# Advanced NLP
# Assignment 3 - Report

**Theory Questions**

**Question 1**
**Concept of Soft Prompts:** How does the introduction of "soft prompts" address the limitations of discrete text prompts in large language models? Why might soft prompts be considered a more flexible and efficient approach for task-specific conditioning?

**Answer:**

The concept of "soft prompts" offers a way to address the limitations of discrete text prompts when using large language models (LLMs) by introducing a learnable, continuous embedding layer instead of relying on hard-coded or human-designed textual prompts. Here's how soft prompts improve upon discrete text prompts:

**Limitations of Discrete Text Prompts**

1. **Limited Flexibility and Expressiveness:** Discrete prompts are often hand-crafted, limiting the model's flexibility. The wording and structure of these prompts can significantly impact the model's response, yet designing optimal prompts requires trial and error.
2. **Difficulty in Optimization:** Discrete prompts are not directly optimized through gradient-based learning since they are fixed tokens. As a result, their effectiveness can be inconsistent, especially for complex or nuanced tasks where exact phrasing matters.
3. **Dependency on Natural Language:** Discrete prompts rely on human language to encode information, which may not always align with the model's internal representations, leading to inefficiencies in communicating the intended task.

**Advantages of Soft Prompts**

1. **Task-Specific Adaptability:** Soft prompts are learnable embeddings, enabling the model to fine-tune task-specific representations directly through training, which improves adaptability and precision for various tasks.

2. **Better Utilization of Latent Space:** Since soft prompts operate in the continuous embedding space, they leverage the rich semantic relationships within the model's internal representations more effectively than discrete tokens, leading to enhanced task conditioning.

3. **Reduced Memory and Computational Overhead:** Soft prompts typically require fewer parameters and less training time than full fine-tuning, making them more memory-efficient and computationally efficient for tasks that don't require model-wide adjustments.

4. **Robustness Across Tasks:** By bypassing the need for precise linguistic phrasing, soft prompts can generalize better across tasks, reducing the need for extensive prompt engineering. This makes them particularly useful for multi-task learning and cross-domain applications.

In summary, soft prompts allow for a more nuanced, efficient, and adaptable conditioning approach than discrete text prompts, aligning closer with the model's representational capabilities and enabling more consistent performance across various tasks.

**Question 2:**

**Scaling and Efficiency in Prompt Tuning:** How does the efficiency of prompt tuning relate to the scale of the language model? Discuss the implications of this relationship for future developments in large-scale language models and their adaptability to specific tasks.

**Answer:**

Prompt tuning becomes increasingly efficient as language models scale up due to the inherent design of large models, which capture rich, multi-faceted representations that can be selectively activated or conditioned by smaller, task-specific prompt embeddings. This relationship has important implications for how we develop and apply large-scale language models to specific tasks in the future. Here's how scaling and efficiency in prompt tuning relate:

**Relationship between Prompt Tuning Efficiency and Model Scale**

1. **Enhanced Representation Capacity:** Larger language models have more parameters, resulting in richer latent spaces and a greater ability to encode a diverse set of concepts and relationships. This means that a relatively small prompt, when applied to a large model, can activate relevant parts of the model's knowledge without the need to fine-tune all parameters, making prompt tuning more efficient as the model scale grows.

2. **Reduction in Task-Specific Training Requirements:** As model size increases, the need for extensive training on task-specific data decreases. In large models, even lightweight prompt tuning can sufficiently leverage pre-existing representations to adapt to new tasks, minimizing the computational cost and time required for fine-tuning.

3. **Parameter-Efficient Tuning:** For very large models, prompt tuning only requires adjusting a small fraction of the model's parameters or embeddings. This drastically reduces storage and computational needs, which is especially beneficial for deploying these models on limited hardware or at scale across multiple applications.

**Implications for Future Developments**

1. **Scalability of Customization:** As models grow, prompt tuning provides a scalable way to customize them for numerous tasks without needing separate, fully fine-tuned models. This approach enables deploying a single large model in various domains and applications by simply swapping prompt embeddings, reducing redundancy and enhancing model reuse.
2. **Lower Barrier to Task Adaptation:** For both research and commercial use, prompt tuning on large models lowers the barrier to adapting these models for specific, even niche, tasks. With efficient prompt tuning, end-users and developers can create task-specific adaptations with minimal resources.
3. **Foundation for Few-Shot and Zero-Shot Learning:** In the context of large models, prompt tuning aligns well with few-shot or zero-shot learning approaches, where models rely heavily on their pre-trained knowledge base to handle new tasks. As models scale further, this method of task conditioning could allow models to generalize better with even less task-specific data.

In summary, as large-scale language models continue to evolve, prompt tuning's efficiency enhances the adaptability and scalability of these models, reducing the need for resource-intensive fine-tuning and supporting broader use across specialized tasks. This efficiency positions prompt tuning as a critical technique in the practical deployment and utilization of future LLMs in both industry and academia.

**Question 3:**

**Understanding LoRA:** What are the key principles behind Low-Rank Adaptation (LoRA) in fine-tuning large language models? How does LoRA improve upon traditional fine-tuning methods regarding efficiency and performance?

**Answer:**

Low-Rank Adaptation (LoRA) is an innovative approach to fine-tuning large language models (LLMs) that reduces computational cost and memory usage while maintaining strong performance. LoRA achieves this by introducing low-rank matrices that adaptively modify a model's weights without updating the full set of model parameters, offering significant advantages over traditional fine-tuning methods.

**Key Principles Behind LoRA**

1. **Low-Rank Decomposition**: LoRA fine-tunes a model by decomposing the weight updates into two low-rank matrices, AAA and BBB, rather than modifying the full weight matrix directly. The rank of these matrices is much smaller than the original, which limits the number of parameters that need updating.
   - **Parameter Efficiency**: Instead of adjusting all model weights, LoRA inserts two small, learnable matrices into specific layers of the model (typically, the attention layers in transformer models) that capture task-specific information without affecting the full weight structure.

2. **Frozen Base Model with Adaptation Layers**: In LoRA, the original weights of the model remain frozen, and only the low-rank adaptation matrices AAA and BBB are trained. This means that LoRA captures task-specific information through these additional layers rather than overhauling the entire model.

3. **Linear Transformation of Weight Updates**: By representing the weight changes in a low-rank format, LoRA approximates the parameter updates in a way that maximizes information retention while limiting the scale of updates. This allows LoRA to learn specialized knowledge for a task without requiring a full fine-tuning of the model's vast parameter space.

**Advantages of LoRA Over Traditional Fine-Tuning**

1. **Memory and Computational Efficiency**: Since LoRA only updates a small fraction of the model parameters, it requires significantly less GPU memory and computational power compared to traditional fine-tuning, which alters the entire model. This is especially beneficial for large models, where full fine-tuning would demand massive resources.
2. **Parameter Efficiency and Storage**: LoRA's low-rank matrices are orders of magnitude smaller than the full model parameters, resulting in fewer parameters to store and manage. For practical deployment, this means that task-specific adaptations can be stored and loaded as compact parameter sets, making it feasible to deploy multiple task-specific variations of a model without redundant weight storage.
3. **Maintained Performance with Fewer Updates**: Despite modifying only a small portion of the model's parameters, LoRA often achieves comparable or even superior performance on task-specific benchmarks. This efficiency arises from leveraging the pre-trained knowledge in the frozen base model and only introducing targeted, low-rank updates that enhance task-relevant features.
4. **Modular Adaptability**: LoRA enables modularity by allowing different low-rank adaptation matrices for different tasks to be added as layers on top of a single base model. This modular approach enables multiple task adaptations without re-training the entire model, allowing for efficient multi-tasking and deployment in diverse applications.

**Summary**

LoRA represents a significant step forward in efficient model fine-tuning, especially for very large models. By employing low-rank adaptation matrices that reduce the need for full parameter updates, LoRA achieves both computational and storage efficiency, facilitating faster and cheaper task-specific fine-tuning. This efficiency makes LoRA highly practical for applications where resource constraints or the need for multiple task-specific adaptations are primary concerns.

**Question 4:**

**Theoretical Implications of LoRA:** Discuss the theoretical implications of introducing low-rank adaptations to the parameter space of large language models. How does this affect the expressiveness and generalization capabilities of the model compared to standard fine-tuning?

**Answer:**

Introducing low-rank adaptations with LoRA has intriguing theoretical implications for large language models, particularly concerning their expressiveness and generalization capabilities. By limiting parameter updates to low-rank matrices, LoRA constrains the model's adaptive capacity to task-specific signals in a structured way, impacting both expressiveness and generalization in distinct ways compared to standard fine-tuning.

**1. Constraining the Parameter Space**

- **Inductive Bias through Low-Rank Structure**: LoRA's low-rank constraint introduces an inductive bias, focusing adaptation on only the most salient features for a given task rather than modifying all parameters. This aligns with the assumption that many tasks require

only a small subset of the model's representational capacity, thus streamlining the learning process.

- **Controlled Expressiveness**: By updating parameters within a low-rank subspace, LoRA limits expressiveness to changes that can be captured by the low-rank matrices. While this could theoretically reduce the ability to fit extremely complex or non-linear patterns, in practice, it encourages the model to focus on efficient representations that generalize better to the task domain without overfitting.

## 2. Enhanced Generalization Capabilities

- **Reduced Overfitting**: Traditional fine-tuning often risks overfitting, especially when adapting large models to small datasets. LoRA's low-rank approach constrains the extent of updates, reducing the model's capacity to memorize irrelevant patterns. This often results in better generalization since LoRA focuses on fine-tuning only a subset of the parameters that are most relevant to the task.
- **Alignment with Pre-Trained Knowledge**: Since LoRA retains the pre-trained model weights and only learns task-specific low-rank updates, it effectively operates as a specialized, fine-grained adjustment on top of the existing knowledge. This selective modification makes LoRA well-suited for tasks where generalization is more valuable than task-specific expressiveness, as it can rely on the frozen weights' broad representational knowledge while adapting minimally to new data.

## 3. Theoretical Trade-offs in Expressiveness

- **Loss of Full-Model Adaptability**: In theory, by constraining updates to a low-rank space, the model sacrifices some adaptability that full fine-tuning offers, potentially limiting its ability to learn highly nuanced or domain-specific features if these cannot be captured within the low-rank adaptations. However, LoRA's selective adaptation often suffices for most tasks, especially those that do not require extensive, model-wide changes.

- **Localized Task-Specific Representations**: The low-rank constraints in LoRA encourage the model to localize task-specific adaptations within a restricted portion of its overall representational space. This structured limitation may reduce the risk of altering core model representations that could be valuable across different tasks, thus preserving the model's broader utility.

**4. Implications for MultiTask and Transfer Learning**

- **Modular Adaptability Across Tasks**: LoRA's low-rank approach inherently aligns with modular adaptations, meaning that different low-rank matrices can be tailored for different tasks while leveraging a shared base model. This modularity is beneficial for multi-task learning and transfer learning, as it enables models to adapt across various tasks without compromising generalization.
- **Efficient Reuse of Knowledge**: By isolating task-specific adjustments to low-rank matrices, LoRA allows for efficient knowledge transfer, where the same model can quickly adapt to new tasks by adding small, low-rank updates. This reusability supports both theoretical and practical advantages for deploying adaptable, multi-task models with minimal retraining.

**Summary**

LoRA's low-rank adaptation introduces a structured constraint on expressiveness, favoring generalization and efficient use of pre-trained knowledge over exhaustive task-specific adjustments. This approach allows large models to be more adaptable and modular, making LoRA particularly suited for scenarios where task-specific adaptability must balance with broad, general-purpose utility. Theoretically, LoRA promotes better generalization by focusing on core, transferable features rather than overfitting to task idiosyncrasies, positioning it as a powerful fine-tuning strategy for scalable, efficient adaptation in large models.

# Report on Fine-Tuning GPT-2 for Summarization Task

## 1. Introduction

The goal of this architecture was to fine-tune GPT-2 for a summarization task using the CNN/Daily Mail dataset, balancing efficient adaptation to the new task with computationally feasible strategies. To achieve this, three key methods were implemented:

1. **Soft Prompt Tuning**: Using a soft prompt embedding layer that enables the model to condition on task-specific prompts without updating the entire model.
2. **Low-Rank Adaptation (LoRA)**: Adding low-rank matrices to specific layers of GPT-2 to introduce task-specific adaptations with minimal parameter updates.
3. **Traditional Fine-Tuning (Last Layers Only)**: Updating only the last classifier layers of GPT-2, allowing the model to retain core language understanding while focusing on summarization.

---

## 2. Architecture Overview

### 2.1 Soft Prompt Embedding Layer (Prompt Tuning)

- **Layer Setup**: An embedding layer was created specifically for the soft prompts, separate from the main GPT-2 model. The embedding layer dimensions are `[num_prompts, embedding_size]`, where `embedding_size = 768` for GPT-2 Small.
- **Functionality**: This layer allows the model to condition on task-specific prompt tokens (e.g., "[SUMMARIZE]") that are prepended to the input sequence. The embeddings for these tokens are learned while the main GPT-2 model parameters remain frozen.

- **Advantages**:
  - **Efficiency**: Only the prompt embeddings are updated, minimizing computational cost.
  - **Flexibility**: Allows reusability of prompts across related tasks.
- **Challenges**:
  - **Expressiveness**: Limited by the number and variability of prompt tokens. May require careful prompt selection and tuning to maximize effectiveness.

## 2.2 LoRA (Low-Rank Adaptation) Architecture

- **LoRA Layer Design**: LoRA introduces two low-rank matrices, AAA and BBB, in specific layers of the GPT-2 model. The matrices have dimensions `[input_dim, rank]` and `[rank, output_dim]`, where `rank` is a hyperparameter controlling the compression factor.
- **Implementation**: The GPT-2 model remains frozen, with task-specific fine-tuning accomplished by updating only the LoRA layers.


- **Advantages**:
  - **Reduced Memory and Computational Requirements**: Low-rank matrices are much smaller than the full weight matrices, resulting in minimal memory usage.
  - **Maintains Model Integrity**: Since only low-rank matrices are updated, core model parameters remain unchanged, preserving pre-trained knowledge while allowing task-specific customization.


- **Challenges**:
  - **Limited Expressiveness**: Although LoRA adapts efficiently, the rank setting may constrain the model's capacity to learn highly nuanced information in the target task.

## 2.3 Traditional Fine-Tuning (Last Layers Only)

- **Layer-Freezing Strategy**: In this method, all parameters are frozen except those in the last few layers (e.g., the pre-classifier and classifier layers). These layers are specifically fine-tuned for the summarization task.
- **Implementation**: Only the `pre_classifier` and `classifier` parameters are allowed to update, while the rest of the model remains frozen. The example code provided achieves this by setting `requires_grad` to `False` for all layers except these last layers.

- **Advantages**:
  - **Computationally Feasible**: Reduces the parameters requiring updates, allowing fine-tuning with lower resource requirements.
  - **Targeted Task Adaptation**: The last layers are trained to map high-level representations to the summarization task effectively.
- **Challenges**:
  - **Expressive Constraints**: Since the updates are limited to the last layers, the model might not fully adapt to complex patterns within the summarization data.

---

## 3. Fine-Tuning and Hyperparameters

**Hyperparameter Choices**

1. **Batch Size**: A larger batch size enhances training stability. For memory-constrained environments, gradient accumulation is employed to effectively increase the batch size.
2. **Epochs**: Up to 10 epochs were used with early stopping based on validation loss, allowing the model to adapt without overfitting.
3. **Gradient Clipping**: Gradient clipping at a norm of 1.0 prevents exploding gradients, particularly useful in large models.

## 4. Training and Implementation Details

During training, each method was applied as follows:

- **Soft Prompt Tuning**: Prompt embeddings were learned by backpropagating gradients only through the soft prompt layer, keeping GPT-2 parameters frozen. This allowed for efficient task conditioning.
- **LoRA Fine-Tuning**: Low-rank matrices were introduced to select layers of GPT-2, capturing task-specific variations without full parameter updates. The pre-trained weights remained frozen, reducing the risk of overfitting.
- **Traditional Fine-Tuning (Last Layers Only)**: Only the last layers were updated, reducing computational load while adapting model output to the summarization task.

## 5. Performance and Evaluation

**Metrics and Dataset**

- **Dataset**: The CNN/Daily Mail dataset was utilized for training and evaluation. It provides a robust benchmark for summarization.
- **Evaluation Metrics**: ROUGE scores (ROUGE-1, ROUGE-2, ROUGE-L) were used to evaluate summarization performance.

**Results and Observations**

1. **Soft Prompt Tuning**:
   - **ROUGE Scores**: Achieved ROUGE scores: {'rouge1': 0.084, 'rouge2': 0.019, 'rougeL': 0.053}
2. **LoRA Fine-Tuning**:
   - **ROUGE Scores**: Achieved ROUGE scores: {'rouge1': 0.254, 'rouge2': 0.19, 'rougeL': 0.21}
   -

3. **Traditional Fine-Tuning (Last Layers Only)**:
   - **ROUGE Scores**: **ROUGE Scores**: Achieved ROUGE scores: {'rouge1': 0.209, 'rouge2': 0.205, 'rougeL': 0.209, 'rougeLsum': 0.209}

## 6. Conclusion

Each fine-tuning method offers unique advantages:

- **Soft Prompt Tuning** is highly efficient and suitable for limited hardware, though it may be less expressive for complex tasks.
- **LoRA Fine-Tuning** balances adaptation capability with computational efficiency, achieving strong task performance with minimal memory requirements.
- **Traditional Fine-Tuning (Last Layers Only)** allows for computationally inexpensive tuning but may lack the expressiveness of LoRA for intricate tasks.

## Future Considerations

- **Rank Optimization in LoRA**: Experimenting with different ranks may optimize performance further.
- **Hybrid Methods**: Combining LoRA and soft prompts could allow adaptable yet efficient task conditioning.
- **Extending to Multi-Task Learning**: Given LoRA's modular design, the model could adapt to multiple summarization domains by loading different low-rank matrices for each.

This architecture presents a versatile framework for deploying task-specific adaptations of GPT-2, balancing efficiency and performance across different fine-tuning methods.