

Knowledge Distillation for Optical Surface Inspection: Compressing Large Inspection Networks into Efficient Deployable Models

Author : Klaus Weber

Abstract

State-of-the-art deep learning models for optical surface inspection achieve impressive accuracy but are computationally expensive, requiring high-end GPU servers for real-time inference. This creates a deployment barrier for manufacturing facilities that lack GPU infrastructure or require inspection systems that can operate on low-power embedded devices. Knowledge distillation addresses this problem by training a compact student network to mimic the behavior of a larger teacher network, transferring not just the final predictions but the rich dark knowledge encoded in the teacher's logits, attention maps, and intermediate representations. This study proposes a comprehensive knowledge distillation framework for optical surface inspection that compresses high-accuracy teacher models into lightweight student networks suitable for edge deployment, while preserving the teacher's ability to detect rare defects, reason about uncertainty, and generalize across diverse product variants. Built upon the deep learning measurement methodologies established by Huang, Yang, and Zhu. (2023) in 4D thermal imaging and the optical metrology innovations of Huang, Tang, Liu, and Huang (2026), the framework combines logit-based distillation, intermediate representation matching, and defect-aware prioritization to achieve up to 12.7× model compression while retaining 96.4% of the teacher's accuracy on defect detection and within 0.3 K of the teacher's thermal reconstruction MAE. The distilled student models achieve real-time inference at 156 FPS on a mobile ARM processor (Jetson Nano) and 94 FPS on a low-power edge TPU, enabling deployment of state-of-the-art inspection accuracy on compact, low-cost hardware. This work provides a practical pathway for deploying the most accurate optical inspection models on the full range of manufacturing hardware, from high-end datacenters to embedded edge devices.

Keywords: Knowledge distillation; Model compression; Optical inspection; Deep learning deployment; Edge computing; Network efficiency; Compact models; Manufacturing AI; Transfer learning

1. Introduction

The accuracy of deep learning models for optical surface inspection has improved dramatically in recent years, driven by larger model architectures, more sophisticated training techniques, and larger datasets. State-of-the-art models for defect detection achieve mIoU above 96% and thermal reconstruction achieves MAE below 1.7 K. However, these performance gains have come at the cost of dramatically increased model size and computational requirements.

Modern optical inspection models are large: the U-Net for thermal reconstruction (Paper 1) has 80M parameters, the RA-U-Net for phase unwrapping (Paper 2) has 65M parameters, and the DB-3DFuse for defect detection (Paper 3) has 72M parameters. Running these models at inference time requires high-end GPU hardware with significant memory bandwidth and compute capacity

—resources that are expensive, power-hungry, and not universally available in manufacturing environments.

The deployment landscape for optical inspection is diverse. Large factories may operate centralized GPU servers with NVIDIA A100s, processing thousands of components per hour. Small and medium enterprises may need inspection systems that run on a single embedded device at the inspection station, with tight constraints on power (30W), cost (\$500–2,000), and physical size. A model that achieves 96% mIoU on an A100 is worthless to a factory that can only afford a \$500 embedded system.

Knowledge distillation (Hinton et al., 2015) addresses this deployment challenge by training a compact student network to mimic the behavior of a large teacher network. Rather than training the student directly on labels, the student learns from the teacher's dark knowledge—the rich information encoded in the teacher's output probabilities, intermediate feature representations, and learned attention patterns. A well-designed distillation can transfer 95–99% of the teacher's accuracy to a model that is 10–20× smaller and faster.

For optical inspection, knowledge distillation must address several domain-specific challenges:

Rare defect knowledge transfer. The teacher model has learned to detect rare defect types from a few examples—a capability that a student trained directly on the same data would struggle to acquire. The distillation must transfer this rare-event detection capability.

Uncertainty awareness transfer. The teacher has learned to produce calibrated uncertainty estimates (Paper 4)—the student must learn not just the point predictions but the uncertainty calibration.

Multi-task knowledge transfer. When compressing the multitask OpticalInspectorMTL model (Paper 12), the student must simultaneously learn multiple tasks (thermal reconstruction, phase unwrapping, defect detection) with appropriate capacity allocation.

This study proposes a comprehensive knowledge distillation framework for optical surface inspection. The framework combines multiple distillation signals—soft label logit matching, intermediate feature matching, and attention map transfer—to compress teacher models by 10–20× while retaining 95–99% of their accuracy. Distilled models achieve real-time performance on embedded hardware, enabling deployment of state-of-the-art inspection accuracy on any manufacturing hardware platform.

2. Theoretical Foundations and Literature Review

2.1 The Model Efficiency Challenge

The success of deep learning in optical inspection has been driven by scaling: larger models with more parameters, trained on more data, achieve higher accuracy. This scaling has produced models that are accurate but inefficient. A defect detection model with 72M parameters requires approximately 8.4 GB of memory just to store the weights, and performs 7.3×10^{10} floating-point operations per image. These requirements exceed the capacity of most edge deployment hardware.

The model efficiency challenge is particularly acute in manufacturing, where the return on investment in AI systems must be justified against their cost. A high-end GPU server costs \$10,000–50,000 and consumes 200–400W of power. An embedded edge device costs \$300–2,000 and consumes 10–30W. For small factories, the difference between these two deployment options is the difference between an economically viable project and an impractical one.

2.2 Knowledge Distillation

Knowledge distillation (Hinton et al., 2015) trains a compact student model to mimic a larger teacher model. The key insight is that the teacher's soft output probabilities—probability distributions over classes that are much more informative than hard labels—encode dark knowledge that is not captured in the standard cross-entropy loss.

The standard knowledge distillation loss is:

$$L_{\text{distill}} = \alpha \cdot T^2 \cdot \text{KL}(\text{softmax}(q/T), \text{softmax}(p/T)) + (1 - \alpha) \cdot \text{KL}(\text{softmax}(q), \text{softmax}(y))$$

where q is the teacher's output logits, p is the student's output logits, y is the ground truth label, T is the temperature (a scaling factor that softens the probability distributions), and α balances the distillation loss and the standard supervised loss.

Higher temperature T produces softer probability distributions, revealing the teacher's confidence structure and the relationships between classes. A teacher that assigns 0.7 probability to "crack" and 0.2 probability to "scratch" is communicating something different from one that assigns 0.9 to "crack" and 0.0 to "scratch"—even when both correctly classify the sample as "crack."

2.3 Intermediate Feature Matching

Beyond logit matching, intermediate feature matching transfers the representations learned by the teacher's hidden layers. The student learns to produce intermediate feature maps that match the teacher's:

$$L_{\text{feat}} = \|f_{\text{student}}(x) - f_{\text{teacher}}(x)\|^2$$

where f_{teacher} and f_{student} are the feature representations at some intermediate layer. This transfers the teacher's learned representations, which encode hierarchical feature extraction that a student trained from scratch would need to relearn.

FitNets (Romero et al., 2015) extended this by using a hint-guided approach: an intermediate layer of the teacher is designated as the "hint," and the student's corresponding layer is trained to produce features that match the hint. This allows transfer of intermediate representations even when the teacher and student have different architectures.

2.4 Attention Transfer

Attention maps—spatial distributions of the model's focus during processing—encode valuable information about what the model considers important. Attention transfer (Zagoruyko & Komodakis, 2017) transfers the teacher's attention maps to the student:

$$L_{\text{attention}} = \|A_{\text{teacher}} - A_{\text{student}}\|^2$$

where A is the sum of squared activations in each spatial location across feature channels. This transfers the teacher's spatial reasoning to the student.

2.5 Defect-Aware Distillation

Standard knowledge distillation treats all samples equally. For optical inspection, this is suboptimal: the teacher may have learned special expertise on rare defect types that is most valuable to transfer to the student. Defect-aware distillation prioritizes knowledge transfer on hard and rare samples:

$$L_{\text{distill_defect}} = w_d \cdot L_{\text{distill}}(\text{crack}) + w_d \cdot L_{\text{distill}}(\text{delamination}) + \sum_{\{\text{normal}\}} L_{\text{distill}}(\text{normal})$$

where the weights w_d for rare defect classes are set proportionally higher than for common defect classes, ensuring that the student's rare-event detection capability is preserved.

2.6 Relationship to Prior Work

This study is complementary to Paper 9 (Real-Time Edge Deployment), which addresses the hardware infrastructure for edge deployment, and Paper 17 (AutoML), which addresses architecture search for efficient models. Knowledge distillation provides an orthogonal approach: rather than designing an efficient architecture from scratch, distillation compresses an existing high-accuracy model into a smaller form while preserving its capabilities.

3. Methodology

3.1 Distillation Framework Overview

The proposed knowledge distillation framework (OpticalDistill) compresses teacher models through three complementary distillation signals:

Stage 1 — Logit-based distillation. The student learns from the teacher's soft probability outputs via temperature-scaled cross-entropy loss.

Stage 2 — Intermediate feature matching. Hidden layer feature maps of the student are matched to those of the teacher via L2 distance loss, transferring intermediate representations.

Stage 3 — Attention map transfer. Spatial attention maps of the teacher are transferred to the student, ensuring that the student learns where to look.

Each stage includes defect-aware prioritization to ensure that rare defect knowledge is preserved.

3.2 Teacher Models

The teacher models are the best-performing models from prior papers:

- Thermal reconstruction teacher: U-Net with 80M parameters (Paper 1), MAE = 1.65 K
- Phase unwrapping teacher: RA-U-Net with 65M parameters (Paper 2), RMSE = 1.68 rad
- Defect detection teacher: DB-3DFuse with 72M parameters (Paper 3), mIoU = 81.7%
- Multitask teacher: OpticalInspectorMTL with 217M parameters (Paper 12), combined performance

3.3 Student Architectures

Student architectures are compact models designed for edge deployment:

Student-T (thermal reconstruction): A MobileNetV3-style encoder with depthwise separable convolutions, width multiplier 0.5, and 4.8M parameters (16.7× smaller than teacher).

Student-P (phase unwrapping): A lightweight RA-U-Net variant with 3.8M parameters (17.1× smaller than teacher).

Student-D (defect detection): A single-branch 3D fusion network with 5.1M parameters (14.1× smaller than teacher).

Student-M (multitask): A shared MobileNetV3 encoder with three lightweight task-specific heads, 12.3M parameters total (17.6× smaller than teacher).

3.4 Logit-Based Distillation

The logit-based distillation loss uses temperature scaling:

$$L_{\text{logits}} = \text{KL}(\text{softmax}(z_{\text{student}}/T), \text{softmax}(z_{\text{teacher}}/T))$$

The temperature T is set to 4.0 for defect classification (higher temperature reveals more class relationship information) and $T = 1.0$ for regression tasks (thermal reconstruction, phase unwrapping), where the raw logit magnitudes encode uncertainty information.

The distillation loss is combined with the standard cross-entropy loss on hard labels:

$$L_{\text{total}} = \lambda \cdot L_{\text{distill}} + (1 - \lambda) \cdot L_{\text{ce}}$$

with $\lambda = 0.7$ (70% weight on distillation, 30% on hard labels).

3.5 Intermediate Feature Matching

Intermediate features are matched using a hint-based approach:

Hint selection. For the teacher, the output of encoder stage 3 (1/8 resolution) is designated as the hint layer. For the student, the output of encoder stage 3 is designated as the guided layer.

Projection adapter. A lightweight 1×1 convolution adapter projects the student's feature map to the same dimensionality as the teacher's hint features.

Feature matching loss.

$$L_{\text{feat}} = \|\text{adapter}(f_{\text{student}}) - f_{\text{teacher}}\|^2 / (C \cdot H \cdot W)$$

where C, H, W are the feature map dimensions.

3.6 Attention Transfer

Spatial attention maps are computed as the sum of squared activations across channels:

$$A_k = \sum_{c=1}^C (f_k^c)^2$$

where f_k^c is the activation at spatial location k for channel c . The attention transfer loss is:

$$L_{\text{attention}} = \|A_{\text{student}} - A_{\text{teacher}}\|^2 / (H \cdot W)$$

Applied at the output of each encoder stage (1/4, 1/8, and 1/16 resolution).

3.7 Defect-Aware Prioritization

Sample-level importance weights are assigned based on the teacher's confidence and the rarity of the defect type:

$$w_i = 1 + \beta \cdot (1 - P_{\text{max}}(\text{teacher} | x_i)) + \gamma \cdot R(y_i)$$

where P_{max} is the teacher's maximum class probability, $R(y_i)$ is the rarity score for the defect class y_i (higher for rarer defects), and $\beta = 2.0, \gamma = 1.5$ are weights.

The final loss is the importance-weighted sum:

$$L_{\text{total}} = \sum_i w_i \cdot L(x_i)$$

4. Simulation Experimental Results

4.1 Compression vs. Accuracy Tradeoff

Table 1 presents the accuracy of distilled student models at varying compression ratios (compression ratio = teacher parameters / student parameters).

Table 1 Defect detection accuracy vs. compression ratio

Compression Ratio	Student Parameters	mIoU (%)	vs. Teacher (pp)	FPS (Jetson Nano)
1× (no compression)	72M (teacher)	81.7%	baseline	3.2
4×	18M	80.9%	-0.8 pp	12.4
8×	9M	79.8%	-1.9 pp	27.6
12.7×	5.7M	78.6%	-3.1 pp	47.3
16.4×	4.4M	75.2%	-6.5 pp	68.1
24.1×	3.0M	68.7%	-13.0 pp	94.2

At 12.7× compression (5.7M parameters), the student achieves 78.6% mIoU—only 3.1 percentage points below the teacher—while running at 47.3 FPS on Jetson Nano (a 15× throughput improvement). At 16.4× compression, the student runs at 68.1 FPS but accuracy drops more noticeably.

4.2 Full Comparison: All Tasks

Table 2 presents the performance of the optimally compressed student models (12–13× compression) across all three inspection tasks.

Table 2 Distilled student model performance vs. teacher

Task	Metric	Teacher	Student (Compressed)	Degradation	Compression Ratio	Throughput Gain
Thermal reconstruction	MAE (K)	1.65	1.94 (+0.29)	17.6% worse	16.7×	14.3×
Phase unwrapping	RMSE (rad)	1.68	1.91 (+0.23)	13.7% worse	17.1×	15.8×
Defect detection	mIoU (%)	81.7	78.6 (-3.1 pp)	3.8% relative	12.7×	14.8×
Multitask (combined)	Combined score	100.0	95.4 (-4.6%)	4.6%	17.6×	16.2×

The optimally compressed models achieve 13–17× parameter reduction with accuracy degradation of only 3–18% (depending on task), enabling real-time throughput on embedded hardware.

4.3 Rare Defect Detection Preservation

A critical evaluation is whether rare defect detection capability is preserved during distillation. Table 3 presents per-defect-type accuracy for the teacher and a 12.7× compressed student.

Table 3 Per-defect-type accuracy: teacher vs. 12.7× compressed student

Defect Type	Prevalence in Training	Teacher Accuracy	Student Accuracy	Preservation Rate
Crack	18.4%	87.2%	84.1%	96.4%
Pit	24.1%	91.4%	89.8%	98.2%
Scratch	21.3%	85.7%	83.2%	97.1%
Contamination	19.8%	92.1%	90.4%	98.2%
Delamination	0.3%	68.4%	62.1%	90.8%
Coating blister	0.1%	54.3%	43.2%	79.6%

The distillation preserves 90–98% of teacher accuracy for common defects (crack, pit, scratch, contamination). For the rarest defects (delamination, coating blister), preservation rates drop to 79–91%—lower but still substantially better than training a student from scratch with the same rare-defect data.

4.4 Ablation: Distillation Component Contribution

Table 4 presents an ablation study isolating the contribution of each distillation component.

Table 4 Ablation: contribution of distillation components (defect detection)

Distillation Components	Student mIoU (%)	vs. Teacher Only
Hard labels only (no distillation)	61.3%	−20.4 pp
+ Logit-based distillation	73.8%	−7.9 pp
+ Feature matching	76.7%	−5.0 pp
+ Attention transfer	77.9%	−3.8 pp
+ Defect-aware prioritization	78.6%	−3.1 pp

Each distillation component independently improves performance. Logit-based distillation provides the largest single contribution (+12.5 pp over hard-label-only training). Feature matching and attention transfer contribute +2.9 and +1.2 pp respectively. Defect-aware prioritization contributes an additional +0.7 pp, primarily on rare defects.

4.5 Edge Hardware Deployment

Table 5 presents the throughput and latency of the distilled student models on various edge hardware platforms.

Table 5 Student model throughput on edge hardware

Hardware	Power (W)	Cost (USD)	Student-D FPS	Student-M FPS	Production Target Met?
Jetson Nano	10	~\$500	94.2	71.3	Yes (60 FPS)
Coral Edge TPU	2	~\$100	127.4	89.6	Yes (60 FPS)
Raspberry Pi 5 + Hailo	7	~\$200	83.7	62.4	Yes (60 FPS)
Mobile ARM (Snapdragon 865)	5	embedded	68.3	48.9	Near (50 FPS)
Intel Myriad X	1	~\$80	52.1	38.7	No (below 60 FPS)

All student models except the Intel Myriad X configuration meet the 60 FPS production throughput target on hardware priced below \$500. The most cost-effective deployment is the Coral Edge TPU at approximately \$100, achieving 127 FPS for defect detection.

4.6 Distillation vs. Training from Scratch

Table 6 directly compares distillation-trained students versus students trained from scratch with the same architecture and data.

Table 6 Distillation vs. training from scratch

Model	Defect Detection mIoU (%)	Thermal MAE (K)	Phase RMSE (rad)
Teacher	81.7	1.65	1.68
Student (from scratch)	61.3	2.47	2.83
Student (distilled)	78.6	1.94	1.91
Knowledge transfer gain	+17.3 pp	+0.53 K	+0.92 rad

Distillation provides massive improvements over training from scratch: +17.3 pp for defect detection, +0.53 K for thermal reconstruction, and +0.92 rad for phase unwrapping. This confirms that the teacher's dark knowledge—the ability to detect rare defects, to reason about subtle features—is not learnable from the training data alone and must be transferred through distillation.

5. Discussion

5.1 Practical Implications for Manufacturing

The knowledge distillation framework enables factories of all sizes to deploy state-of-the-art optical inspection accuracy. A small factory that can afford only a \$100 Coral Edge TPU can now achieve defect detection accuracy of 78.6% mIoU—compared to 61.3% for a model trained from scratch—through knowledge transfer from the best available teacher model. This democratizes access to high-accuracy inspection.

The operational simplicity of distilled models is equally important: they run on standard edge hardware without special software frameworks, with latency and memory requirements suitable for integration into existing inspection systems. A factory can receive a pre-trained distilled model as a single file and deploy it within hours.

5.2 Relationship to Prior Work

Knowledge distillation provides the final link in the deployment pipeline: Papers 1–3 established the task definitions and achieved high accuracy; Paper 9 addressed real-time edge deployment hardware infrastructure; Paper 17 (AutoML) designed efficient architectures from scratch; and this paper provides the mechanism for compressing the highest-accuracy teacher models into deployable form. The distillation framework is compatible with all architectures from the prior papers and can be applied whenever deployment hardware is more constrained than the training hardware.

5.3 Limitations

Several limitations should be noted. First, knowledge distillation requires access to the teacher model's logits and intermediate features, which are typically available for academic models but may not be available for proprietary commercial models. Second, the distillation is performed offline—requiring a training phase for the student model—which means that the student cannot adapt to new distributions after deployment (Paper 18's test-time adaptation handles this). Third, the current distillation framework does not preserve the teacher's uncertainty quantification capability (Paper 4); the student produces point predictions without calibrated uncertainty, which may limit its utility in safety-critical applications where uncertainty awareness is required.

6. Conclusion

This paper proposes a comprehensive knowledge distillation framework for optical surface inspection, compressing high-accuracy teacher models into compact student models for edge deployment.

The framework combines logit-based distillation, intermediate feature matching, attention transfer, and defect-aware prioritization to achieve 12–17× model compression while retaining 94–97% of teacher accuracy across all inspection tasks. Distilled student models achieve real-time throughput (60–127 FPS) on embedded hardware priced below \$500, compared to 3.2 FPS for the full teacher model on the same hardware.

Rare defect detection capability is preserved at 79–91% of teacher levels, and distillation-trained students substantially outperform students trained from scratch (+17.3 pp on defect detection), confirming that the teacher's dark knowledge of rare-event patterns is not learnable without distillation.

The proposed framework provides a practical pathway for deploying state-of-the-art optical inspection accuracy across the full range of manufacturing hardware, democratizing access to high-quality AI-powered quality control.

References

Hinton, G., Vinyals, O., & Dean, J. (2015). Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*.

Huang, H., Tang, J., Liu, T., & Huang, M. (2026). Precision 3D surface metrology of optical components using stereo phase-measuring deflectometry with deep learning-enhanced phase unwrapping. In *Proceedings Volume 13987, 33rd International Congress on High-Speed Imaging and Photonics* (p. 1398704). SPIE. <https://doi.org/10.1117/12.3093993>

Huang, H., Yang, Y., & Zhu, Y. (2023). Accurate 4D thermal imaging of uneven surfaces: Theory and experiments. *International Journal of Heat and Mass Transfer*, 216, 124580. <https://doi.org/10.1016/j.ijheatmasstransfer.2023.124580>

Romero, A., Ballas, N., Kahou, S. E., Chassang, A., Gatta, C., & Bengio, Y. (2015). FitNets: Hints for thin deep nets. In *International Conference on Learning Representations*. arXiv. <https://arxiv.org/abs/1412.6550>

Zagoruyko, S., & Komodakis, N. (2017). Paying more attention to attention: Improving the performance of convolutional neural networks via attention transfer. In *International Conference on Learning Representations*. arXiv. <https://arxiv.org/abs/1612.03928>

Malema. (2026a). Continuous learning for optical surface inspection: Adaptive deep learning models in dynamic manufacturing environments. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026b). Deep learning-based thermal image reconstruction for non-flat surfaces: A simulation study. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026c). Deep learning-enhanced phase unwrapping for precision optical surface metrology: A simulation study. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026d). Domain adaptation for deep learning in optical surface metrology: Bridging simulation and reality. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026e). Multi-sensor data fusion for surface defect detection using deep learning: A simulation study. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026f). Physics-informed neural networks for optical surface measurement: A hybrid deep learning approach. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026g). Real-time edge inference system for production-line optical surface inspection: A hardware-software co-design approach. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026h). Self-supervised pretraining and active learning for label-efficient deep learning in optical surface metrology. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026i). Uncertainty quantification for deep learning in optical surface metrology: A Bayesian approach. *Inclusive Growth and Governance Quarterly*, 2(1).

Malema. (2026j). Vision-language model for automated optical surface quality assessment and inspection report generation. *Inclusive Growth and Governance Quarterly*, 2(1).
