# Probabilistic generative models

Advanced Techniques and Optimization of Transformer Models

https://tinyurl.com/vbmk4p9n

# Advanced Fine-Tuning of Transformers: LORA – Context and Motivation for Fine-Tuning

**Adapting Large Pre-trained Transformer Models to Specific Tasks with Limited Data:**

- Pre-trained transformer models like BERT, GPT-2, and others are often trained on vast corpora of data, but significant challenges arise when adapting them to specific tasks with limited available data.

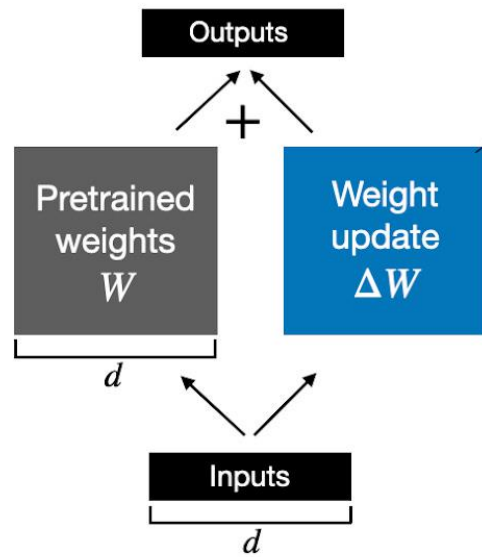**Importance of Optimization for High Performance:**

- Optimizing models to maximize their performance is crucial. This includes not only adjusting model weights and biases but also regularizing to prevent overfitting and optimizing learning rates for different layers.

# Advanced Fine-Tuning of Transformers: LORA – Context and Motivation for Fine-Tuning
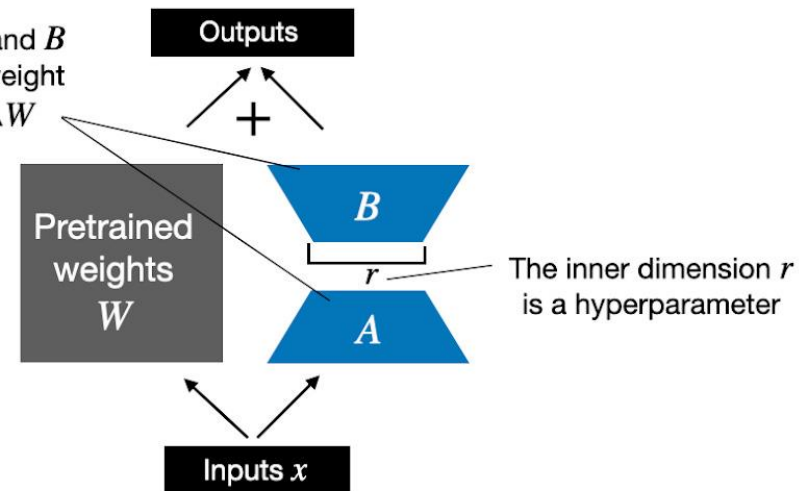
**LORA (Learning Optimization and Regularization Approach):**

- Technique to Improve Fine-Tuning of Large Pre-trained Models

- LORA combines selective regularization and optimization techniques to improve the efficiency of fine-tuning transformer models.



HU, Edward J., SHEN, Yelong, WALLIS, Phillip, *et al.* Lora: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*, 2021.

# Advanced Fine-Tuning of Transformers: LORA – Principles of LORA

**Selective Regularization:**

- Apply different regularization techniques according to the model layers to prevent overfitting and improve generalization.

- Use Dropout on the final layers where overfitting is more likely. Add Dropout to high-level layers to prevent overfitting, e.g., a high Dropout rate (0.5) on the final layer.

- Apply L2 Regularization on fully connected layers to limit model complexity. Apply L2 regularization penalty to certain layers to prevent weights from becoming excessively large.

**Selective Optimization:**

- Adapt learning rates selectively for different model layers. Adjust learning rates for different layers, e.g., a lower learning rate for the initial layers (deep layers) compared to the final layers (shallow layers).

- Dynamically adjust hyperparameters based on observed performance during fine-tuning. Incorporate techniques to dynamically adjust hyperparameters during training in response to validation performance.

# Case Studies and Practical Implementations – Text Classification

**Implementing LORA on a BERT Model for Text Classification with a Small Dataset**

- **Data Preparation:**

- **Loading and Preparing the Model:**

- **Applying Dropout and L2 Regularization:**

```python
from datasets import load_dataset
dataset = load_dataset("imdb")
```

```python
from transformers import AutoTokenizer, AutoModelForSequenceClassification

tokenizer = AutoTokenizer.from_pretrained("bert-base-uncased")
model = AutoModelForSequenceClassification.from_pretrained("bert-base-uncased")

def tokenize_function(examples):
    return tokenizer(examples["text"], padding="max_length", truncation=True)

tokenized_datasets = dataset.map(tokenize_function, batched=True)
tokenized_datasets = tokenized_datasets.shuffle(seed=42)
```

```python
from torch.nn import Dropout, Linear

class BERTWithDropout(nn.Module):
    def __init__(self, model):
        super(BERTWithDropout, self).__init__()
        self.bert = model.bert
        self.dropout = Dropout(0.3)
        self.classifier = Linear(model.config.hidden_size,
        model.config.num_labels)

    def forward(self, input_ids, attention_mask=None, labels=None):
        outputs = self.bert(input_ids=input_ids,
        attention_mask=attention_mask)
        pooled_output = outputs[1]
        pooled_output = self.dropout(pooled_output)
        logits = self.classifier(pooled_output)
        return logits

model = BERTWithDropout(model)
```

# Case Studies and Practical Implementations – Text Classification

**Implementing LORA on a BERT Model for Text Classification with a Small Dataset**

- **Selective Optimization:**

```python
from transformers import AdamW, get_linear_schedule_with_warmup

optimizer = AdamW(model.parameters(), lr=2e-5)

def set_optimizer_grouped_parameters(model, base_model_name, lr, epsilon,
weight_decay):
            no_decay = ["bias", "LayerNorm.weight"]
            optimizer_grouped_parameters = [
                        { "params": [p for n, p in model.named_parameters()
                        if not any(nd in n for nd in no_decay)],
                        "weight_decay": weight_decay,
                        "lr": lr,
                        "epsilon": epsilon,
                        },
                        { "params": [p for n, p in model.named_parameters()
                        if any(nd in n for nd in no_decay)],
                        "weight_decay": 0.0,
                        "lr": lr,
                        "epsilon": epsilon, },
            ]
            return optimizer_grouped_parameters

optimizer_grouped_parameters = set_optimizer_grouped_parameters(
model, base_model_name="bert", lr=2e-5, epsilon=1e-8, weight_decay=0.01 )
```

# Case Studies and Practical Implementations – Text Classification

**Implementing LORA on a BERT Model for Text Classification with a Small Dataset**

- **Training:**

```python
from transformers import TrainingArguments, Trainer

training_args = TrainingArguments(
            output_dir="./results",
            evaluation_strategy="epoch",
            learning_rate=2e-5,
            per_device_train_batch_size=16,
            per_device_eval_batch_size=16,
            num_train_epochs=3,
            weight_decay=0.01, )

trainer = Trainer(
            model=model,
            args=training_args,
            train_dataset=tokenized_datasets["train"],
            eval_dataset=tokenized_datasets["test"], )

trainer.train()
```

# Case Studies and Practical Implementations – Text Classification

**Results and Comparison with Traditional
Fine-Tuning Approaches:**

- **Measuring Performance:**

```python
eval_results = trainer.evaluate()
print(f"Accuracy: {eval_results['eval_accuracy']}")
```

## Analysis:

- Compare performance metrics (accuracy, recall, F1-score) with and without applying LORA to evaluate the effectiveness of the technique.

# Case Studies and Practical Implementations – Text Classification

**Applying LORA to GPT-2 for Conditioned Text Generation**

- **Data Preparation:**

- **Loading and Preparing the Model:**

- **Selective Optimization and Regularization:**

```python
from datasets import load_dataset
dataset = load_dataset("wikitext", "wikitext-2-raw-v1")
```

```python
from transformers import GPT2Tokenizer, GPT2LMHeadModel

tokenizer = GPT2Tokenizer.from_pretrained("gpt2")
model = GPT2LMHeadModel.from_pretrained("gpt2")
```

```python
from transformers import Trainer, TrainingArguments, AdamW

training_args = TrainingArguments( output_dir='./results',
                learning_rate=5e-5,
                weight_decay=0.01,
                per_device_train_batch_size=2,
                per_device_eval_batch_size=2,
                num_train_epochs=3,
                logging_dir='./logs',
                logging_steps=10,
                evaluation_strategy="epoch" )

optimizer = AdamW(set_optimizer_grouped_parameters(model, "gpt2", 5e-5, 1e-8, 0.01))

trainer = Trainer( model=model,
                args=training_args,
                tokenizer=tokenizer,
                train_dataset=dataset["train"],
                eval_dataset=dataset["test"] )

trainer.train()
```

# Case Studies and Practical Implementations – Text Classification

**Results and Analysis of Performance Gains
and Quality of Generated Text**

- **Model Evaluation:**

```
eval_results = trainer.evaluate()
print(f"Perplexity: {eval_results['eval_loss']}")
```

**Analysis:**

- Evaluate the fluidity and coherence of the generated texts.

- Compare results obtained with and without applying LORA.

**Conclusion**

Using LORA (Learning Optimization and Regularization Approach) can significantly improve the performance of transformers during fine-tuning on small specific tasks, by combining selective regularization and optimization techniques. The case studies demonstrate potential performance gains in text classification and conditioned text generation tasks.

# Optimization and Scalability – Principles of Optimization in Transformer Architectures

**Optimization Challenges:**

- **Challenges with Transformers and Long Sequences :**
    - Current language models, such as Transformers, have limitations regarding the context length they can handle. In practice, they are often limited to a few thousand tokens, which may be insufficient for tasks requiring an understanding of very long contexts, such as continuous dialogues or multi-page texts.
    - These limitations can lead to a loss of coherence and relevance when language models need to process long input sequences.

# Optimization and Scalability – Principles of Optimization in Transformer Architectures

**Introduction to SAMBA (Scalable Adaptive Mutual Backpropagation):**

- **Advanced Techniques for Enhancing Transformer Efficiency and Scalability:**
  - SAMBA aims to optimize memory consumption and accelerate computational time while ensuring rapid and stable convergence for very large models.
  - It implements advanced techniques that adapt backpropagation processes and memory management.

# Advanced Methods (SAMBA) – SAMBA Architecture

**Adaptive and Mutualized Backpropagation:**

- **Gradient Adjustment:**
  - SAMBA dynamically adjusts gradients based on each layer's specific needs.
  - This adaptation is done by monitoring and adjusting learning rates locally for each layer, allowing for finer and faster convergence.

- **Backpropagation Information Sharing:**
  - Layers exchange backpropagation information, allowing for mutualization of gradients.
  - This strategy aims to avoid redundancy in calculations and harmonize gradient adaptation across different layers of the model.

**Memory Optimization:**

- **Gradient Calculation Strategies:**
  - Implement advanced gradient calculation methods such as partial gradient reuse to save memory during backpropagation steps.

- **Reversible Backpropagation:**
  - A technique that allows recalculating necessary activations for backpropagation instead of storing them, significantly reducing memory usage.

- **Distributed Computing:**
  - Distribute gradient calculations across multiple GPUs or nodes using distributed computing frameworks like Horovod or distributed memory-sharing mechanisms.
  - This allows for better usage of global resources and accelerates training times.

# Advanced Methods (SAMBA) – Scalability and Application to Large Models

**Model Partitioning Techniques:**

- **Dividing the Model into Sub-models:**
  - Splitting the model into multiple sub-models, each trained independently or in a coordinated manner.

- **Parallel Pipelines:**
  - Using parallel pipelines where different parts of the model are trained simultaneously on different hardware.
  - This approach not only accelerates training but also enables handling larger models by resolving memory constraints.

**Distributed Training:**

- **Using GPU Clusters and Cloud Infrastructures:**
  - Utilizing GPU clusters to distribute computations among multiple processing units. For example, using clusters available on cloud platforms like AWS, Google Cloud, or Azure.

- **Horovod and Other Distributed Frameworks:**
  - Adopting distributed computing frameworks like Horovod to facilitate distributed training of transformer models across multiple GPUs, significantly increasing training speed.

# Advanced Methods (SAMBA) – Case Study and Practical Implementation

## Example Implementation with SAMBA

- **Preparation and Configuration:**

- **Example Code with Distributed Computing:**

```
pip install torch distributed
pip install horovod
```

```python
import torch
import torch.optim as optim
import horovod.torch as hvd
from transformers import BertForSequenceClassification, BertTokenizer

# Initialize Horovod
hvd.init()
# Configure the GPU to use torch.cuda.set_device(hvd.local_rank())
device = torch.device("cuda")
# Load the model and tokenizer
model_name = "bert-base-uncased"
model = BertForSequenceClassification.from_pretrained(model_name)
tokenizer = BertTokenizer.from_pretrained(model_name) model.to(device)
# Optimizer
optimizer = optim.Adam(model.parameters(), lr=1e-5 * hvd.size())
# Distributed gradient reduction technique
optimizer = hvd.DistributedOptimizer(optimizer,named_parameters=model.named_parameters())
# Initialization of parameter broadcast
hvd.broadcast_parameters(model.state_dict(), root_rank=0)

# Simplified training example
def train_step(data, targets):
        model.train()
        outputs = model(data.to(device), labels=targets.to(device))
        loss = outputs.loss optimizer.zero_grad() loss.backward()
        optimizer.step()
        return loss.item()
```

# Advanced Methods (SAMBA) – Case Study and Practical Implementation

## Computation Offloading and Memory Optimization

- **Activation Checkpointing:**

```python
import torch

def checkpointing(layer, *args):
        def custom_forward(*inputs):
                return layer(*inputs)
        return torch.utils.checkpoint.checkpoint(custom_forward, *args)

for module in model.modules():
        if isinstance(module, BertLayer):
                module.forward = checkpointing
```

- **Conclusion**

Advanced optimization and scalability of transformers involve adopting techniques like SAMBA, which combine adaptive gradient management and memory reuse. Distributed computing practices and model partitioning not only facilitate training very large models but also their application in production. When used together, LORA and SAMBA can maximize the efficiency and performance of transformers, even with limited resources.

# Advanced Projects and Applications – Selection of an Advanced Project Topic

**Identifying Project Topics:**

Consider current trends and critical applications in the field of AI and transformers.

**Project Proposals:**

- **Project 1: Fake News Detection with BERT**
  - **Objective:** Use BERT to classify articles as true or false.
  - **Methodology:** Analyze news articles and identify relevant features. Fine-tune BERT on a fake news dataset.
  - **Evaluation:** Measure accuracy, recall, and F1 score.

- **Project 2: Art Generation from Text Descriptions with GPT-3**
  - **Objective:** Use GPT-3 to generate artistic descriptions and graphical representations based on textual prompts.
  - **Methodology:** Create specific prompts and analyze the results generated by GPT-3.
  - **Evaluation:** Use qualitative and quantitative evaluation techniques (e.g., CSS).

- **Project 3: Automatic Summarization of Long Documents with Transformer Model**
  - **Objective:** Use a transformer model to generate automatic summaries of long texts.
  - **Methodology:** Preprocess long documents, fine-tune a summarization model like BART.
  - **Evaluation:** Use metrics such as BLEU, ROUGE to assess summarization performance.

# Advanced Projects and Applications – Design, Implementation, and Optimization

## Elaboration of Project Plan

- **Defining Objectives:** Clarify the project goal and expected outcomes.

- **Implementation Steps:** Identify key project stages, including data collection and preprocessing, model fine-tuning, and result evaluation.

- **Implementation Methodologies:** Choose appropriate techniques and tools for each project stage (e.g., using Hugging Face, PyTorch).

# Advanced Projects and Applications – Preparation of a Results Presentation

**Structuring the Presentation:**

- **Context and Problem Statement:** Present the project context and the problem addressed.

- **Methodology:** Describe the techniques and algorithms used, including LORA and SAMBA.

- **Results:** Present the obtained results, including evaluation metrics.

- **Discussion:** Analyze performance, encountered challenges, and provided solutions.

- Create **slides** to summarize the main project steps and results. Example Presentation Structure :

```python
import matplotlib.pyplot as plt

# Performance visualization
epochs = range(1, len(trainer.state.log_history) + 1)
loss_values = [x['loss'] for x in trainer.state.log_history if 'loss' in x]

plt.plot(epochs, loss_values, label="Loss")
plt.xlabel("Epochs")
plt.ylabel("Loss")
plt.legend()
plt.show()
```

**Introduction:**
   **Context and Problem Statement**
**Methodology:**
   **Techniques Used**
   **Trained Models**
**Results:**
   **Obtained Results**
   **Evaluations**
**Discussion:**
   **Analysis**
   **Future Improvements**