

Lightweight Transformer-Based Context-Aware Text Generation

Guangyan An*
Department of CISE
University of Florida
Gainesville, USA
anguangyan@ufl.edu

Xiaolin Zheng*
Department of CISE
University of Florida
Gainesville, USA
xiaolinzheng@ufl.edu

I. INTRODUCTION

Text generation systems play a critical role in improving user efficiency in applications like typing assistants, chatbots, and news summarization tools. Traditional methods, often based on n-gram or rule-based models, struggle with understanding complex sentences and providing accurate suggestions, especially in scenarios requiring nuanced context or real-time performance.

The emergence of transformer-based models has transformed Natural Language Processing (NLP). Models like BERT, BART, and GPT-2 leverage bidirectional context, enabling more accurate predictions. However, despite their effectiveness, transformers are computationally intensive, which can limit their feasibility in real-time applications like typing assistance. Additionally, token-based transformer models, such as GPT-2, often struggle with incomplete input (e.g., partial words), reducing their utility for text completion tasks.

Our project initially focused on creating a lightweight, context-aware text completion and summarization system by leveraging pre-trained transformers optimized through techniques such as knowledge distillation, movement pruning, and quantization. While summarization tasks demonstrated promising results, we observed significant limitations in using transformer models like GPT-2 for real-time typing assistance, particularly for incomplete word completions. These challenges motivated a pivot to a hybrid approach: combining a lightweight, character-level LSTM model for word completion with GPT-2 for phrase predictions. This hybrid design aims to balance efficiency, accuracy, and real-time performance.

This report outlines the development and outcomes of this dual-focus project. Section II discusses related work on model optimization techniques, including knowledge distillation, pruning, and quantization. Section III describes the methodologies used for optimizing transformers, the LSTM-based character model, and the integration of both into a hybrid system. Section IV presents the evaluation results for both summarization and typing assistance tasks, with a focus on the hybrid system's performance. Finally, Section V concludes the project and suggests future improvements to enhance both summarization and real-time text completion capabilities.

II. RELATED WORK

A. Knowledge Distillation

Hinton et al. [1] introduced the concept of knowledge distillation, providing the basis for model compression and optimization. Knowledge distillation uses soft targets of the teacher model as the learning target for the smaller student model by minimizing the KL divergence between the output probability distributions of the teacher and the student. Building on this, FitNets [2] proposed using intermediate representations, such as feature maps and hidden states, of the teacher model as a guidance for students to learn internal representations. Students not only imitate the final outputs of the teachers, but also replicate their details. Then, DistilBERT [3] applied this technique to BERT, a Transformer model, reducing its parameters by 40% percent while retaining more than 97% of the performance on various workloads. Subsequently, Shleifer and Rush [4] extended knowledge distillation techniques to the field of text summarizations. Their approach introduced task-related techniques to distill large summarization models, such as T5 and BART, into smaller models, with a smaller model inference time. These researches demonstrate that by effectively extracting and transferring knowledge from different layers of the teacher model, knowledge distillation can significantly reduce model footprint and inference time, making models adaptable to various application scenarios.

B. Pruning

Early pruning strategy focused on sparsifying weight matrices by removing weights with small values, which is unstructured pruning [5]. For Transformers, structured pruning is applied to optimize the self-attention and feed-forward neural network components. This approach calculates importance scores to select necessary attention heads and drop the others [6] without significantly degrading the model performance. After that, movement pruning [7] was proposed to dynamically adjust the pruning strategy during training or fine-tuning process according to gradient information, improving the overall model precision and robustness after pruning. This strategy can be applied to both unstructured and structured pruning, which is flexible. Famous real-world model pruning implementations include `torch.nn.utils.prune`, which

* These authors contributed equally to this research.

contains both unstructured and structured pruning. For unstructured pruning, the implementation leverages random masks or weight thresholds to drop weights in a parameter-by-parameter manner. Structured pruning is a larger-granularity pruning that removes neurons, convolution kernels, or attention heads.

C. Quantization

Quantization includes weight quantization and activation quantization. In recent years, 8-bit weight and activation quantization has become mainstream in Transformer models, since quantizing a 32-bit float point Transformer to an 8-bit model can largely reduce the model footprint and inference cost while retaining high precision. For instance, Q8BERT [8] introduced a quantized 8-bit BERT model, reducing the model size by approximately 4x by applying quantization-aware training. Furthermore, TernaryBERT [9] uses a ternary quantization mechanism, and the weights of the model only take values in $\{-1, 0, 1\}$, further achieving a higher compression rate while retaining the performance. Real-world implementations of quantization include TorchQuantization, which provides implementations of post-training quantization and quantization-aware training. Post-training quantization requires the offline collection of calibration data to generate precise scale factors after converting the model. For quantization-aware training, TorchQuantization supports inserting fake quantization modules in the model and significantly improves the performance of low-precision models by simulating and propagating quantized model loss to optimize the weights of the quantized model.

III. METHOD

A. Hybrid Model Design for Text Prediction

To address the limitations of token-based models like GPT-2 in handling incomplete inputs and achieving real-time performance, we implemented a hybrid system that combines a lightweight, character-level LSTM model with GPT-2. This hybrid design leverages the strengths of both models, enabling efficient and accurate predictions for typing assistance tasks.

The LSTM model focuses on character-level predictions, which are particularly effective for completing partial words. It takes an input sequence of characters, processes them through recurrent layers, and predicts the most likely next character at each timestep. For example, given an input like “ligh”, the LSTM predicts the following characters (‘t’, ‘w’, ‘e’, ‘i’, ‘g’, ‘h’, ‘t’) iteratively until a space character or punctuation is encountered.

The prediction process involves:

- 1) Converting each character into an index in the vocabulary (English alphabet, numbers, and special characters), and embedding the indices into a dense representation.
- 2) Processing the embedded sequence through LSTM layers that retain context over the input.
- 3) Generating a probability distribution over the vocabulary for the next character at each timestep, selecting the character with the highest probability.

This iterative mechanism allows the LSTM to generate highly accurate character-level completions, even for incomplete inputs.

In the hybrid system, inputs are first processed by the LSTM. If the LSTM produces a non-empty result (non-space or punctuation character), it is returned directly to the user. Otherwise, the input is passed to GPT-2, which provides context-aware phrase completions or generates predictions for longer, more complex inputs. By delegating simpler tasks to the LSTM, the system achieves lower latency and reduces computational load on GPT-2.

B. Knowledge Distillation

Knowledge Distillation is a technique for transferring knowledge from a large, complex model (teacher model) to a smaller, more efficient model (student model). The key idea is that the student model can achieve comparable performance to the teacher model by learning soft labels of the teacher: probability distributions instead of learning the true labels only. This allows the student to learn intermediate representations from the teacher. The loss function of knowledge distillation can be expressed as:

$$L_{KD} = \alpha L_{CE}(y, \hat{y}) + (1 - \alpha) T^2 L_{KL}(q_t, q_s) \quad (1)$$

where $L_{CE}(y, \hat{y})$ represents the loss between prediction probability distribution \hat{y} and true labels y using traditional cross entropy loss. $L_{KL}(q_t, q_s)$ expresses the difference between the output distributions of the student and the teacher. T is the temperature hyperparameter, used to soften softmax distributions of q_t and q_s . Finally, α controls the balance between different kinds of losses.

BART-large-CNN [10] was pre-trained and fine-tuned on the CNN Daily Mail dataset and has a powerful ability of text summarization. It has 12 Transformer encoder layers and 12 Transformer decoder layers, the model being in total of 1.51GB. GPT-2 [11] was a pre-trained autoregressive decoder model using a causal language masking objective, containing 12 Transformer decoder layers with a total size of 0.548 GB. We extracted the layers 0, 2, 5, 7, 9, 11 from the encoder/decoder to form the student’s encoder/decoder, thereby halving the model’s storage. Weight state dicts of these layers are also loaded to avoid pre-training from scratch.

In our implementation, we chose to use MSE to match the logits output at each hidden layer as described in [4] since this can encourage the student to even match teacher hidden states. The loss function is updated as follows, where we want the student to generate hidden layers weights as comparable to the teacher as possible:

$$L_{KD} = \alpha L_{CE}(y, \hat{y}) + \beta L_{MSE}(hid_t, hid_s) + \gamma T^2 L_{KL}(q_t, q_s) \quad (2)$$

For the α, β, γ, T hyperparameters, we set α to 0.1, β to 3, γ to 0.8, and T to 2, since we aim to maximize the hidden state learning by the student from the teacher.

As shown in Figure 1, this image illustrates an example encoder-decoder Transformer (e.g. BART) teacher model with

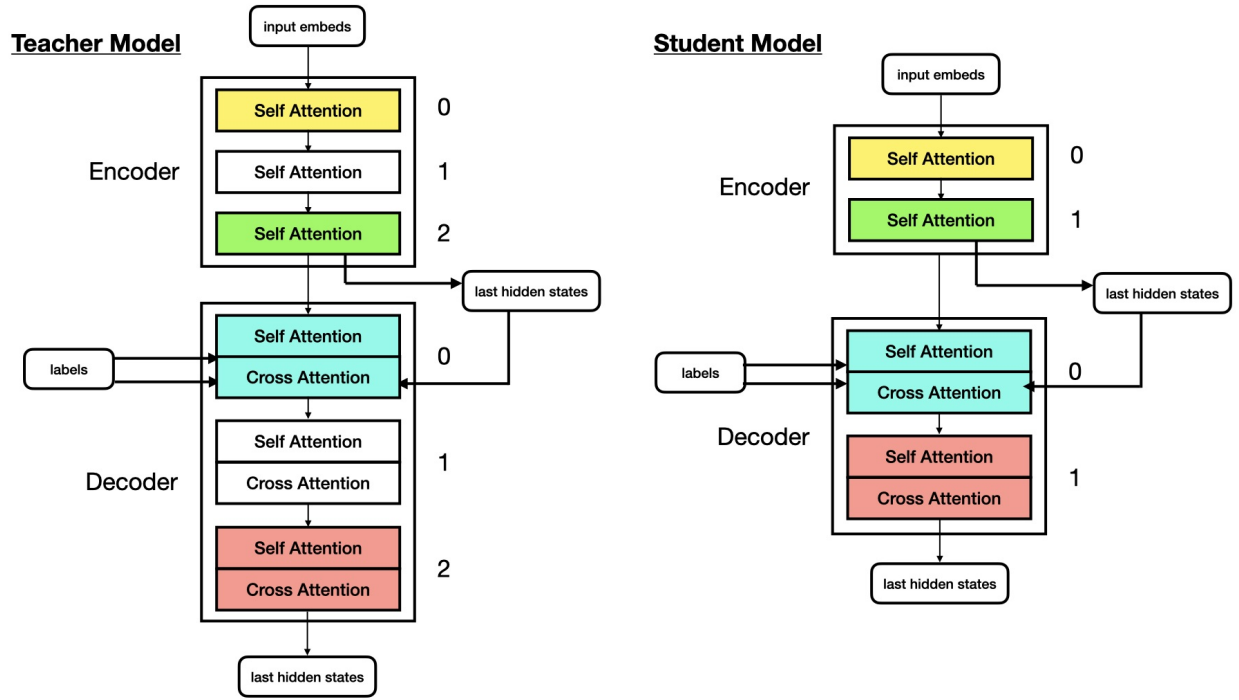


Fig. 1: Example of distillation of Encoder-Decoder model, e.g. BART-large

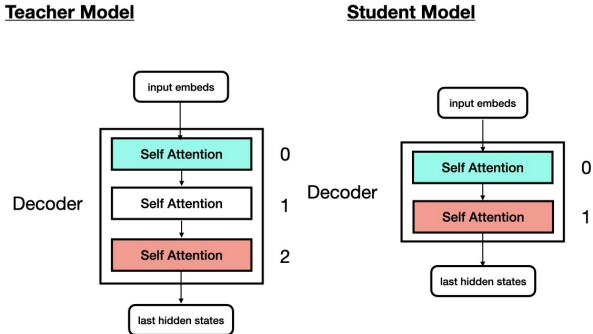


Fig. 2: Example of distillation of Decoder-Only model, e.g. GPT-2

3 layers for each on the left. On the right a distilled student model extracts the first and last layers of the teacher’s encoder and decoder (using the same colors) and forms a new encoder and decoder. Arrows in this figure show the data flow: `input embeds` are the text embedding input, which will become Q, K, V in encoders. Encoder’s last hidden layer states will be transformed to K, V used in cross-attention of decoder layers, and the `labels` are fed as Q, K, V in self-attention but as Q only in cross-attention of decoder layers.

Figure 2 presents an example 3-layer decoder-only Transformer (e.g. GPT-2) teacher model on the left, which is much simpler than the previous one since decoder-only models can choose not to use cross-attention, meaning it only leverages self-attention. On the right a distilled student model extracts the first and last layers of the teacher’s decoder (using same

colors) and forms a new decoder. `input embeds` will become Q, K, V in decoders. Finally the final logits generated by the decoder will be compared with the original text embeddings using cross entropy to calculate the loss because of autoregression.

C. Movement Pruning

Pruning is a model optimization technique aimed at reducing the size and computational complexity of deep learning models by removing less important parameters or structures. There are several different classification approaches for pruning, including by granularity (weight, neuron, channel, layer), by timing (pre-training, during training, post-training), and by target (structural, unstructural). Our pruning strategy focuses on neurons in FFN layers and attention heads in QKV projection layers. Specifically, we train pruning masks for both FFN and QKV projections separately, and finally perform pruning using the masks. Therefore, our approach can be categorized as structural movement pruning, taking neurons and attention heads as the minimum granularity.

We assign a uniformly drawn score to each neuron in an FFN layer and to each head in an attention layer. Each time during forward passes, we use a threshold K to select the top K smallest scores and drop the corresponding neurons and attention heads by transforming the scores into masks. Regarding neurons and attention heads as pruning granularity, we should notice that once we have determined which attention head is useless, we should remove that same head from all QKV projections and the final output projection in the current attention. Also for FFN layers that contain two linear layers

`fc1` and `fc2`, once we have decided which neuron can be pruned in `fc1`, we must remove the corresponding neuron defined in `fc2`. Therefore, we call our trained scores “shared score” or “shared mask”, since they are shared among QKV attention projections, and among FFN layers.

During the training process in order to find the best shared masks in the Transformer, we will apply the current shared masks in forward and backward passes to train learnable shared masks. However, the masking operation is not differentiable since it zeros out all corresponding attention heads and set them to zero as inactive heads during forward passes. To solve this, we applied STE, which stands for straight-through estimator techniques to estimate gradients during backward passes so that the scores can be updated properly. Therefore, the most important neuron or attention heads must have the highest score theoretically due to natural selection. Additionally, we implemented a dynamic threshold scheduler for threshold T for the movement pruning trainer, as shown in Eq. 3. For instance, in our setup, we aim to prune 30% of the attention heads in the Transformer and 60% of the neurons in FFN layers. To minimize abrupt performance degradation, the scheduler gradually increases the total masking threshold from 0% to 30% following a cubic curve for masked attention heads. The thresholds will reach their peak at 85% of the training progress and stay at the peak until the end and use the remaining 15% of the progress for fine-tuning, as shown in Figure 3.

$$T = T_{\text{init}} + (T_{\text{final}} - T_{\text{init}}) \times (1 - \text{progress})^3 \quad (3)$$

Figures 4 and 5 show pruning results for an FFN layer and a decoder layer, respectively. In Figure 4, darker colors indicate the neurons that have been pruned, while in Figure 5, darker colors highlight pruned self-attention and cross-attention heads.

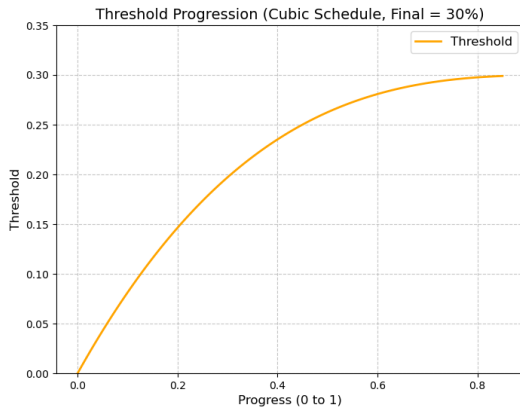


Fig. 3: Movement Threshold Scheduler for Attention Layers

D. Quantization

Quantization can reduce the precision of numbers used to present parameters or activations. By converting high-

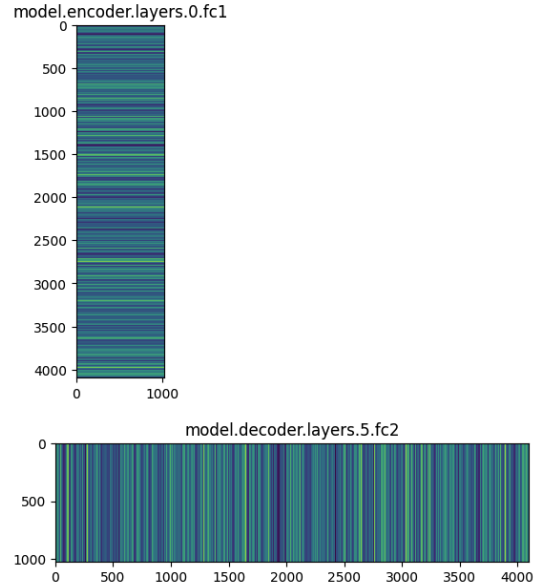


Fig. 4: Sample Neuron Pruning for FFN

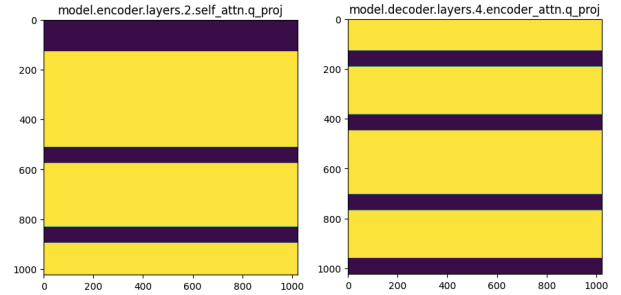


Fig. 5: Sample Attention Head Pruning for Self-Attention and Cross-Attention in Decoder.

precision floating-point numbers (FP32, FP64) to lower precision (INT8, FP16), quantization reduces the memory footprint and computational requirements of a model, making it more efficient for deployment on resource-constrained devices like edge or mobile devices. Types of quantization include post-training quantization (PTQ), and quantization-aware training (QAT). To map FP32, for instance, to INT8, we can convert floating point values into quantized space using linear mapping strategy, as expressed below:

$$\begin{aligned} r &= S(q - Z) \\ q &= \text{round} \left(\frac{r}{S} + Z \right) \end{aligned} \quad (4)$$

where r and q are the number before and after quantization; S and Z are scale and zero-point. The linear mapping is referred to as symmetric or asymmetric mapping, depending on whether Z is zero, as shown in Figure 6, where a symmetric range is converted to a symmetric INT8 range on the left and an asymmetric range is mapped to an asymmetric INT8 range.

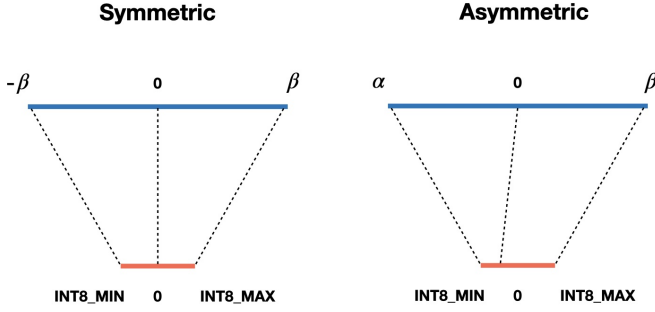


Fig. 6: Linear Mapping with or without Zero-point

After mapping FP32 to INT8, which is PTQ, the memory footprint is significantly reduced; however, this comes at the cost of precision and performance degradation. To compensate for the degradation, calibration can be performed by feeding real data into the quantized model to estimate and adjust weight ranges.

Another approach for the model to adjust for quantization errors is to use online quantization techniques, for instance, QAT, instead of the offline PTQ method. QAT involves training the model with quantization simulated during the process to improve performance after quantization. Despite weights being quantized to INT8, during the forward process weights are still represented as FP32 to adjust for quantization errors. After QAT, the model learns to compensate for quantization errors, ensuring minimal performance degradation, followed by PTQ after training to convert the model to an INT8 format, and inference is performed using INT8 precision.

Our approach follows the above description, performing QAT first, followed by PTQ. To implement QAT, We designed a `QLinear` layer to substitute the original `torch.nn.Linear` layers to perform quantization-aware training. There was a similar problem as described in Subsection III-C: the linear mapping process, as we call it *affine*, leveraging `round` and `clamp` operations, is non-differentiable. We applied the same approach described above: STE to estimate the backward gradient for our `QLinear` layers to learn and gradually adjust for quantization errors. We rewrote the backward function for the $XW + b$ linear operation, plugging in *affine* and *deaffine* logic for W parameters during forward and backward passes. Another operation worth mentioning is that we did not assign zero-points in *affine* and *deaffine* since during experiments we found the performance of the models was better using symmetric quantization, possibly due to weight distribution in the pre-trained models.

IV. EXPERIMENT RESULTS

A. Experiment Setup

We conducted experiments using the OpenWebText dataset [12], an open-source replication of OpenAI’s WebText, for autoregressive text prediction tasks and training the LSTM model for character-level text completion. The BBC

news dataset, containing original posts and their corresponding summaries, was used for text summarization tasks with Transformer-based models.

For the Transformer-based models, we utilized HuggingFace’s *transformers* [13] library, which provided pre-trained definitions for BART and GPT-2. The library’s `Trainer` and `Seq2SeqTrainer` APIs were used for fine-tuning, streamlining the training process. Experiments were run on an A100 GPU with 40GB of memory using Google Colab. Due to computational constraints, only a small subset of the original datasets was used, requiring approximately 5 hours of runtime. The pre-trained models were fine-tuned on task-specific data to facilitate performance comparisons post-optimization.

For the LSTM model, a subset of OpenWebText dataset (10K samples) was preprocessed into overlapping character-level sequences. Texts were split into chunks of up to 100 characters, with a one-third overlap between consecutive chunks to preserve context. The vocabulary included lowercase and uppercase English letters, punctuation symbols (e.g., ., ,, !, ?), and spaces, with unknown or non-standard characters handled as special tokens.

We trained the LSTM model with a batch size of 1024. Each sequence was split into inputs (characters up to the second-last in the chunk) and targets (characters from the second character onward). This setup allowed the LSTM to learn character-level dependencies effectively and generate predictions for incomplete text inputs. On an Nvidia 4090 GPU, it took about 2.5 hours to train the model.

B. Knowledge Distillation

We have evaluated the performance, inference time, and the reduction of storage for our distiller with the settings discussed above. As the table I shows below, the distiller for BART-large-CNN on text summarization had a good result: not only did our student have nearly the same rouge scores as the teachers’, but it also had 26.7% inference time reduction and 43% model size reduction. The rouge scores range from 0 to 1, with higher values indicating better quality. Please note that the model size reduction is less than 50% but that is normal since Transformer models have a shared embedding layer above the encoder to transform text ids to text embeddings, and another linear language modeling head layer below the decoder to map hidden states to vocabulary logits.

Name	BART-large-CNN	
	Teacher Model	Student Model
Encoder Layers	12	6
Decoder Layers	12	6
ROUGE-1 Score	0.717	0.716
ROUGE-2 Score	0.617	0.617
ROUGE-L Score	0.541	0.551
ROUGE-Lsum Score	0.559	0.604
Inference Time	59.25s	43.43s
Model Size	1.51GB	0.86GB

TABLE I: Comparison of Teacher and Student Models for BART-large-CNN

The result of distiller for GPT-2 are shown below in the table II. Since we use GPT-2 for text prediction and it is an autoregressive decoder-only model, perplexity was used to evaluate model performance. The lower the perplexity score, the better the model’s predicted probability distribution matches the true data distribution, indicating higher quality in generated text under certain conditions. The table suggests that the perplexity score downgraded from 18.25 to a “more perplexing” 31.98, which we believe it is due to the small size of the dataset, and because on text generation we did not train it to full convergence due to a lack of resources. However, the inference time shows 19% reduction, with the model reduced by 43% (the latter one is same as the result from BART-large-CNN).

GPT-2		
Name		
Metric	Teacher Model	Student Model
Decoder Layers	12	6
Perplexity	18.25	31.98
Inference Time	43.18s	34.95s
Model Size	0.548GB	0.31GB

TABLE II: Comparison of Teacher and Student Models for GPT-2

C. Movement Pruning

We applied only our pruner on the fine-tuned BART-large-CNN model without applying distillation to show pruning results. The table III below shows 35% of storage reduction and about 4% of rouge score loss, which means our pruning strategy works. A reflection on the pruning strategy is that, we use learnable scores (masks), train those masks, and apply those masks during training to get the most valuable neurons / attention headers, which has similar semantics as QAT using quantization-aware weights during training. Besides, the STE strategy is applied as well both in pruning and in QAT.

BART-large-CNN		
Name		
Metric	Original	Pruned
Encoder Layers	12	12
Decoder Layers	12	12
ROUGE-1 Score	0.717	0.713
ROUGE-2 Score	0.617	0.612
ROUGE-L Score	0.541	0.530
ROUGE-Lsum Score	0.559	0.535
Inference Time	59.25s	59.2s
Model Size	1.51GB	0.98GB

TABLE III: Comparison of Original and Pruned Models for BART-large-CNN

The next table IV shows the pruning result of the GPT-2 model, with a 25% reduction in storage and an increase in perplexity from 18.25 to 39.84. The reduction of model size reduction is smaller for GPT-2 compared to BART-large-CNN, as decoder layers in BART-large-CNN include both self-attention and cross-attention, which can be pruned, whereas only self-attention is prunable in GPT-2’s decoder layers. Hence the storage reduction on BART-large-CNN is greater.

GPT-2		
Name		
Metric	Original	Pruned
Decoder Layers	12	12
Perplexity	18.25	39.84
Inference Time	43.18s	42.1s
Model Size	0.548GB	0.406GB

TABLE IV: Comparison of Original and Pruned Models for GPT-2

D. Quantization

We applied only our quantizer on models, and the table III below shows 62% of storage reduction with a loss of less than 2% on the BART-large-CNN model. However, the inference time remains the same: after PTQ weights are transformed into quantized INT8 formats, and the model should have applied INT8 by INT8 matrix multiplication when doing linear calculation; however, it is challenging to find highly optimized INT8 by INT8 matrix third-party APIs both working on CUDA chips and MPS devices (Macbook M2). To keep the performance data consistent on both platforms, we employed a straightforward approach by converting INT8 to FP32 during inference. While the results remain consistent, this approach is slower. The results of GPT-2 are comparable to those of BART-large-CNN; therefore, they are not presented here for brevity.

BART-large-CNN		
Name		
Metric	Original	Quantized
Encoder Layers	12	12
Decoder Layers	12	12
ROUGE-1 Score	0.717	0.717
ROUGE-2 Score	0.617	0.616
ROUGE-L Score	0.541	0.536
ROUGE-Lsum Score	0.559	0.552
Inference Time	59.25s	59.3s
Model Size	1.51GB	0.58GB

TABLE V: Comparison of Original and Quantized Models for BART-large-CNN

E. Pipelines

We combined these optimizations in a distillation-pruning-quantization order as a pipeline and performed the final optimized models on text summarization and text prediction. After each optimization was done, we performed model size calculation and performance measurement for that optimization to observe the effect.

The table VI presents the results of the BART-large-CNN pipeline, which reduces the memory footprint by 77.5% and inference time by 26.2%, with a performance loss of less than 3.5%. However, there is still more space to optimize storage and inference time. For example, we can extract 3 or 4 layers from the teacher during distillation instead of 6, drop attention heads by more than 30% as we set, or implement INT8 by INT8 matrix multiplication logic for quantized model inference, but we would also like to balance the loss to get

a good performance figure so we made some trade-offs on hyperparameter selections. From the table, we may find that the attention layers and the FFN layers have been sufficiently optimized, and the remaining storage is mainly due to the existence of the shared embedding layer and the language modeling head layers, which are both huge with a size of (50264×1024) in BART-large-CNN, indicating that these two layers, i.e. the vocabulary size is the current bottleneck. In comparison, unoptimized Q projection in BART-large-CNN only has a size of (1024×1024) . One way to handle this is to prune the vocabulary table: we can scan all data from datasets and drop never-used words from the vocabulary table. We researched and found using our BBC dataset, we could shrink the size of the vocabulary table to nearly half of the original. However, since this is highly dataset-related, so we finally did not apply it.

For the GPT-2 pipeline, the results are as the below table VII, showing that the pipeline reduces the memory footprint by 69% and inference time by 21.5%, with an increase of perplexity of 47.41%. Therefore, the results in BART-large-CNN and GPT-2 optimizations are consistent and effective.

F. Task-Specific Evaluation

To provide a practical perspective on model performance, we conducted evaluations specifically designed for real-world tasks in a typing assistance context. While perplexity is commonly used to measure model quality, it does not directly assess usability, particularly for incomplete inputs. To address this, we devised two tests:

The **Incomplete Word Test** involved truncating the last word of an approximately 60-character sentence to a partially typed state, requiring the models to predict the correct completion. The **Complete Word Test** removed the last full word from the same sentence, tasking the models to predict the missing word. Each test used 1,000 samples derived from the OpenWebText-100k dataset.

The results, as shown in Table VIII, demonstrate distinct strengths for each model. For incomplete word predictions, the LSTM achieved 62.24% accuracy, far surpassing GPT-2's 5.51%. This aligns with the LSTM's design, which focuses on character-level dependencies and is particularly effective for partially typed inputs. GPT-2's poor performance in this test highlights its reliance on token-based predictions, which are less suited for incomplete inputs.

Conversely, in the complete word test, GPT-2 outperformed the LSTM with 14.87% accuracy compared to 8.71%. This reflects GPT-2's ability to leverage its contextual understanding and token-based architecture for generating complete words in structured inputs. However, the relatively low accuracy for both models in this test indicates room for improvement in handling broader word prediction tasks.

These results justify the hybrid design, leveraging the LSTM's strengths for incomplete inputs and GPT-2's contextual capabilities for complete inputs. By evaluating models in a task-specific manner, we bridge the gap between standard

metrics like perplexity and real-world usability, ensuring the system performs effectively in practical scenarios.

V. CONCLUSIONS AND FUTURE WORK

The experiments and evaluations conducted in this study demonstrate the trade-offs inherent in optimizing lightweight models for practical tasks like typing assistance and summarization. While our current optimizations and pipelines have shown measurable improvements in model size and inference time, several avenues remain to enhance performance and usability further.

One promising direction is exploring frameworks such as ONNX, which can streamline operations like INT8 matrix multiplication and improve runtime compatibility across diverse hardware. Integrating ONNX into our workflow could significantly reduce inference latency for quantized models. Another area of improvement lies in better utilization of PyTorch's capabilities. While Python's inherent performance overhead has been mitigated to some extent by PyTorch's C++ backend, leveraging Just-In-Time (JIT) compilation and converting models to TorchScript can further eliminate Python-specific bottlenecks. This approach would ensure faster execution during inference without requiring extensive changes to existing training workflows. For fine-tuning, emerging techniques like Low-Rank Adaptation (LoRA) offer a potential path forward. By splitting weight matrices into low-rank components, LoRA reduces the computational cost of gradient updates during fine-tuning while maintaining high precision. Applying such methods to the LSTM and GPT-2 models may lead to faster convergence and better resource utilization. Finally, future efforts should investigate optimizing the hybrid system itself. For example, dynamic task allocation mechanisms could be developed to improve the decision-making process between LSTM and GPT-2, ensuring optimal resource use based on input complexity. Additionally, refining the LSTM to handle broader contexts or partially fine-tuning GPT-2 on character-level tasks might further enhance the system's adaptability.

In conclusion, advancing lightweight Transformer models requires a balanced approach, combining innovative algorithmic strategies with practical implementation refinements. By leveraging emerging technologies and techniques, the potential for creating even more efficient and capable systems is substantial.

REFERENCES

- [1] Geoffrey Hinton. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- [2] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. *arXiv preprint arXiv:1412.6550*, 2014.
- [3] V Sanh. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [4] Sam Shleifer and Alexander M Rush. Pre-trained summarization distillation. *arXiv preprint arXiv:2010.13002*, 2020.
- [5] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28, 2015.

Pipeline	Fine-Tuned	Distilled	Pruned	Quantized	Percentage (Q / FT)
Encoder Layers	12	6	6	6	-
Decoder Layers	12	6	6	6	-
ROUGE-1 Score	0.717	0.716	0.712	0.707	98.6%
ROUGE-2 Score	0.617	0.617	0.611	0.604	97.8%
ROUGE-L Score	0.541	0.551	0.528	0.523	96.6%
ROUGE-Lsum Score	0.599	0.604	0.578	0.579	96.6%
Inference Time (s)	59.25	43.42	43.44	43.75	73.8%
Model Size (GB)	1.51	0.86	0.56	0.34	22.5%

TABLE VI: Pipeline of BART-large-CNN

Pipeline	Fine-Tuned	Distilled	Pruned	Quantized	Percentage (Q / FT)
Decoder Layers	12	6	6	6	-
Perplexity	18.25	31.98	68.27	65.66	+47.41
Inference Time (s)	43.18	34.95	34.96	33.92	78.5%
Model Size (GB)	0.548	0.31	0.23	0.17	31.0%

TABLE VII: Pipeline of GPT-2

Test Type	LSTM Accuracy (%)	GPT-2 Accuracy (%)
Incomplete Word Test	62.24	5.51
Complete Word Test	8.71	14.87

TABLE VIII: Accuracy results for task-specific evaluations of LSTM and GPT-2 models.

- [6] Paul Michel, Omer Levy, and Graham Neubig. Are sixteen heads really better than one? *Advances in neural information processing systems*, 32, 2019.
- [7] Victor Sanh, Thomas Wolf, and Alexander Rush. Movement pruning: Adaptive sparsity by fine-tuning. *Advances in neural information processing systems*, 33:20378–20389, 2020.
- [8] Ofir Zafrir, Guy Boudoukh, Peter Izsak, and Moshe Wasserblat. Q8bert: Quantized 8bit bert. In *2019 Fifth Workshop on Energy Efficient Machine Learning and Cognitive Computing-NeurIPS Edition (EMC2-NIPS)*, pages 36–39. IEEE, 2019.
- [9] Wei Zhang, Lu Hou, Yichun Yin, Lifeng Shang, Xiao Chen, Xin Jiang, and Qun Liu. Ternarybert: Distillation-aware ultra-low bit bert. *arXiv preprint arXiv:2009.12812*, 2020.
- [10] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*, 2019.
- [11] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, Ilya Sutskever, et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.
- [12] Aaron Gokaslan and Vanya Cohen. Openwebtext corpus, 2019.
- [13] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing: system demonstrations*, pages 38–45, 2020.

APPENDIX A

LINK TO CODE REPOSITORY

The GitHub repository containing the code is available at [Phoslight/CAP6617-Project](https://github.com/Phoslight/CAP6617-Project).