# Training Technical Terms
### *feat. ChatGPT* 2023/03/17

- **LR Scheduler & Optimizer**

In machine learning, the `Scheduler` is a hyperparameter that controls the step size of the updates to the weights of the model during training. The `Optimizer` is the algorithm used to update the weights of the model based on the computed gradients during training. Below are some common options for both of the settings:

## Scheduler:

- `Constant:` The learning rate remains fixed throughout training.

- `Linear:` The learning rate decreases linearly over time, allowing for larger updates at the beginning of training and smaller updates as the model approaches convergence.

- `Cosine:` The learning rate follows a cosine curve, gradually decreasing to a small value at the end of training. This is often used in conjunction with a warm up period, during which the learning rate increases linearly from a small initial value to the maximum value.

## Optimizer:

- `Adam:` A popular adaptive optimizer that maintains separate learning rates for each parameter and updates the learning rates based on the historical gradients.

- `AdamW:` A variant of `Adam` that includes weight decay regularization to prevent overfitting.

- `AdamW8bit:` A variant of the `AdamW` with reduced memory requirement but lost precision.

- `SGD:` Stochastic Gradient Descent computes the gradient of the loss function with respect to the parameters and updates the parameters in the opposite direction of the gradient.

- `Adagrad:` An adaptive optimizer that scales the learning rate based on the historical gradients.

Choosing the right learning rate and optimizer can have a significant impact on the training of your model. It's often a matter of trial and error to find the optimal settings for your particular dataset and model architecture.

- **Learning Rate**

    The `learning rate` is a hyperparameter that controls the step size of the updates to the model's weights during training. It's an important hyperparameter to tune, as a value too large can cause the optimization process to overshoot the optimal point; while a value too small can cause the optimization process to converge too slowly or get stuck in local minima.

    The value for the learning rate depends on the specific problem you're trying to solve, as well as the architecture of your neural network. However, a common range for the learning rate is between `0.001` and `0.1`. When choosing a learning rate, you should consider the following:

    ➢ If the learning rate is too small, your model may take a long time to converge, or it may get stuck in a local minimum. If you find that your model's training loss is not decreasing even after many epochs, you may need to increase the learning rate.

    ➢ If the learning rate is too large, your model's loss may oscillate or even diverge. In this case, you should reduce the learning rate.

    ➢ If you're using an adaptive optimizer like Adam, you may be able to use a larger learning rate than you would with a non-adaptive optimizer like SGD, as the optimizer can adapt the learning rate based on the gradients.

    Some deep learning architectures may require different learning rates for different layers, as the gradients may have different magnitudes in different layers. In this case, you can try using a learning rate schedule that adjusts the learning rate for different layers.

    In general, it's a good idea to start with a relatively small learning rate and gradually increase it until you find a value that works well for your problem. You can also try using a learning rate schedule that decreases the learning rate over time, as this can help your model converge more smoothly.


- **LR Warmup**

    The value for the `LR Warmup` depends on the specifics of your model and training data, but typically ranges between `5%` and `20%` of the total number of training steps or epochs. The purpose of the warmup is to gradually increase the learning rate from a small value to its maximum value, so that the model can more easily explore the search space and avoid getting stuck in local minima.

    If you are using the `cosine` learning rate scheduler, the warmup phase should be followed by a cosine annealing phase, where the learning rate gradually decreases towards zero. The length of the annealing phase can also vary depending on the specifics of your model and training data, but typically ranges from `80%` to `95%` of the total number of training steps or epochs.

    As a starting point, you could try setting the warmup phase to `5%` to `10%` of the total number of training steps or epochs, and adjust it based on the performance of the model during training. It's also a good idea to monitor the learning rate during training to ensure that it is not too high or too low, as this can affect the stability and convergence of the training process.

- **Overfitting**

Training more steps is not always better for training. The number of steps required to train a model depends on the complexity of the problem, the size of the dataset, the architecture of the model, and the learning rate, among other factors. In some cases, training for too many steps can lead to overfitting, where the model becomes too specialized to the training data and does not generalize well to new, unseen data.

`Overfitting` occurs when a model becomes too complex and starts to fit the noise in the training data instead of the underlying pattern. This can happen when the model has too many parameters relative to the size of the dataset, or when the model is trained for too long. In some cases, overfitting can be mitigated by using regularization techniques such as dropout or weight decay, which penalize the model for using too many parameters or for fitting the noise in the data.

`Steps` and `Epochs` are often used interchangeably in the context of training neural networks. An epoch refers to one complete pass through the entire training dataset, while a step refers to one update of the model's weights based on a mini-batch of data. In practice, the number of steps required to complete one epoch depends on the batch size, which is the number of examples in each mini-batch. *For example, if you have a dataset of 1,000 examples and a batch size of 10, you would need 100 steps to complete one epoch.*

Training a model for too many epochs can cause overfitting, especially if the model is too complex relative to the size of the dataset. When a model is trained for too many epochs, it can start to memorize the training data instead of learning the underlying patterns. This can lead to a model that performs well on the training data but does not generalize well to new, unseen data.

Overfitting can be mitigated by using early stopping, which involves monitoring the model's performance on a separate validation set and stopping the training process when the performance on the validation set starts to deteriorate. Early stopping helps to prevent the model from overfitting to the training data by stopping the training process before the model has had a chance to memorize the training data.

It's worth noting that the number of epochs required to train a model depends on the complexity of the problem, the size of the dataset, and the architecture of the model. In some cases, a smaller number of epochs may be sufficient to achieve good performance, while in other cases, a larger number of epochs may be required. It's important to monitor the model's performance on a separate validation set and adjust the number of epochs accordingly to avoid overfitting.

- **Network Rank (Dimension) & Network Alpha**

`Network Rank` and `Network Alpha` are settings used in some deep learning models that use low-rank factorization to reduce the number of parameters in the model. Low-rank factorization is a technique used to approximate a matrix by factorizing it into two lower-dimensional matrices.

The **Network Rank** setting controls the dimensionality of the factorization. A higher rank corresponds to a higher-dimensional factorization and a larger number of parameters in the model. In general, a higher rank allows the model to capture more complex patterns in the data, but may also increase the risk of overfitting.

The **Network Alpha** setting controls the amount of regularization applied to the low-rank factorization. A higher value of alpha corresponds to more regularization and a smaller number of parameters in the model. Regularization helps to prevent overfitting by adding a penalty term to the loss function that encourages the model to have smaller weights.

The optimal values for **Network Rank** and **Network Alpha** depend on the specific problem and the size of the dataset. In general, a higher rank and a lower value of alpha may work well for larger datasets with more complex patterns, while a lower rank and a higher value of alpha may be more appropriate for smaller datasets with simpler patterns. It's often necessary to experiment with different values of these settings to find the optimal configuration for your specific problem.

- **Gradient Accumulation**

**Gradient Accumulation** is a technique used in deep learning to simulate training with larger batch sizes than can fit in GPU memory. When training a model, the parameters are updated based on the gradients of the loss function with respect to the parameters. In SGD and its variants, the gradients are computed based on a subset of the training data, known as a mini-batch.

With gradient accumulation, instead of updating the parameters after each mini-batch, the gradients are accumulated over multiple mini-batches before the parameters are updated. For example, if the gradient accumulation steps are set to 4, the gradients are computed and accumulated over four mini-batches before the parameters are updated based on the accumulated gradients. This has the effect of simulating a larger batch size, which can lead to more stable and accurate updates of the parameters.

Gradient accumulation can be useful when training very large models or when working with limited GPU memory. By using gradient accumulation, you can increase the effective batch size without requiring additional memory. However, gradient accumulation can also increase the training time, since more mini-batches need to be processed before each parameter update. The optimal setting for the gradient accumulation steps depends on the specific problem and the available hardware, and may need to be determined through experimentation.

- **Shuffle Caption**

Enabling **Shuffle Caption** during training of a text encoder model shuffles the order of the words or tokens in the captions that are used to train the model. This can be a useful technique to prevent the model from overfitting to the specific word order or syntax of the captions, and to encourage it to learn more general representations of the text.

By shuffling the captions, the model is forced to learn the meaning and context of each word or token independently of its position in the sentence, and to extract more robust and informative features from the text. This can also help to reduce the impact of language patterns or biases that may be present in the training data, and to improve the overall quality and diversity of the generated captions.

However, it's worth noting that shuffling the captions can also make the training process more challenging and time-consuming, as the model needs to learn to decode and reconstruct the original sentence from the shuffled tokens. Additionally, shuffling the captions may not be appropriate or effective for all types of text encoder models, and may require some experimentation and tuning to achieve the desired results.

## • Training without Captions

Training a feature recognition AI without providing captions means that the model will not have any text information to learn from, and will have to rely solely on the visual features in the images to learn how to recognize the clothing.

In this case, the model will use a purely visual approach to recognize the clothing, which can be effective for certain types of datasets and tasks. However, it's important to note that using only visual features may limit the model's ability to generalize to new or unseen examples, as it may not be able to recognize the clothing in different contexts or environments where the visual features are different.

Additionally, without text information, the model will not be able to generate any textual descriptions or captions of the clothing, and will be limited to recognizing and classifying the clothing based on its visual features alone.

Overall, the effectiveness of training a feature recognition AI without captions will depend on the specifics of the dataset and the task at hand, and it may be necessary to experiment with different approaches and techniques to achieve the desired results.

## • Precision

➢ **float:** refers to `32-bit` floating point precision, which is the most common format used in deep learning. It provides high precision and accuracy, but requires more memory and computation time compared to lower precision formats.

➢ **fp16:** refers to `16-bit` floating point precision, which is a lower precision format that can reduce memory usage and speed up computation. However, it can result in numerical instability and loss of accuracy, especially for large or complex models.

➢ **bf16:** refers to `16-bit` floating point format designed to improve numerical stability while still providing some benefits of lower precision formats. It is used in some specialized hardware and software environments.