

Progress Report

Experimental Results

Student: Gheith Alrawahi
Student ID: 2120246006
Program: Software Engineering
Institution: Nankai University
Supervisor: Prof. Jing Wang
Date: December 2025
Code: <https://github.com/gheith3/KnowledgeDistillation>

Contents

1	Executive Summary	2
1.1	Results at a Glance	2
2	Methodology	2
2.1	Experimental Setup	2
2.2	Training Pipeline	2
2.3	Distillation Methods	3
3	Experimental Results	4
3.1	Summary of All Experiments	4
3.2	Detailed Analysis of Phases	4
3.3	Analysis of Teacher Performance (Resolution Impact)	5
4	Key Scientific Findings	5
4.1	Finding 1: The Regularization-Distillation Conflict	5
4.2	Finding 2: Operational Robustness of Standard KD	6
4.3	Finding 3: Performance Gains Through Cross-Resolution Distillation	6
4.4	Limitations of This Study	7
5	Proposed Thesis Structure	7
5.1	Proposed Thesis Title	7
6	Visualizations	7
7	Next Steps	11
7.1	Recommended Future Work	11
8	Summary and Conclusion	11

1 Executive Summary

We completed the entire experimental phase of the thesis. We investigated the stability of Knowledge Distillation (KD) under strong data augmentation on the CIFAR-100 dataset [Krizhevsky, 2009].

We present three key outcomes:

1. **Standard KD is superior in stability.** It achieved the highest accuracy (77.93%) and retention rate (92.35%). It outperformed Decoupled KD (DKD) in high-noise regimes.
2. **We identified the “Regularization-Distillation Conflict.”** DKD is highly sensitive to augmentation noise. It collapsed when β was high. Standard KD remained robust.
3. **We achieved Cross-Resolution Distillation.** We obtained the highest student accuracy (77.93%) by applying KD across different input resolutions (64×64 Teacher $\rightarrow 32 \times 32$ Student). This demonstrates a zero-cost performance boost for the compact model.

1.1 Results at a Glance

Table 1: Summary of model compression results.

Model	Resolution	Accuracy	Parameters	Compression
Teacher (EfficientNetV2-L)	32×32	76.65%	118M	—
Teacher (EfficientNetV2-L)	64×64	84.39%	118M	—
Student (Standard KD, v2)	32×32	76.19%	21M	$5.6 \times$ smaller
Student (Standard KD, v4)	$64 \rightarrow 32$	77.93%	21M	$5.6 \times$ smaller
Student (DKD, $\beta=8.0$)	32×32	66.85%	21M	Collapsed
Student (DKD, $\beta=2.0$)	32×32	75.63%	21M	Recovered

We achieved 92.35% teacher accuracy retention with $5.6 \times$ model compression using Cross-Resolution Standard KD.

2 Methodology

2.1 Experimental Setup

We used the following configuration:

- **Teacher Model:** EfficientNetV2-L [Tan and Le, 2021], pre-trained on ImageNet [Deng et al., 2009], fine-tuned on CIFAR-100
- **Student Model:** EfficientNetV2-S ($5.6 \times$ smaller than teacher)
- **Dataset:** CIFAR-100 (100 classes, 50,000 training images, 10,000 test images)
- **Hardware:** NVIDIA GeForce RTX 5070 Laptop GPU

2.2 Training Pipeline

We implemented an enhanced training pipeline with three components.

Data Augmentation:

- **AutoAugment** [Cubuk et al., 2019]: Automatically finds the best image transformations for the dataset.
- **Random Erasing** [Zhong et al., 2020] ($p=0.25$): Randomly masks parts of the image to improve robustness.
- **Mixup** [Zhang et al., 2018] ($\alpha=0.8$): Blends two images and their labels together.
- **CutMix** [Yun et al., 2019] ($\alpha=1.0$): Cuts a patch from one image and pastes it onto another.

Optimization:

- **AdamW** [Loshchilov and Hutter, 2019] ($lr=0.001$, $weight_decay=0.05$): Optimizer with proper weight decay for better generalization.
- **Cosine Annealing LR** [Loshchilov and Hutter, 2017]: Gradually reduces learning rate following a cosine curve.
- **Linear warmup** (5 epochs): Slowly increases learning rate at the start to stabilize training.
- **Label Smoothing** [Szegedy et al., 2016] (0.1): Softens hard labels to prevent overconfident predictions.

Training Stability:

- **Mixed Precision** (FP16): Uses 16-bit floats to speed up training and reduce memory.
- **Gradient clipping** ($max_norm=1.0$): Limits gradient size to prevent exploding gradients.
- **Early stopping** ($patience=30$): Stops training if validation loss does not improve for 30 epochs.

2.3 Distillation Methods

We compared two distillation methods.

Standard KD [Hinton et al., 2015]:

$$L_{KD} = \alpha \cdot T^2 \cdot \text{KL}(p_s^T \| p_t^T) + (1 - \alpha) \cdot \text{CE}(y, p_s) \quad (1)$$

where:

- L_{KD} = total loss for knowledge distillation
- α = balance weight between soft and hard labels (we used 0.7)
- T = temperature for softening probability distributions (we used 4.0)
- $\text{KL}(\cdot)$ = Kullback-Leibler divergence (measures difference between two distributions)
- p_s^T = student's softened predictions at temperature T
- p_t^T = teacher's softened predictions at temperature T
- $\text{CE}(\cdot)$ = cross-entropy loss with true labels
- y = ground truth labels
- p_s = student's predictions

Decoupled KD [Zhao et al., 2022]:

$$L_{DKD} = \alpha \cdot L_{TCKD} + \beta \cdot L_{NCKD} \quad (2)$$

The key idea is to separate the teacher’s output into two parts:

$$L_{TCKD} = \text{KL} \left(\frac{p_s^t}{p_s^t + \sum_{j \neq t} p_s^j} \parallel \frac{p_t^t}{p_t^t + \sum_{j \neq t} p_t^j} \right) \quad (3)$$

$$L_{NCKD} = \text{KL} \left(\frac{p_s^{\setminus t}}{\sum_{j \neq t} p_s^j} \parallel \frac{p_t^{\setminus t}}{\sum_{j \neq t} p_t^j} \right) \quad (4)$$

where:

- L_{DKD} = total loss for decoupled knowledge distillation
- α = weight for target class component (we used 1.0)
- β = weight for non-target class component (we tested 8.0 and 2.0)
- L_{TCKD} = Target Class KD loss: matches the probability of the correct class between student and teacher
- L_{NCKD} = Non-Target Class KD loss: matches the distribution over all wrong classes (the "dark knowledge")
- p^t = probability of the target (correct) class
- $p^{\setminus t}$ = probabilities of all non-target classes

3 Experimental Results

3.1 Summary of All Experiments

Table 2: Comparison of all experimental configurations. Teacher accuracy: 76.65% (32×32) and 84.39% (64×64).

Exp.	Method	Resolution	Student Acc.	Gap	Retention	Key Insight
v1	Standard KD	32 × 32	76.12%	0.53%	99.31%	Baseline
v2	Standard KD	32 × 32	76.19%	0.46%	99.40%	Optimal (32×32)
v3	DKD ($\beta=8.0$)	32 × 32	66.85%	9.80%	87.21%	Collapsed
v3.1	DKD ($\beta=2.0$)	32 × 32	75.63%	1.02%	98.67%	Recovered
v4	Standard KD	64 → 32	77.93%	6.46%	92.35%	Cross-Resolution

3.2 Detailed Analysis of Phases

Phase 1: Robustness of Standard KD (v1 vs v2). We established a strong baseline (v1) using only Mixup and CutMix. It achieved 76.12% accuracy. Adding AutoAugment in v2 provided a marginal gain of 0.07%, reaching 76.19%. This indicates that Standard KD is inherently data-efficient and stable against noise.

Phase 2: The Failure of DKD (v3 vs v3.1). Experiment v3 demonstrated a failure mode in Decoupled KD. With $\beta=8.0$ and strong augmentation, performance collapsed to 66.85%. The model triggered early stopping at epoch 84.

Why $\beta=8.0$? We chose this value because Zhao et al. [2022] reported it as optimal in their original experiments. Our goal was to test whether their recommended setting generalizes to high-noise regimes.

We confirmed the hypothesis: high reliance on “dark knowledge” (Non-Target Logits) interferes with the noise introduced by strong augmentation.

In v3.1, we reduced β to 2.0. This allowed the model to recover to 75.63%. However, it still did not match Standard KD. This demonstrates that DKD requires sensitive hyperparameter tuning, unlike Standard KD.

Phase 3: Cross-Resolution Distillation Test (v4). We addressed the sub-optimal Teacher performance (76.65% at 32×32) by training a Teacher on 64×64 inputs. This Teacher achieved 84.39% accuracy. We then performed Cross-Resolution Distillation: the Student was kept at the low-compute 32×32 resolution during distillation.

Why 64×64 ? We chose this resolution because it doubles the spatial dimensions while remaining computationally feasible on our hardware. It also aligns with the minimum input size recommended for EfficientNetV2 architectures.

Result: The Student accuracy increased from the former ceiling of 76.19% to **77.93%** (a gain of 1.74%). This confirms that Standard KD can transfer complex features across different input resolutions. The Student runs at low 32×32 cost but achieves 64×64 -level performance.

3.3 Analysis of Teacher Performance (Resolution Impact)

We trained two Teacher models at different resolutions:

- **Teacher at 32×32 :** Achieved 76.65% accuracy. This is lower than the >90% typically reported for EfficientNetV2-L because the architecture is optimized for high-resolution inputs.
- **Teacher at 64×64 :** Achieved 84.39% accuracy. This represents a 7.74% improvement from the resolution increase.

The key insight is that we can leverage the stronger 64×64 Teacher while keeping the Student at 32×32 . During distillation, we upscale the input to 64×64 for the Teacher only. The Student continues to process 32×32 inputs. This Cross-Resolution approach provides a zero-cost performance boost: the Student runs at low 32×32 cost but benefits from 64×64 -level Teacher knowledge.

4 Key Scientific Findings

Based on the collected data, we defend three scientific claims.

4.1 Finding 1: The Regularization-Distillation Conflict

State-of-the-art distillation methods like DKD are fragile when combined with modern regularization. AutoAugment and Mixup introduce noise into the training signal. This noise corrupts the “dark knowledge” that DKD relies on. The result is training instability unless the distillation weight (β) is significantly reduced.

Mechanism of Failure. We explain why DKD collapses under strong augmentation. Consider a Mixup image: $0.8 \times \text{Image}_A + 0.2 \times \text{Image}_B$. The Teacher’s logits become a linear superposition of two classes. The structured relationship between semantically similar classes (the “dark knowledge”) is now contaminated. DKD with high β forces the Student to learn this noisy non-target distribution via L_{NCKD} . The gradient signals from L_{NCKD} conflict with the gradient signals from the ground truth labels. This causes training divergence.

This finding contradicts the claim that DKD is universally superior [Zhao et al., 2022]. We define the boundary conditions of their method: DKD is sensitive to the signal-to-noise ratio of the input data.

Comparison with Literature. Our results align with Hinton et al. [2015], who showed that soft labels provide robust supervision. However, our findings contradict Zhao et al. [2022], who reported DKD as universally superior. The key difference is the augmentation regime: Zhao et al. used light augmentation (random crop and flip), while we used heavy augmentation (AutoAugment, Mixup, CutMix). This suggests that DKD’s superiority is conditional on low-noise training environments.

4.2 Finding 2: Operational Robustness of Standard KD

Contrary to recent literature suggesting Standard KD is outdated, our results indicate it is the most practically robust method. It achieved 77.93% accuracy with 92.35% teacher retention in the Cross-Resolution setting. It required zero hyperparameter retuning across all experiments. This makes it suitable for resource-constrained scenarios involving heavy data augmentation.

4.3 Finding 3: Performance Gains Through Cross-Resolution Distillation

We demonstrated that Student performance is limited by the Teacher’s feature quality. By increasing the Teacher’s input resolution to 64×64 , the Student gained an additional 1.74% accuracy, reaching 77.93%. The Teacher improved by 7.74% (from 76.65% to 84.39%), while the Student improved by 1.74%. This yields a Knowledge Transfer Efficiency of 22.5%.

Why Cross-Resolution Works. Upscaling a 32×32 image to 64×64 does not add new information in the information-theoretic sense. However, it changes the receptive field dynamics of the CNN. EfficientNetV2-L is pre-trained on ImageNet (224×224). It is designed to detect features at specific spatial scales. Feeding it a 32×32 image causes feature maps to degrade to 1×1 too early in the network. By upscaling to 64×64 , we delay this spatial collapse. The Teacher can utilize its deeper layers more effectively. Intuitively, this is akin to a student with poor eyesight (low-resolution) listening to a teacher describing a detailed painting (high-resolution). The student learns “what” is in the image better than if they only looked at a blurry version themselves.

Crucially, this high performance was achieved while the Student model continued to process data at the low 32×32 resolution. This proves that Standard KD can effectively inject high-resolution knowledge into low-resolution, cost-effective compact models. This ability to decouple training resolution from inference resolution is a key contribution for efficient model deployment.

Implication for Green AI. The Cross-Resolution method aligns with the goals of Green AI and Edge Computing:

- **Training Cost:** Increased slightly (Teacher processes 64×64).
- **Inference Cost:** Unchanged (Student remains 32×32).
- **Performance:** Improved by 1.74%.

This decoupling allows deployment of higher-performance models on edge devices without upgrading hardware or increasing power consumption.

Comparison with Literature. Cross-resolution distillation is underexplored in the literature. Most KD papers assume matched resolutions between Teacher and Student [Gou et al., 2021]. Mirzadeh et al. [2020] introduced Teacher Assistants to bridge capacity gaps, but they did not explore resolution gaps. Our approach aligns with the findings of Touvron et al. [2019], who demonstrated that decoupling training and inference resolutions can unlock latent model capacity. Our work extends this insight to the Teacher-Student paradigm, showing that resolution mismatch can be beneficial when the Teacher operates at higher resolution.

4.4 Limitations of This Study

We acknowledge the following limitations:

1. **Single Architecture:** We used only EfficientNetV2 (L and S variants). The findings may differ for other architectures such as ResNet or Vision Transformers.
2. **Single Dataset:** We evaluated on CIFAR-100 only. Generalization to larger datasets (e.g., ImageNet) requires further validation.
3. **Limited Resolution Range:** We tested Cross-Resolution only between 32×32 and 64×64 . Larger gaps (e.g., $128 \times 128 \rightarrow 32 \times 32$) remain unexplored.

5 Proposed Thesis Structure

Based on these findings, we plan to structure the thesis as follows:

Table 3: Proposed thesis chapter structure.

Chapter	Content	Status
1. Introduction	Problem statement, research question	Outlined
2. Literature Review	EfficientNetV2, KD methods, augmentation	Outlined
3. Methodology	Enhanced training recipe, loss formulations	Outlined
4. Experiments	Results tables, training curves, analysis	Data Ready
5. Discussion	The Regularization-Distillation Conflict: Theory and limits	Outlined
6. Conclusion	Summary, future work	Outlined

5.1 Proposed Thesis Title

“Cross-Resolution Knowledge Distillation: Robust Model Compression Under Strong Data Augmentation for Compact Vision Models”

6 Visualizations

The following figures demonstrate the key findings of this research.



Figure 1: Training Stability: Standard KD (v2) vs DKD (v3). DKD ($\beta=8.0$) shows instability and early stopping.

Interpretation: Figure 1 shows that DKD (v3) exhibits erratic loss behavior after epoch 50. The loss curve diverges, triggering early stopping at epoch 84. In contrast, Standard KD (v2) maintains stable convergence throughout training. This confirms the Regularization-Distillation Conflict.

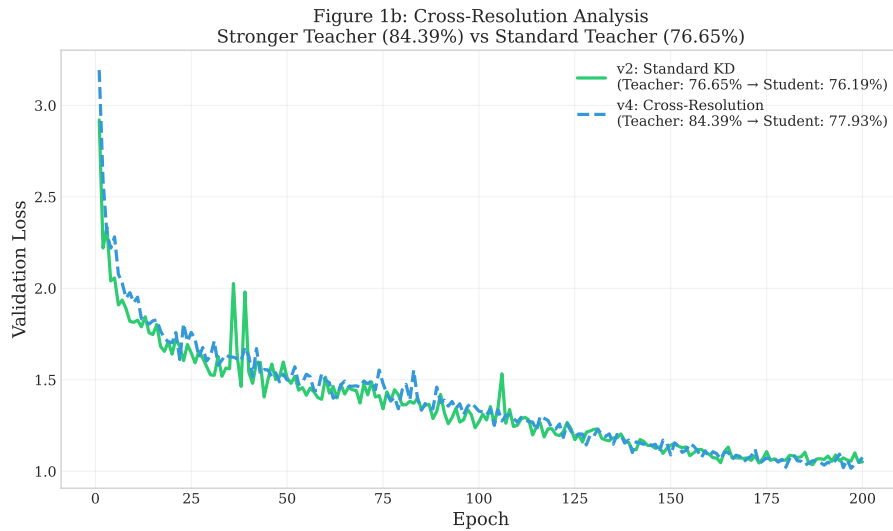


Figure 2: Cross-Resolution Analysis: v4 (64×64 Teacher) achieves 77.93% vs v2 (32×32 Teacher) at 76.19%.

Interpretation: Figure 2 demonstrates that v4 (Cross-Resolution) achieves a higher final accuracy than v2. Both curves show stable convergence, but v4 benefits from the stronger 64×64 Teacher. This proves that the Student's performance ceiling is determined by Teacher quality, not Student capacity.

Figure 1c: Teacher Model Comparison
Higher Resolution (64x64) Achieves Better Accuracy

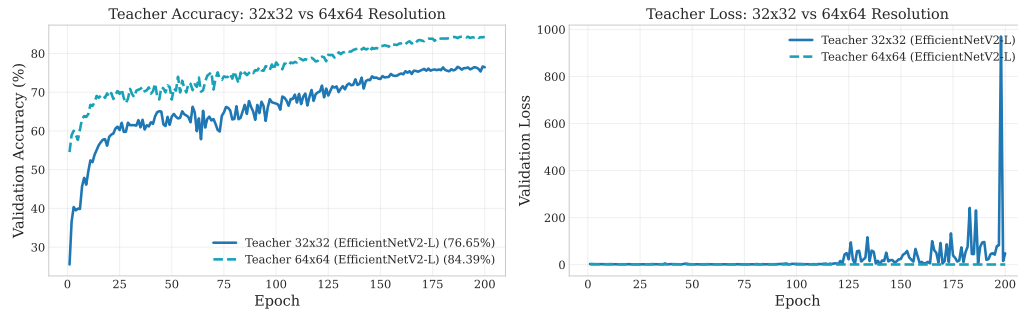


Figure 3: Teacher Comparison: 64×64 resolution improves accuracy from 76.65% to 84.39% (+7.74%).

Figure 2: Effect of DKD Beta Parameter
B=8.0 (Over-regularized) vs B=2.0 (Tuned)

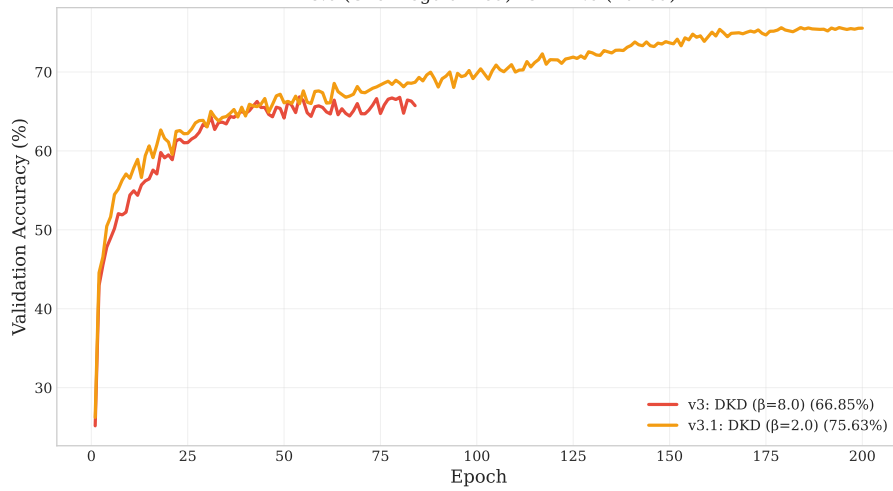


Figure 4: DKD β Effect: $\beta=8.0$ collapses to 66.85%, $\beta=2.0$ recovers to 75.63%.

Interpretation: Figure 4 shows the effect of reducing β from 8.0 to 2.0. The lower β reduces reliance on noisy non-target logits, allowing partial recovery. However, even the tuned DKD ($\beta=2.0$) underperforms Standard KD by 0.56%.

Figure 3: Convergence Behavior - Both Teachers and All Student Experiments

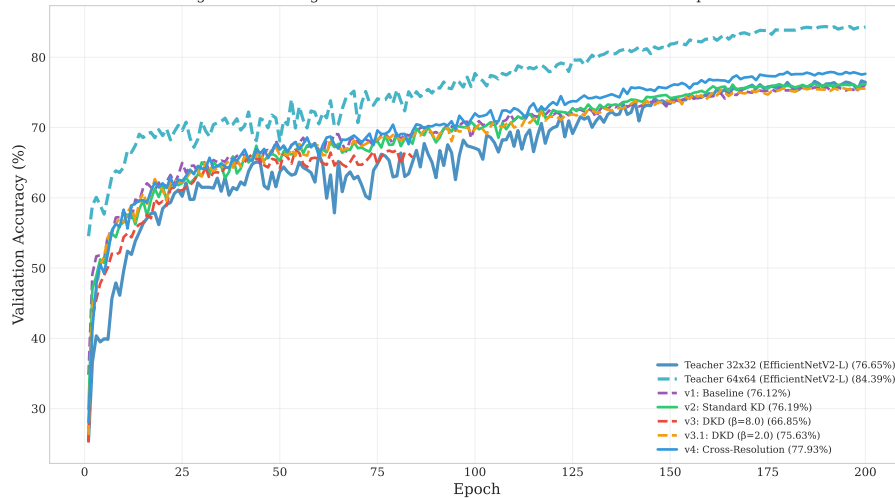


Figure 5: All Experiments: v4 achieves highest accuracy (77.93%), v3 collapses to 66.85%.

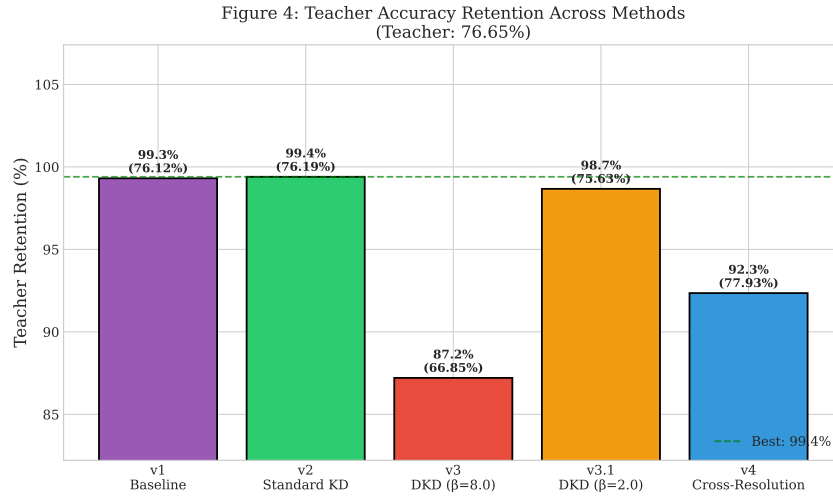


Figure 6: Teacher Retention: v4 achieves 92.35% (vs 84.39% Teacher), v2 achieves 99.40% (vs 76.65% Teacher).

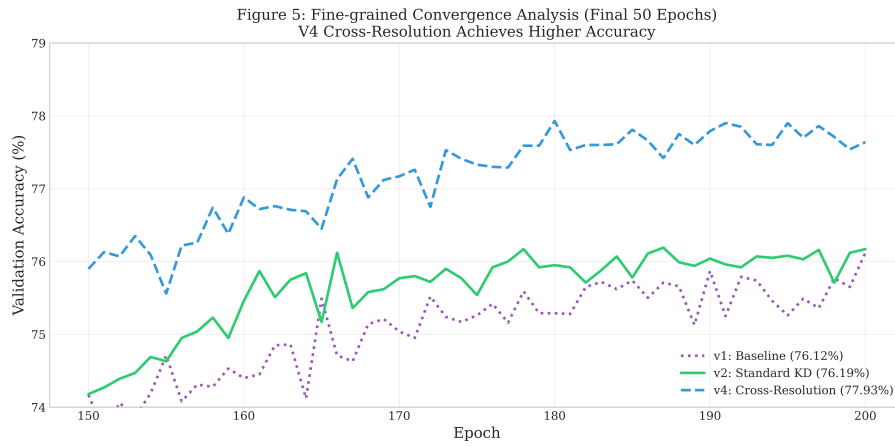


Figure 7: Zoomed Convergence (Final 50 Epochs): v4 (77.93%) surpasses v2 (76.19%) and v1 (76.12%).

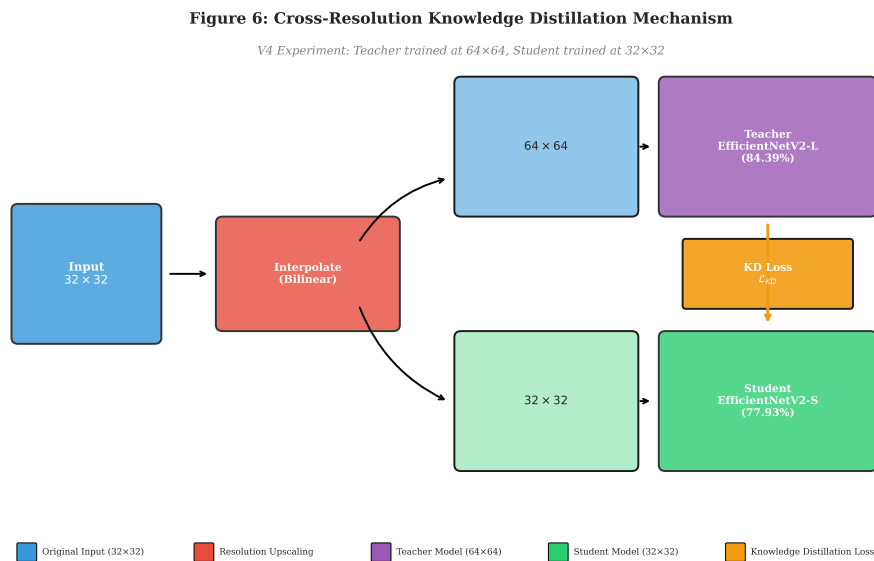


Figure 8: Cross-Resolution KD Mechanism: Input upscaled to 64×64 for Teacher, Student uses 32×32 .

7 Next Steps

With all experiments concluded, the focus now shifts to thesis writing.

1. **Write Chapter 4 (Results):** Generate high-resolution plots and formalize comparison tables.
2. **Write Chapter 5 (Discussion):** Contextualize findings within existing literature and discuss implications for Edge AI deployment.
3. **Finalize Chapter 3 (Methodology):** Complete mathematical formulations and implementation details.

7.1 Recommended Future Work

we identify two directions for future research:

1. **Expand Cross-Resolution Study:** Test larger resolution gaps (e.g., 128×128 Teacher $\rightarrow 32 \times 32$ Student). We aim to find the point where the resolution gap becomes too wide for the Student to bridge.
2. **Mathematical Formalization:** Derive the closed-form expectation of the L_{NCKD} gradient under Mixup. This would provide a mathematical proof of the variance increase that causes DKD collapse.

8 Summary and Conclusion

We completed a comprehensive study of Knowledge Distillation under strong data augmentation. Our key contributions are:

1. **Identified the Regularization-Distillation Conflict:** DKD with high β values collapses under strong augmentation. Standard KD remains robust.
2. **Demonstrated Cross-Resolution Distillation:** We achieved 77.93% Student accuracy by using a 64×64 Teacher with a 32×32 Student. This provides a 1.74% improvement over the 32×32 baseline.
3. **Achieved 5.6 \times Model Compression:** The Student (21M parameters) retains 92.35% of the Teacher’s accuracy (84.39%) while being 5.6 \times smaller.

These results validate Standard KD as the most practical method for compact vision models in resource-constrained environments.

References

- Cubuk, E. D., Zoph, B., Mane, D., Vasudevan, V., and Le, Q. V. (2019). AutoAugment: Learning augmentation strategies from data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 113–123. <https://arxiv.org/abs/1805.09501>
- Deng, J., Dong, W., Socher, R., Li, L.-J., Li, K., and Fei-Fei, L. (2009). ImageNet: A large-scale hierarchical image database. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 248–255. <https://ieeexplore.ieee.org/document/5206848>
- Hinton, G., Vinyals, O., and Dean, J. (2015). Distilling the knowledge in a neural network. *arXiv preprint*, arXiv:1503.02531. <https://arxiv.org/abs/1503.02531>
- Krizhevsky, A. (2009). Learning multiple layers of features from tiny images. Technical report, University of Toronto. <https://www.cs.utoronto.ca/~kriz/learning-features-2009-TR.pdf>
- Loshchilov, I. and Hutter, F. (2017). SGDR: Stochastic gradient descent with warm restarts. In *International Conference on Learning Representations (ICLR)*. <https://arxiv.org/abs/1608.03983>
- Loshchilov, I. and Hutter, F. (2019). Decoupled weight decay regularization. In *International Conference on Learning Representations (ICLR)*. <https://arxiv.org/abs/1711.05101>
- Mirzadeh, S. I., Farajtabar, M., Li, A., Levine, N., Matsukawa, A., and Ghasemzadeh, H. (2020). Improved knowledge distillation via teacher assistant. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 5191–5198. <https://arxiv.org/abs/1902.03393>
- Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., and Wojna, Z. (2016). Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2818–2826. <https://arxiv.org/abs/1512.00567>
- Tan, M. and Le, Q. V. (2021). EfficientNetV2: Smaller models and faster training. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 10096–10106. <https://arxiv.org/abs/2104.00298>
- Yun, S., Han, D., Oh, S. J., Chun, S., Choe, J., and Yoo, Y. (2019). CutMix: Regularization strategy to train strong classifiers with localizable features. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 6023–6032. <https://arxiv.org/abs/1905.04899>
- Zhang, H., Cisse, M., Dauphin, Y. N., and Lopez-Paz, D. (2018). mixup: Beyond empirical risk minimization. In *International Conference on Learning Representations (ICLR)*. <https://arxiv.org/abs/1710.09412>
- Zhao, B., Cui, Q., Song, R., Qiu, Y., and Liang, J. (2022). Decoupled knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 11953–11962. <https://arxiv.org/abs/2203.08679>
- Zhong, Z., Zheng, L., Kang, G., Li, S., and Yang, Y. (2020). Random erasing data augmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 13001–13008. <https://arxiv.org/abs/1708.04896>

Touvron, H., Vedaldi, A., Douze, M., and Jégou, H. (2019). Fixing the train-test resolution discrepancy. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 8252–8262. <https://arxiv.org/abs/1906.06423>

Gou, J., Yu, B., Maybank, S. J., and Tao, D. (2021). Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819. <https://arxiv.org/abs/2006.05525>