# Tiny Titans: Achieving High Accuracy in Small Models through Knowledge Distillation

Sahithi Singireddy
ssingireddy@umass.edu

Sravani Gona
sgona@umass.edu

## Abstract

*Deploying large neural networks on resource-constrained devices, like smartphones and IoT systems, poses challenges due to high memory and computational demands. This project leverages knowledge distillation, a model compression technique, to address these constraints. By transferring knowledge from a large, high-performing teacher model to a smaller, efficient student model, we aim to maintain performance parity while significantly reducing resource demands. Our implementation incorporates two key techniques: logit matching, where the student model learns to replicate the soft predictions of the teacher, and feature matching, which involves aligning intermediate representations between the teacher and the student. Using the CIFAR-10 and ImageNet datasets, we evaluate the student model's perfor- mance across metrics like accuracy, inference speed, and memory usage. Our approach aims to demonstrate that knowledge distillation can effectively balance model effi- ciency and accuracy, providing a viable solution for mobile AI and edge computing applications.*

## 1. Introduction

The project, titled "Tiny Titans: Achieving High Accuracy in Small Models through Knowledge Distillation," addresses the growing need for efficient deployment of neural networks on resource-constrained devices like smartphones and IoT systems. Modern deep learning models such as ResNet and VGG have demonstrated exceptional accuracy and performance across tasks like image recognition, natural language processing, and speech analysis. However, their high computational requirements, memory footprint, and storage needs limit their usability for real-time applications on mobile and edge devices.

To overcome these limitations, the project leverages knowledge distillation (KD)—a technique where a smaller model (student) learns to replicate the performance of a larger pre-trained model (teacher). The smaller model is trained to mimic the outputs and feature representations of the teacher model, allowing it to achieve comparable ac-curacy while being much lighter and computationally efficient.

Using the CIFAR-10 dataset, which consists of 60,000 images across 10 classes, we will train the student model to replicate the teacher model's behavior. Our approach includes logit matching, where the student mimics the teacher's softmax output, and feature matching, where the student replicates the teacher model's intermediate representations.

In addition to CIFAR-10, we will also scale our approach to the ImageNet dataset, which is significantly larger and more complex, containing over 14 million images across 1,000 categories. Testing on ImageNet will allow us to assess the scalability and robustness of our knowledge distillation method in a more challenging environment, thereby evaluating the model's generalization capabilities when faced with diverse and intricate image data.

To evaluate our approach, we will measure the student model's performance based on accuracy, compression ratio, inference time, and memory usage. We anticipate that our results will demonstrate that knowledge distillation can provide an optimal balance between model efficiency and accuracy, making it a viable solution for deploying high-performing models on devices with limited resources.

## 2. Related work

Knowledge Distillation (KD) has become a crucial area of research in deep learning, aiming to reduce the size and complexity of models while preserving their accuracy. This approach has gained prominence due to the increasing demand for deploying AI systems on resource-constrained devices such as smartphones, IoT devices, and embedded systems. These environments often have stringent limitations in terms of computational power, memory, and energy efficiency, making the deployment of large, high-performance neural networks infeasible. By enabling smaller models to inherit the capabilities of larger models, KD provides a practical solution for democratizing AI and extending its reach to edge devices. This section reviews foundational work and subsequent advancements in KD, which inform the design and implementation of our project.

The concept of knowledge distillation was first introduced by Hinton et al. (2015) in the paper "Distilling the Knowledge in a Neural Network" [2], establishing the foundational technique of logit matching. In this method, a smaller student model learns to replicate the behavior of a larger teacher model by mimicking its softened output probabilities rather than relying solely on hard labels.

Key innovations in this approach include:

- **Softened Outputs**: By introducing a temperature parameter to the softmax function, the teacher's logits are smoothed, revealing richer inter-class relationships. For example, in a task involving image classification, a high-confidence prediction of "cat" may still retain a non-negligible probability for "tiger," reflecting their semantic similarity.
- **Learning Rich Relationships**: This additional information allows the student model to generalize better, capturing the nuanced decision-making of the teacher.

This technique is both simple and effective, forming the baseline for many subsequent advancements in KD. It has been successfully applied across tasks such as image classification, language modeling, and object detection, demonstrating its versatility and robustness. In our project, we adopt Hinton's logit matching as the baseline method to help the student model on the CIFAR-10 dataset learn nuanced predictions from the teacher, ensuring a solid foundation for the model's training and performance.

Ba and Caruana (2013) were among the first to explore the concept of compressing neural networks by training shallow student models that replicate the behavior of deeper models [1]. Their work on model compression laid the groundwork for subsequent KD methods, demonstrating that shallow networks could perform well when trained on teacher-generated soft targets, paving the way for research into more efficient, resource-constrained model deployment.

In addition to logit matching, feature matching has emerged as a crucial component in distillation, particularly for tasks involving complex feature representations. Romero et al. (2014) proposed FitNets: Hints for Thin Deep Nets [5], where intermediate layer activations of the teacher are transferred to the student model, improving performance in shallower student networks. By aligning the intermediate features between the teacher and student, this approach allows the student to capture rich hierarchical features, which is especially beneficial in convolutional neural networks (CNNs). Intermediate layers in CNNs encode important visual features such as edges, textures, and object parts, which are essential for accurate image classification. FitNets has been widely recognized for enabling smaller student models to learn effective feature representations from the teacher, leading to improved performance even when the student network has fewer parameters.

Zagoruyko and Komodakis (2016) introduced Attention Transfer, focusing on transferring spatial attention maps from teacher models to student models [8]. This method aligns the attention mechanisms, enabling the student to focus on relevant features in the data. This approach has shown significant promise in vision tasks, such as object detection and segmentation.

Additionally, Yim et al. (2017) proposed "A Gift from Knowledge Distillation: Fast Optimization, Network Minimization, and Transfer Learning" [7], which extends feature matching by exploring the transfer of both low-level and high-level features between teacher and student models. Their approach, which focuses on transferring visual features at different layers, enhances the student's ability to generalize across various image classification tasks. This method also demonstrates the value of intermediate layer matching in improving performance, especially in scenarios where the teacher network is much larger and more complex than the student.

TinyBERT (Jiao et al., 2019) [3], an extension of KD for natural language processing (NLP) models, introduced the concept of multi-level distillation, which goes beyond traditional logit matching. In addition to mimicking the teacher's logits, TinyBERT aligns the feature representations of intermediate layers between the teacher and the student. This approach ensures that the student captures both the high-level outputs and the hierarchical abstractions encoded in the teacher's intermediate layers.

Key contributions of TinyBERT include:

- **Intermediate Feature Matching**: By aligning feature maps at selected layers, the student learns deeper abstractions that improve its understanding of the data.
- **Relevance to CNNs**: Although designed for NLP, TinyBERT's methodology is highly applicable to convolutional neural networks (CNNs), where intermediate layers often encode critical visual features, such as edges, textures, and object parts, which are essential for accurate image classification.

Drawing inspiration from TinyBERT, our project integrates feature matching between selected intermediate layers of the teacher and student models. This enhancement enables the student to learn richer feature representations, complementing the traditional logit-based KD. By incorporating both methods, we aim to achieve a balance between high accuracy and computational efficiency, making the student model suitable for resource-constrained environments.

Tian et al. (2019) introduced Contrastive Representation Distillation (CRD) [6], which bridges the gap between representation learning and KD. By aligning feature embeddings in a contrastive manner, CRD emphasizes preserving semantic relationships in the learned representations, leading to improved generalization in student models, especially on image classification tasks.

One of the challenges in KD arises when there is a significant disparity in size and complexity between the teacher and student models. Mirzadeh et al. (2019) addressed this issue by introducing a teacher assistant model [4], which acts as an intermediary between the teacher and student. The assistant model is smaller than the teacher but larger than the student, enabling smoother and more stable knowledge transfer.

Key aspects of this approach:

- **Stability in Training**: The teacher assistant bridges the complexity gap, reducing the training difficulty for the student.
- **Limitations**: While effective, this method introduces additional complexity by requiring the training of an extra model, which can be computationally expensive and time-consuming.

To maintain simplicity, our project does not incorporate teacher assistants. Instead, we enhance traditional KD by combining logit matching and feature matching, allowing the student model to directly benefit from the teacher's knowledge. This streamlined approach balances simplicity with robust knowledge transfer, ensuring efficiency and practicality in implementation.

The methods discussed above have significantly shaped the design choices in our project. By combining the strengths of logit and feature matching, our approach addresses the need for efficient and effective knowledge transfer. Specifically:

- **Logit Matching**: Provides a strong foundation for the student model to learn nuanced class relationships, making it easier for the student to generalize.
- **Feature Matching**: Captures intermediate-level knowledge, improving the student's ability to replicate the teacher's hierarchical feature representations.
- **Efficiency**: By avoiding the complexity of additional models like teacher assistants, our approach ensures a streamlined and resource-efficient training pipeline.

This structured approach enables our project to achieve high accuracy in the student model while significantly reducing computational requirements. The project's focus on CIFAR-10 serves as an initial benchmark for evaluating the effectiveness of these techniques, with future plans to scale the approach to larger datasets like ImageNet. By building on established KD methods and tailoring them to meet the project's goals, we aim to advance the deployment of deep learning models on resource-constrained devices, making AI more accessible and practical in real-world scenarios.

## 3. Method

### 3.1. Teacher Model Selection and Configuration

The teacher model chosen for this project is a pre-trained ResNet-18, known for its robustness and high accuracy in classification tasks. ResNet-18's use of residual connections enables effective training and feature extraction, making it ideal for guiding a student model. The output layer is modified to match the datasets: 10 classes for CIFAR-10 and 1,000 classes for ImageNet.

To ensure consistency during knowledge distillation, the teacher model operates in evaluation mode, freezing its parameters to provide stable logits and feature representations. This prevents weight updates and ensures the teacher's knowledge remains intact throughout the process.

ResNet-18's pre-trained status on ImageNet allows it to offer high-quality feature representations and nuanced class relationships, which are crucial for guiding the student model. Its balance of accuracy and efficiency makes it an excellent choice for effective knowledge transfer.

### 3.2. Student Model Architecture

The student model is a custom-designed convolutional neural network (CNN), optimized for lightweight efficiency and resource-constrained environments. With fewer layers and parameters compared to the teacher model, the student model is tailored to deliver faster inference times and lower memory usage while retaining sufficient capacity to learn complex patterns.

The output layer of the student model is configured to match the dataset's class structure: 10 classes for CIFAR-10 and 1,000 classes for ImageNet. This ensures compatibility with the teacher's predictions and the target classification tasks.

The student model aims to achieve a compression ratio of approximately 1:5 relative to the teacher model. This balance between size reduction and accuracy preservation is critical for ensuring the student model remains practical for real-world deployment without compromising performance.

### 3.3. Knowledge Distillation Techniques

#### 3.3.1 Logit Matching

Logit matching is a foundational technique in knowledge distillation, introduced by Hinton et al., which involves training the student model to replicate the softened output logits of a larger, pre-trained teacher model. Unlike traditional training that uses hard class labels, logit matching allows the student to learn from the entire probability distribution generated by the teacher. This enables the student to capture richer relationships between classes, as the softened logits reveal not just the most probable class but also the relative likelihoods of other classes. This nuanced supervision helps the student model generalize better, especially when there are similarities between classes.

**Temperature Scaling**: A key component of logit matching is the use of a temperature parameter (e.g., T=3). By

applying this temperature to the logits before applying the softmax function, the probability distribution becomes "smoothed." This softening effect reduces the dominance of the highest-probability class and makes smaller probabilities more informative, providing the student model with subtle guidance on inter-class relationships. For example, if the teacher model predicts "cat" with high confidence, it may also assign a smaller probability to "tiger," indicating their similarity. The student learns from these relationships, which helps improve its ability to make accurate predictions.

**Loss Function**: To measure the difference between the teacher's softened logits and the student's outputs, we use Kullback-Leibler (KL) divergence, a standard metric for comparing probability distributions. The loss function is scaled by ($T^2$) to account for the effect of the temperature parameter on the gradients, ensuring that the optimization process remains stable and effective. This scaling ensures that the softening effect does not interfere with the gradient flow during training.

**Mathematics**: The distillation loss is calculated as:

$$L_{KD} = T^2 \times \text{KL}\left(\text{softmax}(z_{\text{teacher}}/T), \text{softmax}(z_{\text{student}}/T)\right)$$

where z represents the logits of each model, and T is the temperature. This method allows the student to match the teacher's class distribution.

**Key Benefits**: Logit matching enables the student model to learn from the teacher's entire probability distribution, providing more detailed guidance about class relationships than hard labels can offer. This helps the student model develop a more accurate understanding of the decision boundaries, improving its ability to generalize to new, unseen data. In this project, logit matching serves as the baseline method for transferring knowledge from the teacher to the student model, allowing the student to achieve high accuracy while being computationally efficient.

### 3.3.2 Feature Matching

Feature matching is an additional technique that complements logit matching by aligning intermediate layer activations between the teacher and student models. While logit matching focuses on matching the final predictions of the teacher and student, feature matching allows the student to learn from the internal representations of the teacher model at various layers. This method helps the student capture not only high-level abstract features (such as object shapes or textures) but also more detailed, low-level features (such as edges or patterns). By mimicking these intermediate activations, the student model learns a richer, more comprehensive representation of the data, improving its overall performance and ability to generalize.

**Layer Selection**: In ResNet-18, specific intermediate layers, such as selected convolutional blocks, are chosen to be matched between the teacher and student models. These layers are strategically selected to represent a range of feature types:

Shallow layers: These capture low-level features like edges, corners, and simple textures, which are essential for the initial stages of pattern recognition.

Deep layers: These represent higher-level features such as object parts, complex structures, and semantic information, which provide a more abstract understanding of the input data.

By aligning both shallow and deep layers, the student learns not only the high-level decision boundaries but also the fine-grained features that the teacher model has learned, leading to a more holistic understanding of the data.

**Loss Function**: To minimize the difference between the teacher and student feature maps at corresponding layers, Mean Squared Error (MSE) loss is used. The MSE loss measures the squared difference between the activations of the teacher and student models, encouraging the student to replicate the teacher's feature representations as closely as possible. This feature matching loss is then combined with the logit-matching loss to form the overall training objective. A weighting factor $\alpha$ is applied to control the relative importance of the feature matching loss in the overall objective, allowing fine-tuning of the balance between feature and logit learning.

**Mathematics**: The feature matching loss is calculated as:

$$L_{\text{feature}} = \alpha \times \sum_l \text{MSE}(\text{feature}_{\text{teacher}}^l, \text{feature}_{\text{student}}^l)$$

where $l$ represents the chosen layers (both shallow and deep) from the teacher and student models. $\alpha$ is the weighting factor that balances the importance of feature matching relative to logit matching.

**Key Benefits**: Feature matching enables the student model to learn not only the output predictions of the teacher but also the internal representations that lead to those predictions. This enhances the student's ability to generalize from complex patterns and improves its performance on unseen data. By combining feature matching with logit matching, the student model benefits from both high-level guidance and detailed feature alignment, ensuring a more effective transfer of knowledge. This technique is particularly useful in tasks like image classification, where learning both coarse and fine-grained features is essential for accurate predictions.

### 3.4. Combined Loss Function

The combined loss function is designed to effectively guide the training of the student model by incorporating three key

components: Hard Loss, Logit Matching Loss, and Feature Matching Loss. This multi-faceted approach allows the student model to learn from both the ground truth labels and the detailed knowledge provided by the teacher model, ensuring comprehensive knowledge transfer throughout the training process.

Hard Loss: This is the traditional cross-entropy loss that measures the difference between the student model's output and the true labels. It ensures that the student learns to make correct predictions based on the actual data labels. Hard loss serves as the foundational loss, encouraging the student to align with the true class labels, ensuring accuracy on the classification task.

Logit Matching Loss: This component represents the Kullback-Leibler (KL) divergence between the softened logits of the teacher and the student model. By softening the teacher's logits using a temperature parameter, the model is able to capture richer inter-class relationships, allowing the student to learn more nuanced class similarities beyond the hard classification decision. This logit matching helps the student mimic the teacher's class distribution, ensuring it learns from the teacher's predictions.

Feature Matching Loss: The feature matching loss measures the mean squared error (MSE) between the intermediate activations of corresponding layers in the teacher and student models. By aligning the feature representations from key layers of the teacher and student, this loss ensures that the student learns not only the final predictions but also the internal patterns and hierarchies that the teacher has learned. This loss is scaled by a factor , which controls the relative importance of feature matching in the overall loss.

**Final Loss Formula**: The total loss function for the student model is a weighted combination of these three components, formulated as:

$$L_{\text{total}} = (1 - \alpha) \times L_{\text{hard}} + \alpha \times (L_{\text{KD}} + L_{\text{feature}})$$

where $\alpha$ balances the contribution of the hard loss and the distillation losses. This combined loss function ensures that the student learns from both the ground truth labels and the nuanced information from the teacher model.

**Purpose of Combined Loss** The combined loss function ensures that the student model learns from both the true labels (via hard loss) and the valuable insights transferred from the teacher model (via logit matching and feature matching). By balancing these components, the student model is encouraged to:

- Make accurate predictions (hard loss)
- Mimic the teacher's decision-making process (logit matching)
- Replicate the teacher's internal feature representations (feature matching)

This balanced approach enhances the effectiveness of the knowledge distillation process, leading to a student model that is both accurate and efficient, capable of operating in resource-constrained environments while retaining much of the teacher model's high performance.

### 3.5. Training Procedure

The training procedure begins with the CIFAR-10 dataset, chosen for its suitability in rapid iteration and initial model validation. Using this smaller dataset allows quick evaluation of the training dynamics and model performance, providing insights into the student model's capacity to replicate the teacher model effectively.

After achieving satisfactory results on CIFAR-10, the model's scalability and generalizability are evaluated using the ImageNet dataset, a much larger and more complex benchmark. This progression ensures the robustness of the knowledge distillation process across varying dataset complexities.

The optimization process employs the Adam optimizer, chosen for its computational efficiency and adaptive learning rate adjustment capabilities. An initial learning rate is set, with adjustments made during training based on preliminary results to ensure convergence and performance improvement. The batch size of 64 is selected to balance between memory efficiency and stable training dynamics, especially for models evaluated on hardware with limited resources.

The training on CIFAR-10 spans 10 to 20 epochs, sufficient to assess the performance of the distillation process and refine the student model. Upon transitioning to ImageNet, the training procedure is extended and further optimized, with the number of epochs and hyperparameters fine-tuned based on the outcomes of the CIFAR-10 experiments.

To assess the effectiveness of knowledge transfer, we employ metrics such as Top-1 and Top-5 accuracy, inference time, and memory usage, measured consistently across both the teacher and student models. This evaluation allows a comprehensive comparison of the models in terms of accuracy and efficiency, reinforcing the practical applicability of the student model in resource-constrained scenarios.

This training pipeline ensures that the distilled student model not only inherits the teacher model's performance characteristics but also demonstrates superior efficiency and scalability, meeting the goals of knowledge distillation.

### 3.6. Evaluation Metrics

The evaluation of the knowledge distillation process involves assessing the student model's performance and efficiency across the following metrics:

1. **Accuracy**: Accuracy is measured using Top-1 and Top-5 accuracy on both CIFAR-10 and ImageNet datasets.

Top-1 accuracy reflects the percentage of test samples where the highest probability prediction matches the true label, while Top-5 accuracy measures whether the correct label is within the top five predictions. These metrics gauge the student model's ability to closely replicate the teacher model's performance.

2. **Compression Ratio**: This metric quantifies the reduction in model size by comparing the total parameters of the teacher model to those of the student model. A higher compression ratio demonstrates the extent of storage and computational savings achieved, making the student model more suitable for resource-constrained environments.

3. **Inference Time**: The average time taken per forward pass is evaluated to demonstrate real-time efficiency. Lower inference times indicate improved speed and suitability for deployment in latency-sensitive applications like mobile AI and edge computing.

4. **Memory Usage**: Memory consumption during inference is compared between the teacher and student models. This metric is crucial for verifying the feasibility of deploying the student model on devices with limited computational resources.

The student model is expected to achieve accuracy close to that of the teacher while significantly reducing the model's size and inference time. The combined use of logit and feature matching techniques is anticipated to help the student capture complex patterns effectively, making it deployable on devices where traditional models are impractical.

## 4. Results

In this section, we present the evaluation results of the student model in comparison to the teacher model, using both the CIFAR-10 and ImageNet datasets. We analyze performance in terms of Top-1 and Top-5 accuracy, model compression ratio, inference time, and memory usage. The results highlight the impact of knowledge distillation techniques (logit matching and partial feature matching) on improving the efficiency and performance of the student model.

### 4.1. Model Performance on CIFAR-10

For CIFAR-10, we trained the teacher and student models for 10 epochs using logit matching and feature matching techniques. The teacher model, ResNet-18, achieved a Top-1 accuracy of 88.28% and a Top-5 accuracy of 99.60%. The student model, which underwent knowledge distillation using logit matching, achieved a Top-1 accuracy of 83.58% and a Top-5 accuracy of 98.20%. This demonstrates the effectiveness of knowledge distillation in transferring knowledge from the teacher to the student model, although there was a slight reduction in accuracy.

Incorporating feature matching into the knowledge distillation process further improved the student model's performance, boosting the Top-1 accuracy to 84.58% and the Top-5 accuracy to 99.35%. While the accuracy increase was modest, the incorporation of intermediate layer activations enabled the student model to capture more nuanced features from the teacher model.

The student model achieved significant reductions in both inference time and memory usage, making it much more efficient than the teacher model. Specifically, the student model's inference time was reduced by approximately 36.46% (75.30 ms compared to the teacher's 118.50 ms), and memory usage was reduced by a factor of 5.42x (0.12 MB compared to 0.65 MB for the teacher).

| Model | Top-1 Accuracy | Top-5 Accuracy | Inference Time (seconds) | Memory Usage (MB) | Compression Ratio |
|---|---|---|---|---|---|
| Teacher Model | 88.28% | 99.60% | 118.50 ms | 0.65 MB | 1.00x |
| Student Model (Logit Matching) | 83.58% | 98.20% | 75.30 ms | 0.12 MB | 5.42x |
| Student Model (Logit + Feature Matching) | 84.58% | 99.35% | 75.30 ms | 0.12 MB | 5.42x |

Table 1. Model performance comparison on the CIFAR-10 dataset. Results for the teacher model, student model with logit matching, and student model with logit and feature matching, including Top-1 and Top-5 accuracy, inference time, memory usage, and model compression ratio.

### 4.2. Knowledge Distillation Method Impact

Knowledge distillation using logit matching enabled the student model to capture the teacher's class probabilities effectively, leading to a 4.7% reduction in Top-1 accuracy compared to the teacher model. However, this technique resulted in significant efficiency improvements, with a 36.46% reduction in inference time and a 5.42x reduction in memory usage.

To further improve the performance of the student model, we incorporated feature matching between intermediate layers of the teacher and student models. This technique led to a marginal increase in accuracy of 1% in Top-1 accuracy and 1.15% in Top-5 accuracy. The inclusion of feature matching allowed the student model to better capture fine-grained features and improved the overall performance, while the increase in training time per epoch remained manageable.

The combination of logit matching and feature matching resulted in an improved Top-1 accuracy of 84.58%, bringing the student model closer to the teacher's performance. The Top-5 accuracy also improved slightly, reaching 99.35%. This combination of distillation methods proved to be the most effective in improving both the accuracy and efficiency of the student model.

### 4.3. Model Performance on ImageNet

The teacher and student models were also evaluated on the ImageNet dataset. The teacher model, ResNet-18, achieved a Top-1 accuracy of 72.88% and a Top-5 accuracy

of 91.25%. The student model with logit matching alone achieved a Top-1 accuracy of 68.75% and a Top-5 accuracy of 89.75%, showing a similar reduction in accuracy as seen with CIFAR-10, but with significant improvements in inference time and memory usage.

With the addition of feature matching, the student model's performance improved, with the Top-1 accuracy rising to 70.15% and the Top-5 accuracy to 90.85%. Again, the increase in accuracy was modest, but the overall improvements in efficiency were notable. The student model achieved a 38.8% reduction in inference time (82.10 ms compared to the teacher's 134.15 ms) and a 5.17x reduction in memory usage (26.60 MB compared to the teacher's 137.50 MB).

| Model | Top-1 Accuracy | Top-5 Accuracy | Inference Time (seconds) | Memory Usage (MB) | Compression Ratio |
|---|---|---|---|---|---|
| Teacher Model | 72.88% | 91.25% | 134.15 ms | 137.50 MB | 1.00x |
| Student Model (Logit Matching) | 68.75% | 89.75% | 82.10 ms | 26.60 MB | 5.17x |
| Student Model (Logit + Feature Matching) | 70.15% | 90.85% | 82.10 ms | 26.60 MB | 5.17x |

Table 2. Model performance comparison on the ImageNet dataset. Results for the teacher model, student model with logit matching, and student model with logit and feature matching, including Top-1 and Top-5 accuracy, inference time, memory usage, and model compression ratio.

### 4.4. Summary of Results

The student model achieved significant reductions in inference time and memory usage while maintaining a good level of accuracy compared to the teacher model. The incorporation of feature matching alongside logit matching improved the performance of the student model, although the improvement in accuracy was modest. The student model outperformed the teacher model in terms of inference time and memory usage, making it a more suitable choice for deployment in resource-constrained environments.

Overall, knowledge distillation techniques, particularly when combining logit matching and feature matching, proved effective in transferring knowledge from the teacher model to the student model while improving efficiency. The results from both the CIFAR-10 and ImageNet datasets demonstrate the potential of knowledge distillation for building lightweight models with competitive performance.

## 5. Conclusion

This project has successfully demonstrated the effectiveness of knowledge distillation for training compact, efficient models while retaining high accuracy. Through the use of logit matching and feature matching, we were able to transfer knowledge from a larger, pre-trained teacher model to a smaller student model, achieving a significant reduction in model size and computational demands. Our experiments with the CIFAR-10 dataset confirmed that the student model, while much smaller, could replicate much of the teacher model's performance, validating the applicability of knowledge distillation in reducing model complexity without sacrificing accuracy.

The combination of logit and feature matching provided complementary benefits, where logit matching helped the student model learn from the teacher's softened output probabilities, and feature matching enabled the student to capture both high-level abstractions and low-level features from the teacher. Together, these methods enhanced the performance of the student model, making it a viable solution for deployment in resource-constrained environments, such as smartphones and IoT devices.

Based on these findings, we have outlined several next steps, including hyperparameter tuning to refine the balance between logit and feature matching, scaling our approach to the ImageNet dataset to test the model's generalizability to larger and more complex datasets, and exploring additional layers for feature matching to further improve accuracy.

In conclusion, our approach successfully achieves the goal of creating lightweight models that perform well in real-world applications. The results of this project provide a promising pathway for deploying AI in resource-limited environments without compromising on performance, paving the way for more accessible and efficient AI applications across a variety of domains.

## References

[1] Lei Jimmy Ba and Rich Caruana. Do deep nets really need to be deep? 2013. arXiv:1312.6184. 2

[2] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. 2015. arXiv:1503.02531. 2

[3] Xiaoqi Jiao, Yichun Yin, Lifeng Shang, Xin Jiang, Xiao Chen, Linlin Li, Fang Wang, and Qun Liu. Tinybert: Distilling bert for natural language understanding. 2019. arXiv:1909.10351. 2

[4] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, Nir Levine, Akihiro Matsukawa, and Hassan Ghasemzadeh. Improved knowledge distillation via teacher assistant. 2019. arXiv:1902.03393. 3

[5] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. Fitnets: Hints for thin deep nets. 2014. arXiv:1412.6550. 2

[6] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive representation distillation. 2019. arXiv:1910.10699. 2

[7] Junho Yim, Donggyu Joo, Jihoon Bae, and Junmo Kim. A gift from knowledge distillation: Fast optimization, network minimization and transfer learning. 2017. In Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR). 2

[8] Sergey Zagoruyko and Nikos Komodakis. Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. 2016. arXiv:1612.03928. 2