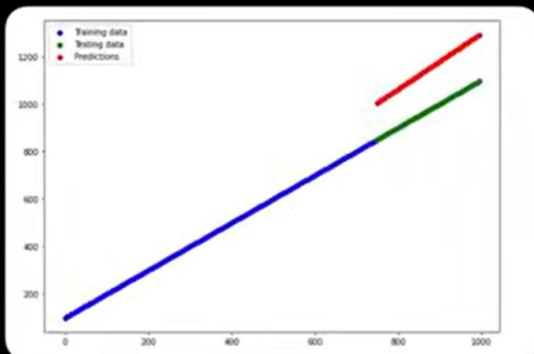
Larger model

## Common ways to improve a deep model:

- Adding layers
- Increase the number of hidden units
- Change/add activation functions
- Change the optimization function
- Change the learning rate (because you can alter each of these, they're hyperparameters)
- Fitting for longer

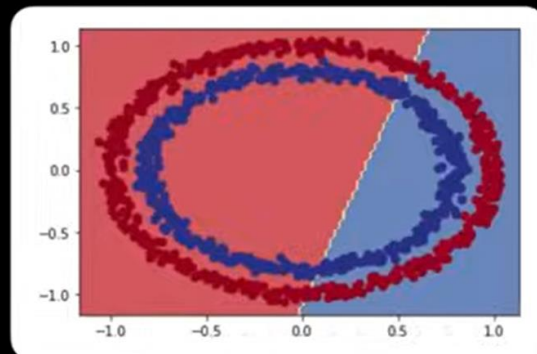
## The missing piece: Non-linearity

🤔 “What could you draw if you had an unlimited amount of straight (linear) and non-straight (non-linear) lines?”



Linear data

(possible to model with straight lines)



Non-linear data

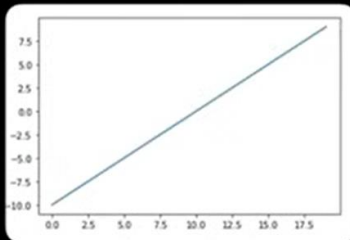
(not possible to model with straight lines)



# The missing piece: Non-linearity

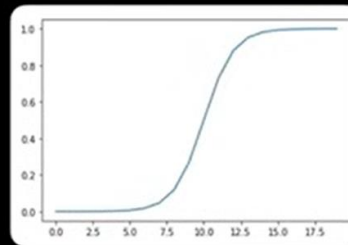
A = torch.arange(-10, 10)

plt.plot(A)



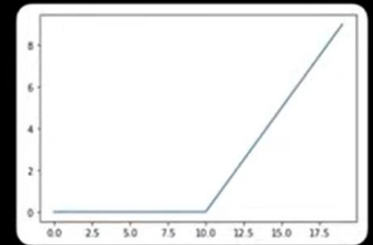
Linear activation

plt.plot(torch.sigmoid(A))



Sigmoid activation  
(non-linear)

plt.plot(torch.relu(A))



ReLU activation  
(non-linear)

## Architecture of a classification model (typical) (we're going to be building lots of these)

Hyperparameter	Binary Classification	Multiclass classification
Input layer shape (in_features)	Same as number of features (e.g. 5 for age, sex, height, weight, smoking status in heart disease prediction)	Same as binary classification
Hidden layer(s)	Problem specific, minimum = 1, maximum = unlimited	Same as binary classification
Neurons per hidden layer	Problem specific, generally 10 to 512	Same as binary classification
Output layer shape (out_features)	1 (one class or the other)	1 per class (e.g. 3 for food, person or dog photo)
Hidden layer activation	Usually ReLU (rectified linear unit) but can be many others	Same as binary classification
Output activation	Sigmoid (torch.sigmoid in PyTorch)	Softmax (torch.softmax in PyTorch)
Loss function	Binary crossentropy (torch.nn.BCELoss in PyTorch)	Cross entropy (torch.nn.CrossEntropyLoss in PyTorch)
Optimizer	SGD (stochastic gradient descent), Adam (see torch.optim for more options)	Same as binary classification

```

1 # Create a model
2 model = nn.Sequential(
3     nn.Linear(in_features=3, out_features=100),
4     nn.Linear(in_features=100, out_features=100),
5     nn.ReLU(),
6     nn.Linear(in_features=100, out_features=3)
7 )
8
9 # Setup a loss function and optimizer
10 loss_fn = nn.BCEWithLogitsLoss()
11 optimizer = torch.optim.SGD(params=model.parameters(),
12                             lr=0.001)
13
14 # Training code...
15
16 # Testing code...

```

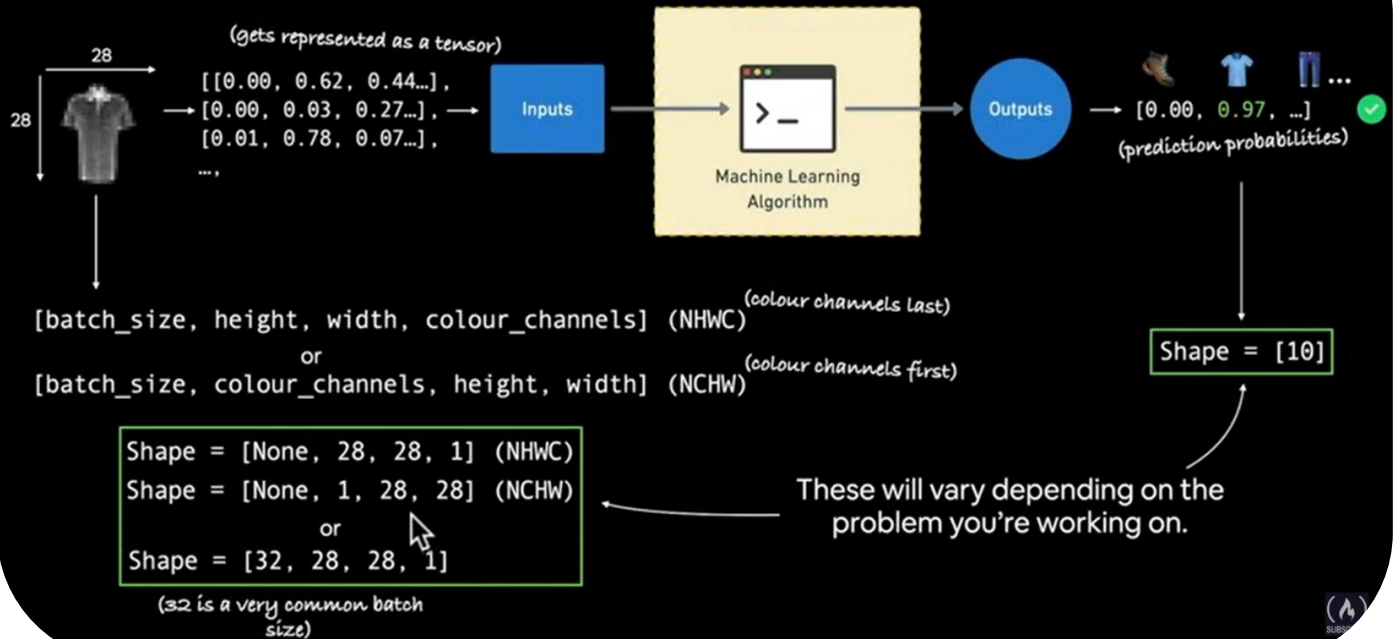
# Classification evaluation methods

Key: **tp** = True Positive, **tn** = True Negative, **fp** = False Positive, **fn** = False Negative

Metric Name	Metric Formula	Code	When to use
Accuracy	$\text{Accuracy} = \frac{tp + tn}{tp + tn + fp + fn}$	<code>torchmetrics.Accuracy()</code> or <code>sklearn.metrics.accuracy_score()</code>	Default metric for classification problems. Not the best for imbalanced classes.
Precision	$\text{Precision} = \frac{tp}{tp + fp}$	<code>torchmetrics.Precision()</code> or <code>sklearn.metrics.precision_score()</code>	Higher precision leads to less false positives.
Recall	$\text{Recall} = \frac{tp}{tp + fn}$	<code>torchmetrics.Recall()</code> or <code>sklearn.metrics.recall_score()</code>	Higher recall leads to less false negatives.
F1-score	$\text{F1-score} = 2 \cdot \frac{\text{precision} \cdot \text{recall}}{\text{precision} + \text{recall}}$	<code>torchmetrics.F1Score()</code> or <code>sklearn.metrics.f1_score()</code>	Combination of precision and recall, usually a good overall metric for a classification model.
Confusion matrix	NA	<code>torchmetrics.ConfusionMatrix()</code>	When comparing predictions to truth labels to see where model gets confused. Can be hard to use with large numbers of classes.

## Computer Vision:

## Input and output shapes



## (typical) Architecture of a CNN

Hyperparameter/Layer type	What does it do?	Typical values
Input image(s)	Target images you'd like to discover patterns in	Whatever you can take a photo (or video) of
Input layer	Takes in target images and preprocesses them for further layers	input_shape = [batch_size, image_height, image_width, color_channels] (channels last) or input_shape = [batch_size, color_channels, image_height, image_width] (channels first)
Convolution layer	Extracts/learns the most important features from target images	Multiple, can create with torch.nn.Conv2d() (X can be multiple values)
Hidden activation/non-linear activation	Adds non-linearity to learned features (non-straight lines)	Usually ReLU (torch.nn.ReLU()), though can be many more
Pooling layer	Reduces the dimensionality of learned image features	Max (torch.nn.MaxPool2d()) or Average (torch.nn.AvgPool2d())
Output layer/linear layer	Takes learned features and outputs them in shape of target labels	torch.nn.Linear(out_features=[number_of_classes]) (e.g. 3 for pizza, steak or sushi)
Output activation	Converts output logits to prediction probabilities	torch.sigmoid() (binary classification) or torch.softmax() (multi-class classification)



```

# Create a Convolutional Neural Network
import torch
from torch import nn
class CNN(nn.Module):
    def __init__(self, input_shape: int, hidden_units: int, output_shape: int):
        super().__init__()
        self.conv_layers = nn.Sequential(
            nn.Conv2d(input_shape,
                      out_channels=hidden_units,
                      kernel_size=5, # how big is the square that's going over the image?
                      stride=1, # take a step one pixel at a time
                      padding=1, # add an extra pixel around the input image
                      bias=False), # non-linear activation
            nn.ReLU(),
            nn.MaxPool2d(kernel_size=2,
                        stride=2) # default stride value is same as kernel_size
        )
        self.classifier = nn.Sequential(
            nn.Flatten(),
            nn.Linear(in_features=hidden_units * 32, # same shape as output of self.conv_layers
                      out_features=output_shape)
        )
    def forward(self, x: torch.Tensor):
        x = self.conv_layers(x)
        x = self.classifier(x)
        return x
cnn_model = CNN_model(input_shape=1, # same as number of input color channels
                      hidden_units=10,
                      output_shape=3 # same as number of classes

```

(what we're working towards building)

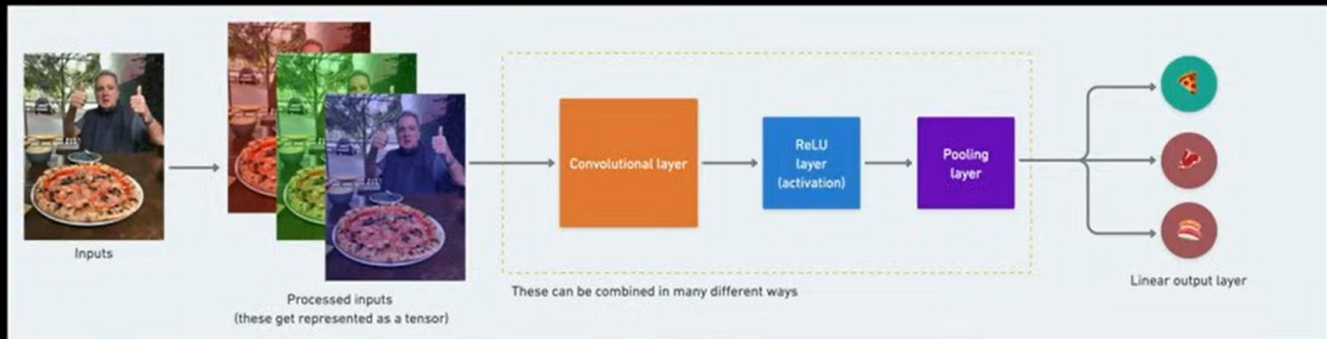
Steak 🍖  
Pizza 🍕



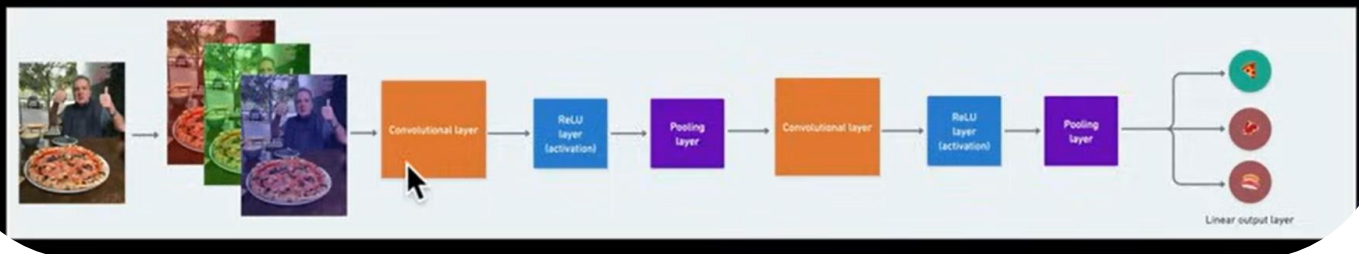
# Typical architecture of a CNN

(coloured block edition)

## Simple CNN



## Deeper CNN



# Breakdown of `torch.nn.Conv2d` layer

**Example code:** `torch.nn.Conv2d(in_channels=3, out_channels=10, kernel_size=(3, 3), stride=(1, 1), padding=0)`

**Example 2 (same as above):** `torch.nnConv2d(in_channels=3, out_channels=10, kernel_size=3, stride=1, padding=0)`

Hyperparameter name	What does it do?	Typical values
<code>in_channels</code>	Defines the number of input channels of the input data.	1 (grayscale), 3 (RGB color images)
<code>out_channels</code>	Defines the number output channels of the layer (could also be called hidden units).	10, 128, 256, 512
<code>kernel_size</code> (also referred to as filter size)	Determines the shape of the kernel (sliding windows) over the input.	3, 5, 7 (lower values learn smaller features, higher values learn larger features)
<code>stride</code>	The number of steps a filter takes across an image at a time (e.g. if <code>strides=1</code> , a filter moves across an image 1 pixel at a time).	1 (default), 2
<code>padding</code>	Pads the target tensor with zeroes (if "same") to preserve input shape. Or leaves in the target tensor as is (if "valid"), lowering output shape.	0, 1, "same", "valid"

■ **Resource:** For an interactive demonstration of the above hyperparameters, see the [CNN Explainer website](#).

# PyTorch Domain Libraries

"Is this a photo of pizza, steak or sushi?"



**TorchVision**

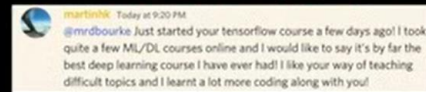
"What song is playing?"



**TorchAudio**

Different domain libraries contain data loading functions for different data sources

"Are these reviews positive or negative?"



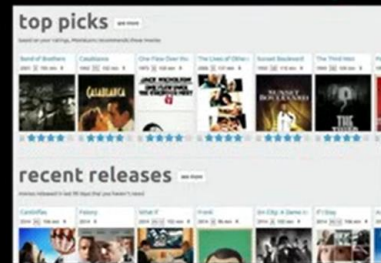
★★★★★

★★★★★



**TorchText**

"How do we recommend similar products?"



**TorchRec**

Source: movielens.org

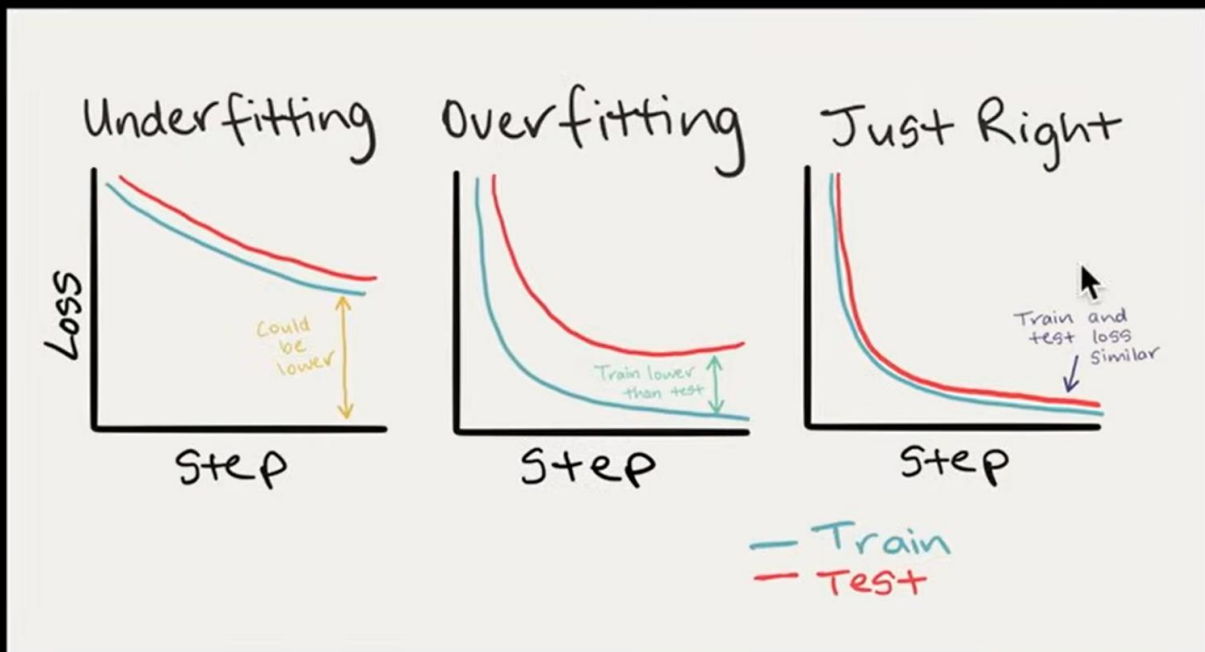
# PyTorch Domain Libraries

Problem Space	Pre-built Datasets and Functions
Vision	<a href="#">torchvision.datasets</a>
Text	<a href="#">torchtext.datasets</a>
Audio	<a href="#">torchaudio.datasets</a>
Recommendation system	<a href="#">torchrec.datasets</a>
Bonus	<a href="#">TorchData*</a>

\*TorchData contains many different helper functions for loading data and is currently in beta as of April 2022.

# Loss curves

(a way to evaluate your model's performance over time)



\*There are more combinations of these, to see them check out [Google's Interpreting Loss Curves guide](#).

## Dealing with overfitting

### Method to improve a model (reduce overfitting)

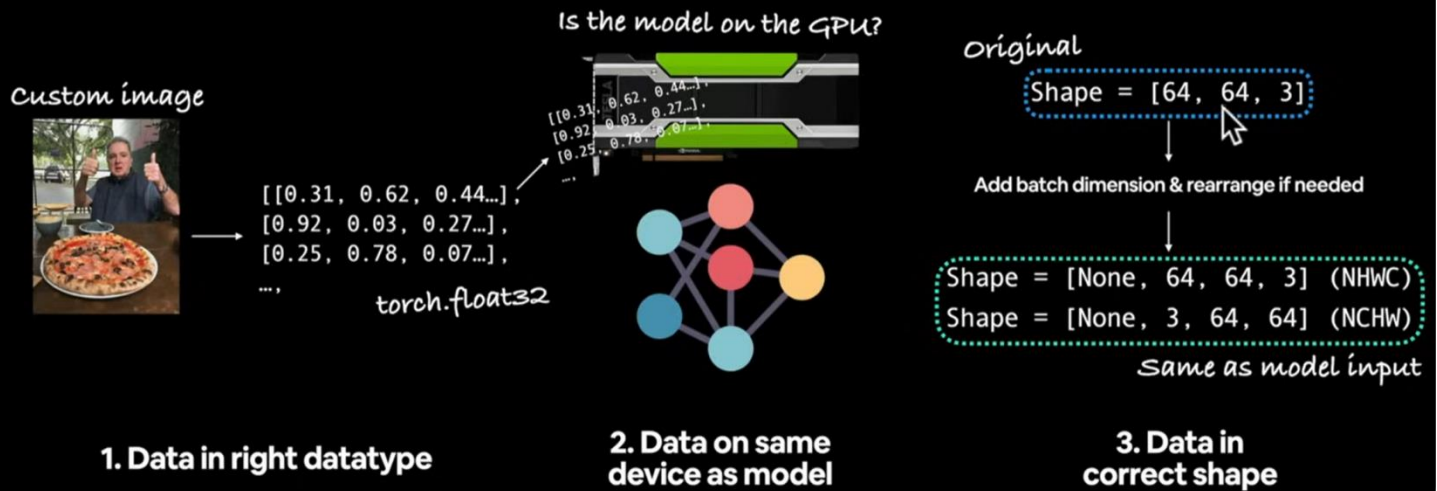
### What does it do?

Get more data	Gives a model more of a chance to learn patterns between samples (e.g. if a model is performing poorly on images of pizza, show it more images of pizza).
Data augmentation	Increase the diversity of your training dataset without collecting more data (e.g. take your photos of pizza and randomly rotate them 30°). Increased diversity forces a model to learn more generalisation patterns.
Better data	Not all data samples are created equally. Removing poor samples from or adding better samples to your dataset can improve your model's performance.
Use transfer learning	Take a model's pre-learned patterns from one problem and tweak them to suit your own problem. For example, take a model trained on pictures of cars to recognise pictures of trucks.
Simplify your model	If the current model is already overfitting the training data, it may be too complicated of a model. This means it's learning the patterns of the data too well and isn't able to generalize well to unseen data. One way to simplify a model is to reduce the number of layers it uses or to reduce the number of hidden units in each layer.
Use learning rate decay	The idea here is to slowly decrease the learning rate as a model trains. This is akin to reaching for a coin at the back of a couch. The closer you get, the smaller your steps. The same with the learning rate, the closer you get to <u>convergence</u> , the smaller you'll want your weight updates to be.
Use early stopping	<u>Early stopping</u> stops model training "before" it begins to overfit. As in, say the model's loss has stopped decreasing for the past 10 epochs (this number is arbitrary), you may want to stop the model training here and go with the model weights that had the lowest loss (10 epochs prior).



# Predicting on custom data

(3 things to make sure of...)



## Why GPUs are efficient for Deep Learning:

### 1. Memory bandwidth:

- ❖ GPUs are optimized for high bandwidth, capable of fetching large amounts of memory at once (up to 750GB/s vs 50GB/s for CPUs).

### 2. Latency hiding:

- ❖ GPUs use thread parallelism to hide their higher latency, effectively providing high bandwidth without significant latency drawbacks for large data operations.

### 3. Register and cache advantage:

- ❖ GPUs have larger aggregate register memory (up to 14MB) operating at higher speeds (80TB/s) compared to CPUs (64-128KB at 10-20TB/s). This allows for efficient storage and reuse of data tiles for matrix operations.

### 4. Parallelism:

- ❖ GPUs have many processing units (stream processors) each with their own small, fast register pack, allowing for efficient parallel computations.

### 5. Memory hierarchy utilization:

- ❖ GPUs can effectively use their memory hierarchy (main memory, L1 cache, registers) to split large matrices into smaller tiles for fast computations.

These factors combine to make GPUs particularly well-suited for the large matrix operations common in deep learning, offering significant performance advantages over CPUs for these tasks.

## Acknowledgement:

These notes were created based on insights gained from a course on YouTube. The images used have been sourced from the corresponding video [lectures](#).

## Reference:

1. [https://www.youtube.com/watch?v=V\\_xro1bcAuA](https://www.youtube.com/watch?v=V_xro1bcAuA)