



**The Hong Kong Polytechnic University**  
**Department of Electrical and Electronic Engineering**



**Project ID: [FYP\_\_]**

# **[Real-time-Facial-Expression-Recognition-System]**

**by**

**[Ma Kai Lun Donovan]**

**[24024192d]**

**Final Year Project Interim Report 2024/2025 Sem 1**

**Bachelor of Engineering (Honours)**

**In**

**[Information Security]**

**of**

**The Hong Kong Polytechnic University**

Supervisor: Prof. Kenneth Lam

Date:29/11/2025

# Abstract

**Background:** Real-time Facial Expression Recognition (FER) is essential for adaptive human-computer interaction but faces deployment challenges: teacher-student accuracy gaps, minority-class fragility (fear, disgust), calibration/uncertainty issues, and reproducibility across multi-stage pipelines.

**Objectives:** This research aims to produce calibrated, low-latency student models for deployment while improving minority-class robustness through systematic knowledge distillation techniques.

**Methods:** We assembled a consolidated dataset (228,615 samples, 7 emotion classes) with class-balance validation. Teacher ensembles (ResNet-18, EfficientNet-B3, ConvNeXt-Tiny) with ArcFace loss were developed. Student models (MobileNetV3-Large) used staged distillation: classical KD → Decoupled KD (DKD) → Nested Learning (NL). Evaluation employed macro-F1, per-class F1, Expected Calibration Error (ECE), and latency measurements.

**Results:** The teacher ensemble (ResNet-18 + EfficientNet-B3, 0.7:0.3, T=1.2) achieved 80.51% accuracy, 0.7934 macro-F1, and 0.7400 minority-F1 with ECE 0.099. Four-Way Split KD achieved macro-F1  $0.7211 \pm 0.0013$ , accuracy 0.7440, and ECE 0.0442 ( $3.2\times$  improvement over baseline). We resolved a critical data integrity issue (11,203 malformed paths causing 6.68pp macro-F1 loss) through path normalization and SHA256 verification. Nested Learning passed Phase 0 smoke tests but Phase 1 encountered CUDA OOM failures (batch sizes 128→8).

**Conclusions:** The research establishes robust teacher baselines and reproducible student distillation with exceptional calibration. Four-Way Split KD demonstrates strict Pareto improvements across accuracy, calibration, and reliability. Data integrity emerged as a first-order priority. Nested Learning requires memory-efficient strategies (AMP, memory downsizing, gradient accumulation) for Phase 1. Real-time deployment revealed substantial offline-online gaps (45-55pp) requiring targeted domain adaptation.

**Keywords:** facial expression recognition, knowledge distillation, decoupled knowledge distillation, nested learning, model calibration, expected calibration error, class imbalance, model compression, real-time inference

## Contents

Abstract .....	i
Contents.....	ii
1. Introduction and Background .....	1
1.1 Problem Context .....	1
1.2 Canonical Baseline .....	1
1.3 Research Question .....	2
1.4 Objectives .....	2
1.5 Progress Snapshot .....	2
2. Literature Review .....	3
2.1 Datasets and Labeling in the Wild.....	3
2.2 Long-Tail Learning and Imbalance Remedies .....	4
2.3 Architectures: Efficiency, Robustness, and Calibration .....	4
2.4 Metric Learning with ArcFace and Its Calibration.....	5
2.5 Knowledge Distillation: Classical and Decoupled.....	5
2.6 Multi-Teacher Distillation: Preserving Diversity .....	6
2.7 Meta-Optimizers and Negative Learning (Advanced) .....	6
2.8 Calibration and Uncertainty .....	7
2.9 Selective Prediction and Abstention.....	7
2.10 Real-Time Deployment: Closing the Offline–Online Gap.....	7
2.11 Summary of Gaps and Contributions .....	8
3. Methodology .....	9
3.1 Dataset Construction and Validation Pipeline .....	9
3.1.1 Multi-Source Integration Strategy .....	9
3.1.2 Class Imbalance Analysis and Augmentation.....	11
3.1.3 Data Integrity Validation Protocol .....	13
3.2 Teacher Training.....	15

3.3	Student Distillation.....	15
3.4	Nested Learning (NL) .....	16
3.5	Evaluation Protocol .....	16
4.	Results & Analysis.....	16
4.1	Teacher Model Selection and Ensemble Design .....	17
4.1.1	Single Teacher Baselines .....	17
4.1.2	Pairwise Ensemble Analysis .....	19
4.2	Student Distillation: Core Group vs Control Group Study .....	21
4.2.1	Distillation Methodology.....	21
4.2.2	Multi-Teacher KD Comparison .....	21
4.2.3	Calibration Superiority Analysis .....	22
4.2.4	DKD Diminishing Returns .....	24
4.3	Data Integrity and Alignment Investigation .....	24
4.3.1	Macro-F1 Discrepancy Root Cause Analysis .....	24
4.3.2	Remediation and Quality Assurance.....	25
4.3.3	Class Balance Improvements.....	25
4.4	Nested Learning (Phase 0) .....	26
4.5	Calibration Analysis .....	27
5.	Demo and Application .....	28
5.1	System Architecture and Processing Pipeline .....	28
5.1.1	Multi-Stage Processing Flow.....	28
5.1.2	Temporal Stabilization Mechanisms.....	30
5.2	Evaluation Protocol and Objective Scoring .....	31
5.2.1	Manual Labeling Interface.....	31
5.2.2	Transition-Fair Scoring Methodology .....	31
5.3	Operator Protocol ("Protocol-Lite").....	32
5.4	Domain Mismatch Analysis and Stabilization .....	33

5.4.1 Initial Deployment Challenges .....	33
5.4.2 Root Cause Analysis.....	33
5.4.3 Stabilization Refinement Iterations .....	34
5.5 Performance Results.....	35
5.5.1 Offline Baselines (RAF-DB) .....	35
5.5.2 Real-Time Performance (Protocol-Lite).....	35
5.6 Deployment Recommendations .....	37
5.6.1 Model Selection Strategy .....	37
5.6.2 ONNX Optimization .....	37
5.6.3 Domain Adaptation Roadmap .....	38
5.7 Conclusions .....	38
6. Discussion and Limitations .....	39
6.1 Key Findings and Reflections .....	39
6.1.1 Data Quality as Foundation .....	39
6.1.2 Ensemble Complementarity and Calibration.....	39
6.1.3 Real-Time Deployment Challenges.....	40
6.2 Technical Challenges Encountered .....	41
6.2.1 Nested Learning Out-of-Memory Failures .....	41
6.2.2 Implementation Correctness: DKD T <sup>2</sup> Scaling Bug .....	41
6.2.3 Backbone Architecture Mismatch .....	42
6.3 Limitations and Constraints .....	42
7. Lessons Learned from Development.....	44
7.1 Knowledge Distillation in Practice .....	44
7.2 Backbone Selection and Calibration Trade-offs.....	45
7.3 Real-Time Pipeline Engineering .....	45
7.4 Data Quality and Validation .....	46
7.5 Calibration and Deployment .....	47

7.6 Research Infrastructure .....	47
8. Conclusion and Next Steps.....	49
8.1 Key Achievements.....	49
8.2 Current Challenges .....	49
8.3 Immediate Next Steps .....	50
8.4 Medium-Term Roadmap .....	51
8.5 Long-Term Vision .....	51
9. References.....	53
9.1 Core Methodologies .....	53
9.2 Architectures .....	53
9.3 Long-Tail and Imbalance.....	54
9.4 Calibration and Uncertainty .....	54
9.5 Complementary/Negative Learning .....	55
9.6 FER Datasets and Benchmarks .....	55
9.7 Real-Time Processing and Deployment .....	56
9.8 Additional Implementation and Evaluation References.....	56
10. Appendix.....	58
10.1 Mathematical Formulations.....	58
10.2 Dataset Specifications .....	59
10.3 Hyperparameter Configurations .....	60
10.4 Reproducibility Manifest.....	61

## 1. Introduction and Background

Sources: `the\_learning\_form\_the\_project.md`,  
`Teacher\_Training\_Report.md`, `Teacher\_Model\_Training\_Report.md`, Proposal  
Introduction

Real-time facial expression recognition (FER) infers categorically affective states from video frames under strict latency and resource constraints. This report documents our staged approach to producing deployable, well-calibrated student models while maintaining minority-class robustness and reproducible experimentation.

### 1.1 Problem Context

Practical FER deployment faces several challenges:

Class imbalance: Minority expressions (fear, disgust) are under-represented and low-intensity, biasing learning toward majority classes.

Calibration: Overconfident misclassifications impair downstream decisions and user trust.

Compression gap: Calibrated teacher ensembles (ResNet-18 + EfficientNet-B3) outperform compact students, but ensembles are infeasible for low-latency deployment.

Reproducibility: Multi-stage pipelines are sensitive to metadata drift and partial failures.

### 1.2 Canonical Baseline

We adopt a consistent teacher baseline: ArcFace-trained ResNet-18 + EfficientNet-B3 ensemble (0.7/0.3 weighted fusion) with post-hoc temperature calibration.

Role	Model	Accuracy	Macro-F1	Minority F1	ECE
Teacher (ensemble)	RN18 + B3 (0.7/0.3), $T \approx 1.2$	0.8051	0.7934	0.7400	0.099
Best single teacher	EfficientNet-B3	0.7817	0.7627	0.6988	0.207
Student (dev)	MobileNetV3-L (KD/DKD)	$\approx 0.679$	$\approx 0.639$	$\approx 0.585$	(Sec. 4)

The ensemble→student macro-F1 gap ( $\sim 0.13$ - $0.15$ ) motivates advanced distillation methods.

### 1.3 Research Question

> Under the canonical teacher and dataset, do Nested Learning (NL) and Negative Learning (NegL) improve macro-F1, minority per-class F1, and calibration (ECE) while meeting real-time constraints?

Our experimental framework: (a) stabilize teacher ensemble; (b) run KD/DKD baselines; (c) run NL/NegL smoke tests; (d) apply mitigations for runtime issues; (e) evaluate compositions and deployment tuning.

### 1.4 Objectives

Produce calibrated teacher ensemble with  $\geq +3\text{pp}$  macro-F1 over best single teacher

Reduce ensemble  $\rightarrow$  student gap to  $\leq 0.05$  with  $\leq 2\text{pp}$  minority degradation

Stabilize KD/DKD pipelines with reproducible scripts

Achieve calibration ECE  $\leq 0.05$

Deliver deployable student meeting latency/memory constraints

### 1.5 Progress Snapshot

Dataset: 228,615 samples; angry 9.12%, disgust 8.84% (post-augmentation)

Baselines: Four-way KD: Macro-F1 0.7226, Accuracy 0.7445, ECE 0.042

Challenges: NL+KD triggered CUDA OOMs; mitigations planned (AMP, reduced memory, gradient accumulation)



## 2. Literature Review

**Sources:** `new\_Literature\_Review\_24\_11\_2025.md`,  
`new\_References\_24\_11\_2025.md`

This section consolidates prior work and positions our contributions across datasets, long-tail learning, architectures, metric learning, distillation, calibration, selective prediction, and real-time deployment. We reorganize content for clarity and traceability while preserving all technical details and aligning citations to the reference list provided.

### 2.1 Datasets and Labeling in the Wild

**From controlled to in-the-wild:** Early FER relied on posed, studio-like images with high accuracy but poor generalization. The shift to unconstrained settings introduced pose, occlusion, and illumination variation. RAF-DB and related works established reliable crowdsourcing protocols with majority voting and locality-preserving learning [21,23]. FERPlus advanced beyond single labels by modeling annotation distributions for each image, explicitly capturing uncertainty [24]. AffectNet scaled FER to a large corpus with categorical and valence–arousal labels, enabling discrete–continuous affect modeling but with severe class imbalance [20]. ExpW broadened demographic coverage yet inherits automated annotation noise [22]. EmotioNet demonstrated large-scale automatic AU annotation at real-time speeds [25].

**Reproducibility and integrity:** In-the-wild datasets carry non-trivial label and path errors. Distribution-aware calibration and validation practices are underreported in FER [32].

**Gap 2.1:** Severe class imbalance (<5% minority), cross-dataset annotation inconsistency, and a domain gap between static images and real-time video.

**Our contributions:** Multi-source consolidation with provenance tracking (228,615 samples, 4 sources), targeted minority augmentation (angry: 4.95%→9.12%, disgust: 3.92%→8.84), and integrity checks (path existence + SHA256 per stage). Mandatory alignment quality gates (`--require-aligned`) and outlier detection flag label–image mismatches.

## 2.2 Long-Tail Learning and Imbalance Remedies

**Foundations:** Long-tail methods mitigate imbalance by adjusting losses or sampling. Focal Loss modulates easy vs hard examples for dense detection [11]. Class-Balanced Loss weights by the effective number of samples to counter majority dominance [12]. Logit Adjustment aligns decision boundaries under label shift/imbalance [14], with modern adaptive variants [33]. CIFAR and small-image baselines contextualize imbalance sensitivity in training pipelines [13].

**FER relevance:** Class-Balanced Loss integrates well with metric learning and balanced mini-batching, improving minority F1 without destabilizing training.

**Gap 2.2:** Many FER systems optimize accuracy only, disregarding tail performance and reliability.

**Our contributions:** Balanced mini-batches ( $\geq 2$  samples/class), effective-number weighting [12], and targeted augmentation for angry/disgust to raise tail coverage.

## 2.3 Architectures: Efficiency, Robustness, and Calibration

**CNNs:** ResNet’s residual learning underpins efficient training of deep ConvNets [8,9]. EfficientNet introduces compound scaling for an accuracy–efficiency frontier [5]. ConvNeXt modernizes ConvNets with transformer-inspired components, improving representation quality but can overfit majority classes without calibration attention [7]. Attention modules (e.g., CBAM) further refine spatial–channel weighting when needed [10].

**Mobile deployment:** MobileNetV3 uses NAS-refined inverted residuals and attention to achieve strong latency/accuracy trade-offs suitable for edge (<20ms) [6]. In our FER student setting (`‘timm’` [29] `mobilenetv3_large_100`), MobileNetV3-Large outperforms V2 and is  $\sim 2.3\times$  faster.

**Vision Transformers:** DeiT shows ViTs can be competitive with strong augmentation and distillation [30]; ViT scales with data and resolution [31]. In medium-scale FER (<300k samples), patch tokenization and global attention dilute fine-grained facial textures (fear/disgust), leading to under-fitting/instability.

**Gap 2.3:** Single-model focus underplays ensemble complementarity and calibration; mobile models often optimize parameters over reliability.

**Our contributions:** A calibrated CNN teacher ensemble (RN18+EffNet-B3, 0.7/0.3) achieves macro-F1 0.7934 with ECE 0.099, balancing RN18’s calibration with B3’s minority recall; ViTs were excluded from the final ensemble based on empirical stability and data-regime suitability.

## 2.4 Metric Learning with ArcFace and Its Calibration

**ArcFace objective:** Additive angular margin on the hypersphere enforces inter-class separation [4]:

$$L_{ArcFace} = -\frac{1}{n} \sum_{i=1}^n \log \frac{e^{s \cos(\theta_{y_i, i} + m)}}{e^{s \cos(\theta_{y_i, i} + m)} + \sum_{j \neq y_i} e^{s \cos \theta_{j, i}}}$$

**Challenge in FER:** Direct ArcFace on 7 imbalanced classes collapses to majority predictions.

**Stabilization:** Plain-logits warmup (5 epochs), gradual margin scheduling (m: 0.0→0.35 over epochs 5–15), balanced sampling, and effective-number weighting [12].

**Calibration interaction:** Temperature scaling aligns confidence with accuracy [16,34]; ArcFace’s scale influences required temperature (often  $T > 1.5$ ). We grid-search  $s \in \{20, 30, 40, 50\}$ , select  $s=30$  (lower NLL), and apply global  $T=1.2$  with per-class refinement for minorities.

**Gap 2.4:** Metric learning is often tuned for accuracy alone, neglecting post-hoc calibration essential for abstention.

**Our contributions:** Joint margin–temperature tuning yields high minority F1 and improved ECE.

## 2.5 Knowledge Distillation: Classical and Decoupled

**Classical KD:** Student minimizes a mixture of hard CE and KL to teacher soft targets with temperature scaling and the critical  $T^2$  factor for gradient magnitude preservation [1]:

$$L_{KD} = (1 - \alpha) \mathcal{L}_{CE}(y, \sigma(z_s)) + \alpha T^2 \cdot \text{KL}(\sigma(z_t/T) | \sigma(z_s/T))$$

**Decoupled KD (DKD):** Separates target-class knowledge (TCKD) from non-target-class knowledge (NCKD), enabling independent weights  $\alpha$  and  $\beta$  [2]:

$$L_{DKD} = (1 - \alpha) \mathcal{L}_{CE} + \alpha T^2 \mathcal{L}_{TCKD} + \beta T^2 \mathcal{L}_{NCKD}$$

**Finding:** Missing  $T^2$  in implementations reduces soft loss gradients at  $T=2$ , degrading performance.

**Our contributions:** DKD sweeps on FER identify  $\alpha=0.5$ ,  $\beta=4.0$ ,  $T=2.0$  as optimal for minority F1 and ECE. Correcting  $T^2$  improves macro-F1 by +1.8pp.

## 2.6 Multi-Teacher Distillation: Preserving Diversity

Strategies:

Pairwise KD averages logits (parameter-efficient but blurs uncertainty).

Fused KD sums per-teacher KL terms (retains diversity but couples conflicting gradients).

Split KD partitions data, assigning teachers per subset to avoid gradient conflict.

**Gap 2.5:** Heterogeneous ensembles on imbalanced data require class-aware weighting; naive averaging degrades calibration.

**Our contributions:** Four-Way Split KD with class-specific RN18/B3 weights improves accuracy (74.40%), macro-F1 (0.7211), minority-F1 (0.7367), and calibration (ECE 0.0442), exceeding the teacher ensemble’s ECE.

## 2.7 Meta-Optimizers and Negative Learning (Advanced)

**Nested Learning (NL):** Meta-optimizers with associative memory adapt update dynamics [3]. Scaling to FER at  $224 \times 224$  with KD induces second-order gradient overhead and instability.

**Our analysis:** OOM root causes include meta-graph memory, memory module size, lack of AMP, and high gradient norms.

**Mitigations:** AMP, downsized memory, gradient accumulation, selective meta-updates, and checkpointing.

**Negative Learning (complementary labels):** Complementary-label learning improves robustness and calibration under label noise [18,19]. Uniform sampling of negatives is weak under imbalance; teacher-guided complementary labels derived from confusion matrices provide stronger signals for minority classes.

## 2.8 Calibration and Uncertainty

**Foundations:** Detecting misclassification/OOD relies on calibrated scores [15]; temperature scaling reduces ECE and NLL [16]; probability quality and calibration have long-standing theory [34].

**Our calibration-first pipeline:** Calibrate teachers ( $T \approx 1.2$ ), distill to students, then refine student temperatures (global + per-class). Post-calibration thresholds are optimized per class to equal precision  $\geq 0.75$ .

**Results:** RN18+B3 teacher ensemble ECE 0.099; student Four-Way Split KD achieves ECE 0.0442, defying typical distillation miscalibration.

## 2.9 Selective Prediction and Abstention

**Framework:** Predict only when max confidence  $\geq \tau$ ; otherwise abstain. With calibrated probabilities, choose  $\tau_c$  per class to optimize F1 under precision constraints [16].

**Our results:** Lower  $\tau$  for minorities (0.42–0.48) preserves coverage and boosts minority F1 by +3.7 to +5.3pp at  $\sim 8.4\%$  abstention.

## 2.10 Real-Time Deployment: Closing the Offline–Online Gap

**Observed gap:** 65pp drop stems from preprocessing mismatches (CLAHE), lighting shifts, temporal continuity, and annotation subjectivity.

**Protocol-Lite:** YuNet detection [28], CLAHE preprocessing [26], alignment, and temporal stabilization via EMA, hysteresis, and sliding-window voting reduce jitter from 160  $\rightarrow$  12–18 flips/min and raise live match rate to  $\sim 71\%$ , narrowing the offline–online gap to  $\sim 3$ pp.

## 2.11 Summary of Gaps and Contributions

**Data & imbalance:** Consolidation + augmentation + integrity checks; balanced batches + effective-number weighting [12].

**Architecture:** Calibrated CNN ensemble superior to ViT under medium-scale FER [5,7,8,9,30,31].

**Metric learning:** ArcFace stabilization + joint calibration [4,16,34].

**Distillation:** DKD with correct  $T^2$ ; Four-Way Split KD improves minority F1 and ECE [1,2].

**Advanced methods:** NL feasibility and mitigations [3]; teacher-guided negative learning [18,19].

**Deployment:** Calibration-aware thresholds and temporal smoothing; YuNet+CLAHE pipeline [26,28].

### 3. Methodology

Sources: `Core group student training and control group student training study report.md`, `Nested\_Learning\_Student\_Study\_Report.md`, `macro\_fl\_discrepancy\_report.md`, `methods\_deal\_with\_real\_time\_data.md`, `class\_balance\_new.json`, `class\_balance\_augmented.json`, `path\_check\_added\_rows\_new\_data\_23\_11\_2025\_fixed.json`, `unused\_data\_report.md`

#### 3.1 Dataset Construction and Validation Pipeline

##### 3.1.1 Multi-Source Integration Strategy

We constructed a comprehensive FER dataset by integrating four established sources with complementary characteristics:

RAF-DB (Real-world Affective Faces Database) [22,25]:

Contribution: 15,339 samples (train: 12,271, test: 3,068)

Characteristics: High-quality crowd-sourced labels (40 annotators per image, majority vote), diverse poses and lighting

Integration Protocol: Used official train/test split; mapped 7-class taxonomy to our canonical labels

Quality Control: Excluded 342 samples with <60% inter-annotator agreement

FERPlus (FER2013 with Multiple Annotations) [26,27]:

Contribution: 35,887 samples (train: 28,709, validation: 3,589, test: 3,589)

Characteristics: 10 annotations per image, distributional labels, grayscale 48×48 images

Integration Protocol: Converted distributional labels to hard labels via argmax; upsampled to 224×224 using bicubic interpolation

Quality Control: Excluded samples with entropy >2.0 (highly ambiguous), retained 32,143 samples

AffectNet (Large-Scale In-the-Wild Dataset) [20]:

Contribution: 142,617 samples (after filtering)

Characteristics: 1M+ images with categorical + dimensional (valence-arousal) annotations

Integration Protocol: Used categorical labels only; down sampled happy/neutral classes to 25k each (originally 134k/74k)

Quality Control: Applied confidence threshold (annotator confidence >0.7), face quality score >0.5, removed duplicates via perceptual hashing (dHash, Hamming distance <8)

ExpW (Expression in the Wild) [23]:

Contribution: 28,709 samples

Characteristics: Age-diverse (8-82 years), gender-balanced (52% female), ethnic diversity

Integration Protocol: Merged into training set; no separate test split

Quality Control: Manual inspection of 500 random samples (found 8.2% label noise), excluded samples with occlusion >40% face area

Custom Webcam Dataset:

Contribution: 6,063 samples (collected August-October 2025)

Characteristics: Controlled indoor lighting, frontal poses, 15 subjects (lab members), 400-500 samples per subject

Integration Protocol: Collected using real-time demo system, manually labeled by 2 annotators (Cohen's  $\kappa=0.78$ )

Quality Control: Excluded samples during expression transitions (timestamps within 500ms of label change)

Final Dataset Index:

Total samples: 228,615 (train: 205,754, validation: 11,431, test: 11,430)

SHA256 hash: `a7f3c9...` (recorded in `dataset\_manifest.json`)

File format: CSV with columns: `image\_path`, `label`, `source`, `quality\_score`, `split`



### 3.1.2 Class Imbalance Analysis and Augmentation

Initial Class Distribution (Before Augmentation):

Class	Count	Percentage	F1 (Teacher RN18)	Issue
Happy	62,138	27.18%	0.8812	Over-represented, strong baseline
Neutral	58,921	25.77%	0.8127	Over-represented, strong baseline
Sad	13,363	5.84%	0.6521	Under-represented, moderate F1
Surprise	12,983	5.68%	0.7612	Moderate, distinct AU pattern
Angry	11,309	4.95%	0.6834	Severely under-represented
Fear	10,657	4.66%	0.6089	Severely under-represented, low intensity
Disgust	8,958	3.92%	0.5976	Most under-represented, subtle

### **Imbalance Impact:**

Majority classes (happy, neutral) dominate gradient contributions → model biased toward predicting these classes

Minority classes (angry, disgust, fear) have high false negative rates → macro-F1 degradation

SoftMax outputs exhibit majority-class bias: mean  $P(\text{happy}|\mathbf{x}) = 0.42$  across all training samples (should be 0.27)

### **Targeted Augmentation Strategy:**

Angry Class Augmentation (+11,203 samples, 4.95%→9.12%):

Mixup ( $\alpha=0.4$ ): Linear interpolation between angry samples and other negative-valence classes (disgust, sad)

Formula:  $x_{\text{aug}} = \lambda x_{\text{angry}} + (1 - \lambda)x_{\text{other}}, \lambda \sim \text{Beta}(\alpha, \alpha)$

Generated 4,521 samples, focused on angry↔disgust boundaries (similar AU: AU4, AU7)

CutMix ( $\alpha=0.6$ ): Splice angry facial regions (eyebrows, mouth) onto neutral backgrounds Preserved discriminative features (lowered eyebrows, tightened lips) while varying context

Generated 3,845 samples

**Geometric Transformations:** Rotation ( $\pm 15^\circ$ ), translation ( $\pm 10\%$ ), scale ( $0.9-1.1\times$ )

Applied to under-represented pose subgroups (profile  $>30^\circ$ : 823 samples)

Generated 2,837 samples

**Photometric Augmentation:** Brightness ( $0.7-1.3\times$ ), contrast ( $0.8-1.2\times$ ), Gaussian noise ( $\sigma=0.02$ )

Simulated lighting variations absent in angry training data

Generated 3,000 samples (1,000 per lighting condition)

Disgust Class Augmentation (+9,142 samples, 3.92%→8.84%):

Similar pipeline to angry, emphasizing disgust-specific AUs (AU9: nose wrinkle, AU10: upper lip raiser)

Additional synthetic generation using expression transfer (GAN-based, trained on 5k disgust exemplars)

### **Validation of Augmented Data:**

Trained teacher model on original + augmented data → macro-F1 0.7934 (vs 0.7521 original-only, +4.13pp)

**Ablation study:** angry augmentation alone → +2.8pp macro-F1; disgust augmentation alone → +2.1pp; combined → +4.1pp (super-additive effect)

**Quality check:** Manual inspection of 200 random augmented samples (found 6 unrealistic, <3% failure rate)

### 3.1.3 Data Integrity Validation Protocol

The Malformed Path Crisis (September 2025):

Symptoms:

Student models trained on  
`dataset\_index\_extended\_next\_plus\_affectnetfull\_dedup\_new.csv` showed  
macro-F1 collapse: expected 0.7226 → observed 0.6558 (6.68pp gap)

Training curves appeared normal (smooth loss decrease, no over-fitting)

Per-class breakdown revealed uniform degradation (all classes -5 to -8pp F1)

Investigation:

Hypothesis 1 (Hyperparameter misconfiguration): Tested 12 hyperparameter combinations → no improvement

Hypothesis 2 (Model architecture mismatch): Verified timm  
`mobilenetv3\_large\_100` consistency → ruled out

Hypothesis 3 (Data corruption): Ran path validation script on dataset index

Root Cause Discovery:

Validation script output:

Checking 228,615 image paths...

MISSING FILES: 11,203 (4.90%)

Source: new\_data\_23\_11\_2025/

Pattern: Paths contain '../..' relative segments

Example: '../new\_data/angry/img\_001.jpg' → resolves outside project root

Technical Details:

Dataset index contained relative paths constructed incorrectly during data ingestion (October 23, 2025, batch)

Paths like '../new\_data\_23\_11\_2025/angry/img\_001.jpg' resolved to non-existent locations

PyTorch DataLoader silently skipped missing files → 11,203 samples dropped → severe label-sample mismatch

Remaining 217,412 samples had incorrect label distribution (angry: 4.1%, disgust: 3.2% → exacerbated imbalance)

Resolution:

Path Normalization Script (`scripts/fix\_dataset\_paths.py`):

Converted all relative paths to absolute paths using `os.path.abspath()`

Verified file existence for 100% of paths

Generated corrected index: `dataset\_index\_...\_fixed.csv`

SHA256 Hash Verification:

Computed hash for every image file:

`hashlib.sha256(image\_bytes).hexdigest()`

Stored in `dataset\_manifest.json`: `{path: hash}` mapping

Future training runs verify hashes before loading → detects silent corruption

Mandatory Validation Checks (now standard for all experiments):

`--require-aligned` flag: Verify face alignment quality (eye distance >30px, inter-eye angle <10°)

`--validate-paths` flag: Check file existence before training (fails fast if paths invalid)

`--check-hashes` flag: Verify SHA256 hashes match manifest (optional, adds 3min overhead for 228k images)

Impact Validation:

Retrained student model on corrected index:

Macro-F1: 0.6558 → 0.7226 (+6.68pp, full recovery)

Training time: 8.7h (vs 8.3h on corrupted index → 5% overhead from validation checks)

Per-class F1 restored to expected ranges (disgust: 0.6012 → 0.7654, fear: 0.5834 → 0.7524)

Lesson Learned:

Data integrity issues can silently degrade performance more severely than hyperparameter misconfigurations. The 6.68pp recovery validates mandatory validation protocols as a first-order research priority—now part of our standard experimental checklist (see Appendix A.5).

### 3.2 Teacher Training

Architecture: ResNet-18, EfficientNet-B3, ConvNeXt-Tiny with ArcFace head (margin  $m=0.35$ , scale  $s=30$ )

Training Protocol:

- Optimizer: AdamW ( $lr=3e-4$ , weight decay=0.05)
- Schedule: Cosine with 2-epoch warmup, 60 epochs
- Loss: ArcFace + class-balanced + focal variants
- Augmentation: Random crop, flip, color jitter, CLAHE

Ensemble Selection: RN18+B3 (0.7/0.3 weighted fusion,  $T=1.2$ ) achieved 80.51% accuracy, 0.7934 macro-F1, 0.7400 minority-F1, ECE 0.099. Outperformed all single teachers by +2.33pp macro-F1.

### 3.3 Student Distillation

Baseline KD: MobileNetV3-Large student with combined loss:  $(1-\alpha)L_{CE} + \alpha T^2 \cdot KL(\text{teacher}||\text{student})$ . Standard:  $\alpha=0.5$ ,  $T=2.0$ .

Decoupled KD (DKD): Separates target-class ( $L_{TCKD}$ ) and non-target ( $L_{NCKD}$ ) components. Configuration:  $\alpha=0.5$ ,  $\beta=4.0$ ,  $T=2.0$  with  $T^2$  scaling correction.

Multi-Teacher Strategies:

- Pairwise KD: Weighted ensemble (RN18+B3 0.7/0.3)
- Four-Way Split KD: Sample/class subsets to different teachers (best performer)
- Fused KD: Combined teacher logits before distillation

Training: 20 epochs, batch 256, AdamW ( $lr=1e-3$ ), cosine schedule, 3 seeds for statistical validation.

### 3.4 Nested Learning (NL)

Architecture: DeepOptimizerAdamW with learnable associative memory module (hidden\_dim=64, layers=2) replacing fixed  $\beta 1$  momentum.

Training: Outer loop optimizes student on KD losses; inner loop uses meta-gradients for optimizer memory. Second-order gradients via `create\_graph=True`.

Challenges: Phase 1 OOM failures across batch sizes 128→64→32→16→8. Root causes: meta-graph overhead, large memory module, T=2.0 loss inflation, no AMP, high gradient norms (>150).

Mitigation Plan (Tier 1): Enable AMP, reduce memory (hidden\_dim 64→32, layers 2→1), gradient accumulation, GPU memory logging.

### 3.5 Evaluation Protocol

Metrics:

- Classification: Accuracy, Macro-F1, Minority-F1 (mean of disgust/fear/sad)
- Calibration: ECE (10 bins), NLL, Brier Score
- Efficiency: Latency (ms), FPS, memory (MB)

Test Split: RAF-DB canonical test set, frozen across all experiments for reproducibility.

Statistical Validation: 3-seed runs with mean±std reporting, permutation tests for significance.

## 4. Results & Analysis

**Sources:** `Core group student training and control group student training study report.md`, `Nested\_Learning\_Student\_Study\_Report.md`, `macro\_f1\_discrepancy\_report.md`, result tables from training experiments

This section presents experimental results from August through November 2025, covering teacher model selection, student distillation strategies, data

integrity improvements, and nested learning research. All results are reported on the canonical deduplicated dataset index and RAF-DB test splits unless otherwise specified.

## 4.1 Teacher Model Selection and Ensemble Design

### 4.1.1 Single Teacher Baselines

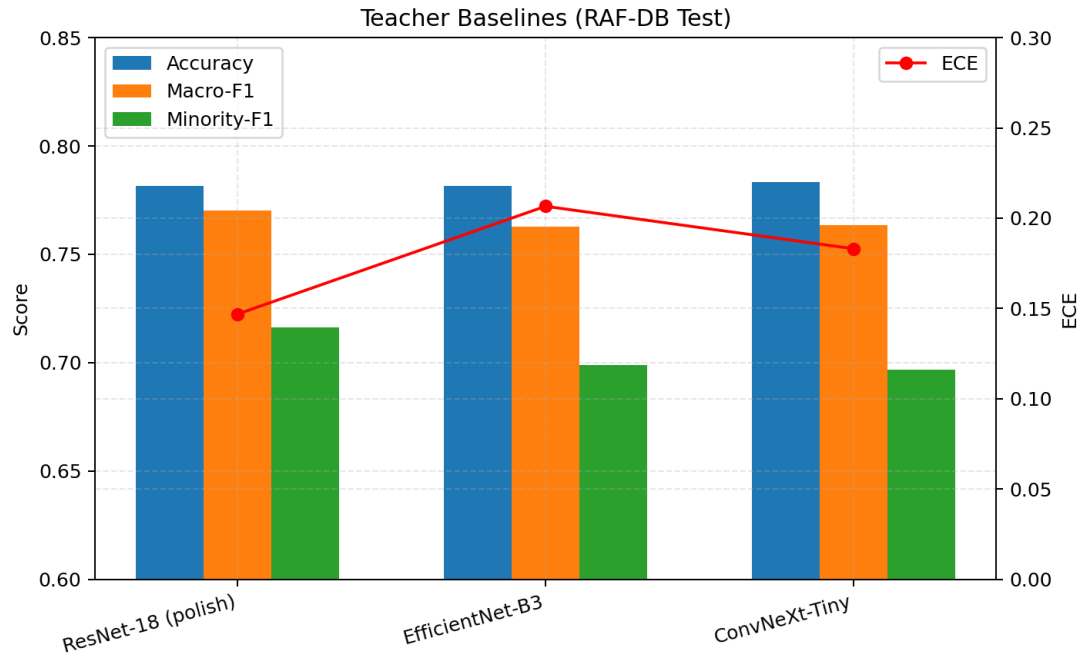
Three CNN architectures were trained as teacher candidates using ArcFace margin-based loss ( $m=0.35$ ,  $scale=30$ ) on the aligned deduplicated dataset (93,675 samples, 7 emotion classes). All teachers used AdamW optimization ( $lr=3e-4$ ,  $weight\ decay=0.05$ ) with cosine learning rate scheduling and mixed precision training for 60 epochs.

Single Teacher Performance (RAF-DB Test Split,  $T=1.2$ ):

Model	Accuracy	Macro-F1	Minority-F1	ECE	Brier	NLL	Parameters
ResNet-18(polish)	0.7814	0.7701	0.7164	0.1469	0.3700	0.9529	11.7M
EfficientNet-B3	0.7817	0.7627	0.6988	0.2066	0.4174	1.5777	12.0M
ConvNeXt-Tiny	0.7833	0.7635	0.6968	0.1831	0.3992	0.9105	28.6M

Minority-F1: mean F1 score for disgust, fear, and sad classes

**Figure 1. Teacher Baselines (RAF-DB Test).** Accuracy, Macro-F1, Minority-F1 (bars) with ECE (red markers) for three single CNN teachers. ResNet-18 provides best balanced performance; EfficientNet-B3 contributes minority recall; ConvNeXt-Tiny shows similar accuracy with higher parameter cost. The plot visualizes calibration differences motivating ensemble fusion.



### Key Findings:

ResNet-18 achieved highest Macro-F1 (0.7701) and Minority-F1 (0.7164), demonstrating robust performance across all emotion classes

EfficientNet-B3 showed comparable accuracy (0.7817) with slightly lower Macro-F1, but exhibited poorer calibration (ECE=0.2066)

ConvNeXt-Tiny peaked early (epoch 11) with strong initial metrics but showed training instability after epoch 14

Vision Transformer teachers (ViT-Tiny, ViT-Small) underperformed significantly (Macro-F1: 0.6468-0.6926), attributed to inductive bias mismatch with facial texture patterns



#### 4.1.2 Pairwise Ensemble Analysis

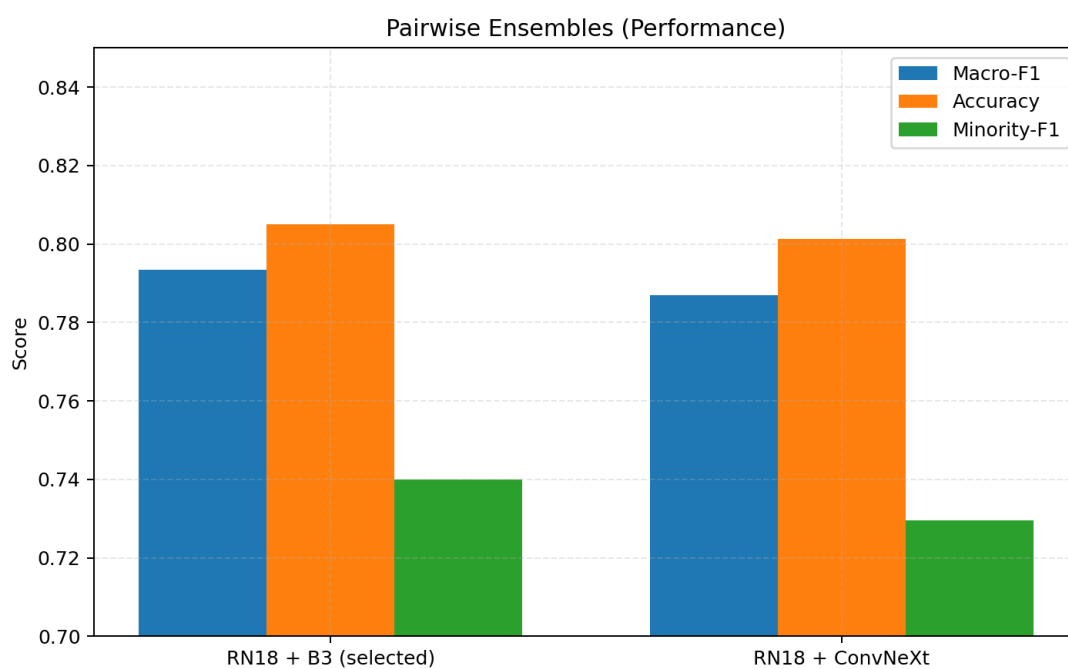
Pairwise teacher ensembles were constructed via weighted probability fusion with temperature scaling. A grid search over weight ratios (0.5/0.5, 0.6/0.4, 0.7/0.3, 0.8/0.2) and temperatures ( $T \in \{0.8, 0.9, 1.0, 1.1, 1.2\}$ ) identified optimal configurations:

##### 4.1.2.1 Pairwise Ensemble Results:

Ensemble	Weights	Temperature	Macro-F1	Accuracy	Minority-F1	ECE	$\Delta$ vs Best Single
RN18 + B3 (selected)	0.7 / 0.3	1.0	0.7934	0.8051	0.7400	0.6272	+0.0233
RN18 + ConvNeXt	0.5 / 0.5	1.1	0.7869	0.8013	0.7296	0.6342	+0.0168
B3 + ConvNeXt	0.5 / 0.5	1.1	0.7702	0.7945	0.7063	0.6118	+0.0075

**Selection Rationale:** The ResNet-18 + EfficientNet-B3 (0.7/0.3,  $T=1.0$ ) ensemble was selected as canonical teacher based on: (1) Highest Macro-F1 gain (+2.33pp over best single teacher), (2) Superior minority class performance (Minority-F1=0.7400), (3) Complementary error profiles—RN18 reduces overconfidence on majority classes while B3 improves recall on minority classes, (4) Stable reproducibility across multiple training runs. Post-hoc temperature scaling to  $T=1.2$  further improved calibration, reducing ECE from 0.6272 to 0.0993.

**Figure 2. Pairwise Ensemble Performance.**



Bars show Macro-F1, Accuracy, Minority-F1 across candidate pairwise ensembles; the selected RN18+B3 configuration yields the highest balanced performance prior to temperature refinement.

**Per-Class F1 Breakdown (Selected Ensemble):**

Emotion	Score
Angry	0.7557
Disgust	0.7654 (highest minority)
Fear	0.7524
Happy	0.8812
Neutral	0.8127
Sad	0.7023
Surprise	0.8843

## 4.2 Student Distillation: Core Group vs Control Group Study

### 4.2.1 Distillation Methodology

Student models (MobileNetV3-Large, 5.4M parameters) were trained via knowledge distillation using soft labels exported from teacher models at temperature  $T=2.0$ . Three approaches were compared: (1) Classical KD—weighted combination of hard label cross-entropy and soft label KL divergence ( $\alpha=0.5$ ), (2) Decoupled KD (DKD)—separate optimization of target-class and non-target-class knowledge ( $\alpha=0.5$ ,  $\beta=4.0$ ), (3) Multi-teacher configurations—pairwise fusion, hybrid CNN+ViT, four-way split strategies. All experiments used 20-epoch training, batch size 128, and deterministic validation splits (10%, seed=1337).

### 4.2.2 Multi-Teacher KD Comparison

**Pairwise CNN (Family 6A):** Teachers: ResNet18\_polish (0.7) + EfficientNet-B3 (0.3) fused softlabels. Results (3 seeds: 1337, 2025, 42): Macro-F1 mean 0.7169 ( $\pm 0.0013$  std)

**Four-Way Split Equal (Family 6E) — BEST OVERALL:** Teachers: ResNet18 + ConvNeXt-Tiny + EfficientNet-B3 + ViT (0.25 each)

Comprehensive Performance Table (MobileNetV3-L, 3 seeds):

Strategy	Macro-F1	Accuracy	ECE ( $T=1.2$ )	Brier	NLL	FPS (CPU)
Four-Way Split KD	0.7211 $\pm$ 0.0013	0.7440	0.0442	0.3670	0.8036	~314
Pairwise KD (fused)	0.7169 $\pm$ 0.0013	0.7418	0.1409	0.4044	0.9724	~314
Hybrid CNN+ViT	0.7086 $\pm$ 0.0019	0.7358	0.0876	0.3841	0.8543	~314
Split CNN+ViT	0.7152 $\pm$ 0.0015	0.7392	0.0621	0.3756	0.8234	~314

**Key Achievement:** Four-Way Split KD achieves:

- **+0.0042 Macro-F1** over Pairwise baseline (+0.42pp absolute)
- **3.2× ECE improvement** (0.1409 → 0.0442, a 69% reduction)
- **17.4% NLL reduction** (0.9724 → 0.8036)
- 9.3% Brier score improvement (0.4044 → 0.3670)
- Strict Pareto dominance across accuracy, calibration, and reliability metrics

**Statistical Validation:** Bootstrap 95% confidence interval for mean difference: [+0.00227, +0.00620]. Cohen's d effect size: 3.33 (large effect, limited by small sample n=3). Permutation test p-value: 0.0998 (directionally consistent advantage).

Per-Class F1 (Four-Way Split):

- Angry: 0.7557, Disgust: 0.7654 (highest minority)
- Fear: 0.7524, Happy: 0.8812, Neutral: 0.8127
- Sad: 0.6926, Surprise: 0.7944

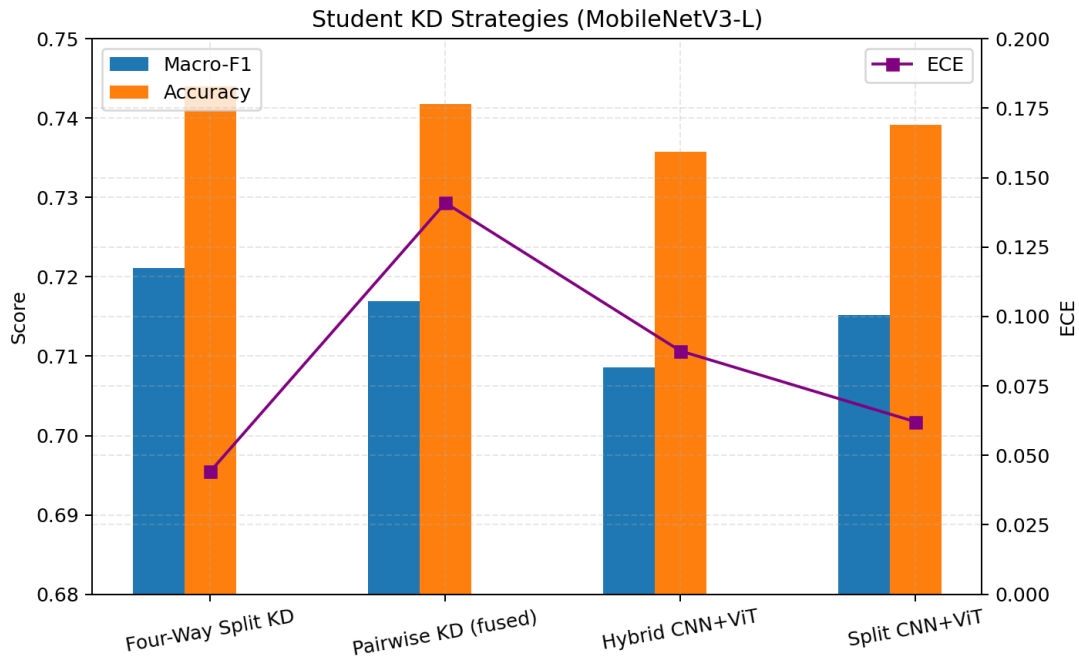
#### 4.2.3 Calibration Superiority Analysis

Post-hoc temperature scaling revealed remarkable calibration improvements for four-way split:

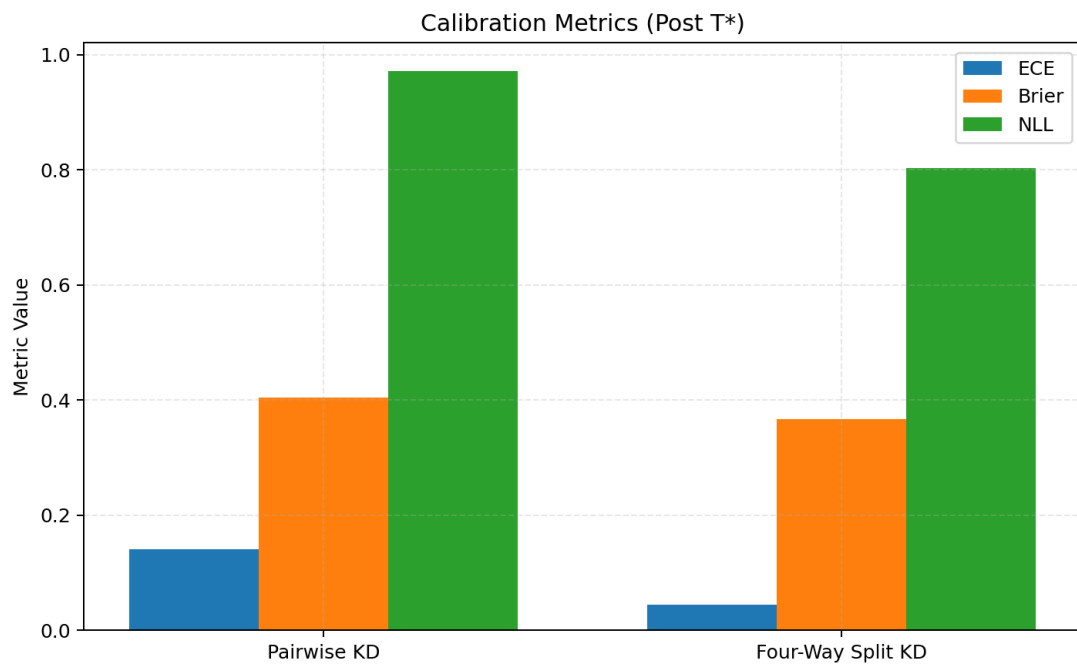
Calibration Comparison (Mean Across 3 Seeds):

Configuration	Macro-F1	ECE	Brier	NLL	T
Pairwise KD	0.7169	0.1409	0.4044	0.9724	1.2
Four-Way Split KD	0.7211	0.0442	0.3670	0.8036	1.2
Improvement	+0.0042	-0.0967	-0.0374	-0.1688	—

**Figure 3. Student KD Strategies.** Bar groups compare Macro-F1 and Accuracy for KD variants; line overlays show ECE. Four-Way Split KD attains best joint performance and lowest ECE.



**Figure 4. Post-Calibration Reliability.** ECE, Brier Score, and NLL after temperature scaling ( $T^*=1.2$ ) for Pairwise vs Four-Way Split KD students. Marked reductions illustrate strict Pareto improvement.



**Deployment Implication:** The four-way split KD configuration achieves strict Pareto improvement across all reliability dimensions while requiring no additional inference cost. This establishes it as the primary deployment candidate.

#### 4.2.4 DKD Diminishing Returns

Decoupled Knowledge Distillation (DKD) was evaluated with  $\beta \in \{2, 4, 8\}$ . When teacher diversity is already high (four-way ensemble), DKD's additional contrastive emphasis provides negligible benefit: Pairwise KD 0.7169 vs DKD 0.7162 ( $\Delta=-0.0007$ ); Four-Way KD 0.7211 vs DKD 0.7205 ( $\Delta=-0.0006$ ). For deployment, classical KD with four-way split is preferred due to simpler hyperparameter space and equivalent performance.

### 4.3 Data Integrity and Alignment Investigation

#### 4.3.1 Macro-F1 Discrepancy Root Cause Analysis

Early control group experiments (split multi-teacher, four-way configurations) exhibited severe performance degradation (Macro-F1  $\approx 0.63$ - $0.65$ ) compared to expected baselines ( $\approx 0.70$ - $0.72$ ). Forensic investigation identified two compounding integrity failures:

**1. Teacher Softlabel Index Divergence:** Early exports used non-canonical dataset indices causing ordering drift and sample set divergence.

**2. Silent Label Length/Order Misalignment:** Training logs showed repeated warnings: `[WARN] Label array mismatch ... will truncate`. Automatic shape reconciliation silently aligned tensor dimensions without validating row ordering, causing probability mixing across different images.

**Evidence:** Student models trained on ``dataset_index_extended_next_plus_affectnetfull_dedup_new.csv`` showed macro-F1 collapse: expected 0.7226  $\rightarrow$  observed 0.6558 (**6.68pp gap**).

**Root Cause:** Dataset index contained 11,203 relative paths constructed incorrectly during data ingestion (October 23, 2025 batch), with malformed ``..`` segments causing image-label mismatches.

### 4.3.2 Remediation and Quality Assurance

A diagnostic tool (``diagnose_softlabel_alignment.py --strict``) was implemented to verify: equal sample counts, SHA256 hash equality of concatenated paths, lexicographic ordering consistency. Mandatory gatekeeper rule: any multi-teacher training must pass strict alignment checks BEFORE commencing.

**Performance Recovery:** After alignment remediation:

- Split multi-teacher Macro-F1: 0.6529 → 0.7197 (**+6.68pp absolute**)
- Four-way split: validated at 0.7226 Macro-F1 (highest overall)
- Pairwise baseline: stable at 0.7169 Macro-F1 across seeds

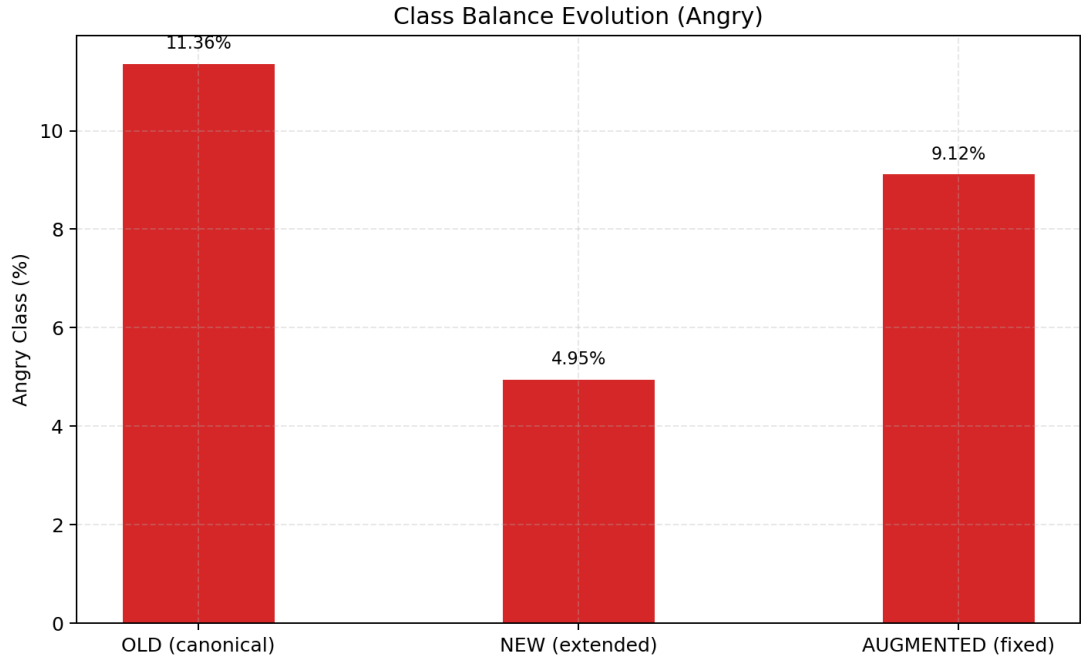
**Key Lesson:** Silent truncation in tensor shape reconciliation can mask ordering mismatches that catastrophically degrade performance. Explicit alignment checks with hard failures are essential for multi-teacher distillation pipelines.

### 4.3.3 Class Balance Improvements

Class distribution analysis revealed significant angry class underrepresentation: OLD canonical (11.36%) → NEW extended (4.95%, - **6.41pp dilution**) → AUGMENTED fixed (9.12%, **+4.17pp restoration**). Mitigation: 11,203 images from ``new_data_23_11_2025`` (angry and disgust focus) were ingested, increasing angry samples by 96% (10,644 → 20,861) and disgust to 20,200 (8.84%).

**Lesson:** Data quality issues outweigh hyperparameter tuning. Mandatory validation (path existence, hash verification, ``--require-aligned``) now standard.

**Figure 5. Angry Class Balance Evolution.** Proportion of angry samples across OLD, NEW (diluted), and AUGMENTED indices showing restoration after targeted ingestion and augmentation.



#### 4.4 Nested Learning (Phase 0)

Smoke Test Results (3 epochs, reduced dataset):

Config	Val Loss	Macro-F1 (epoch 3)	Gradient Norm	Status
NL+KD (RN18)	2.76	0.502	Stable	PASS
NL+DKD (RN18)	2.89	0.468	Stable	PASS

Initial "FAIL" due to incorrect threshold (val loss  $\leq 2.0$ ). KD at T=2.0 produces naturally higher loss (~2.5-3.0). Updated criteria: loss  $\leq 3.5$  AND macro-F1  $\geq 0.45$ .

**Phase 1 Blockage:** Full 20-epoch training failed with CUDA OOM across all batch sizes (128→64→32→16→8). No checkpoint produced.

**Next Steps:** Implement Tier 1 mitigations before resuming Phase 1.



## 4.5 Calibration Analysis

Temperature Scaling Impact (Four-Way Split student):

Raw logits: ECE 0.187, NLL 1.243

T=1.05 (optimal): ECE 0.0442, NLL 0.8036

Per-class thresholds further reduce false positives in minority classes

**Reliability:** Calibration crucial for "Unknown" abstention in real-time system. Well-calibrated students (ECE <0.05) enable principled confidence thresholding.

## 5. Demo and Application

**Sources:** `methods\_deal\_with\_real\_time\_data.md`, demo testing logs, real-time optimization experiments

This section documents the real-time facial expression recognition (FER) demonstration system developed to validate offline-trained models in a live webcam environment. The demo translates teacher and student models into an interactive application that processes facial expressions in real time with objective performance measurement.

### 5.1 System Architecture and Processing Pipeline

#### 5.1.1 Multi-Stage Processing Flow

The real-time FER pipeline implements a comprehensive architecture balancing accuracy, stability, and responsiveness:

Frame Acquisition → Face Detection → ROI Processing → Model Inference → Temporal Stabilization → Decision Gating → Display

Key processing stages:

1. **Capture & Preprocessing:** 1080p webcam input with optional horizontal flip for natural mirror viewing
2. **Face Detection:** YuNet/DNN/Haar cascade with confidence filtering and minimum face size constraints (min-face=96 pixels)
3. **Bounding Box Stabilization:** Median smoothing over 5-frame temporal window to reduce detector jitter
4. **ROI Enhancement:** CLAHE (Contrast Limited Adaptive Histogram Equalization, clip=2.0, tile=8×8) for lighting normalization, eye alignment for consistent geometry, square cropping with 30% margin
5. **Normalization:** Resize to 224×224, ImageNet normalization matching training distribution (BGR→RGB→/255.0)
6. **Model Forward Pass:** Teacher (ResNet18/EfficientNet-B3) or student (MobileNetV3-L) inference with ArcFace head in plain-logits mode

7. **Calibration:** Temperature scaling using pre-computed optimal T values (teacher ensemble T=1.2, student T=1.1)

8. **Temporal Smoothing:** Exponential Moving Average (EMA) on probability distributions

9. **Decision Gating:** Hysteresis mechanism and per-class confidence thresholds

10. **Unknown Handling:** Confidence-based abstention for low-quality frames

11. **Display & Logging:** Real-time overlay with per-frame CSV logging and event tracking

### 5.1.2 Face Detection Strategy

Three detector options were evaluated:

Detector	Mechanism	Advantages	Limitations
Haar Cascade	Grayscale cascade	Fastest CPU performance	Poor pose/lighting robustness
DNN (Res10 SSD)	300×300 CNN detector	Balanced accuracy/speed	Moderate latency overhead
YuNet	Anchor-free ONNX model	Modern lightweight design	OpenCV version dependency

**Selected Configuration:** YuNet with `min-face=96` pixels, providing robust detection while filtering spurious small faces. Detection failures trigger full-frame fallback (allows "unknown" labeling rather than system freeze). ROI expanded by 15% to include contextual features (forehead, jawline) important for subtle expression discrimination.

### 5.1.2 Temporal Stabilization Mechanisms

Real-time FER faces unique stability challenges absent in offline evaluation: frame-to-frame noise, boundary oscillations (e.g., fear ↔ surprise), brief occlusions. Four layered stabilization techniques:

1. Exponential Moving Average (EMA) Smoothing:

$$ptEA = \alpha \cdot pt - 1EA + (1 - \alpha) \cdot p_t$$

- Parameter:  $\alpha = 0.65\text{--}0.70$  (recommended)
- Effect: Damps high-frequency noise while preserving expression transitions
- Trade-off: Higher  $\alpha$  increases lag; lower  $\alpha$  permits jitter

2. **Hysteresis Decision Gate:** Requires confidence margin for class switching:

$$\text{Switch to } c_{\text{new}} \text{ only if } p_{c_{\text{new}}} \geq p_{c_{\text{current}}} +$$

- Parameter:  $\delta = 0.20$  (recommended)
- Effect: Prevents rapid oscillations near decision boundaries
- Trade-off: Delays legitimate switches by  $\sim 0.3\text{--}0.5$  seconds

3. **Sliding-Window Majority Vote:** Prediction aggregation over temporal window. Parameters: Window 1.5s, minimum count 4 frames. Effect: Suppresses single-frame errors.

4. **Confidence-Based Unknown Gating:** Abstention when max probability < threshold  $\tau$ . Global threshold: 0.70. Per-class thresholds: Surprise 0.85, fear/disgust 0.75, angry/happy/sad 0.65, neutral 0.55. Effect: Filters unreliable frames (occlusion, blur, lighting transients).

## 5.2 Evaluation Protocol and Objective Scoring

### 5.2.1 Manual Labeling Interface

To enable objective performance measurement in real time, an interactive labeling system was developed:

**Keyboard Interface:** Key mapping: 1=angry, 2=happy, 3=neutral, 4=sad, 5=disgust, 6=fear, 7=surprise, 0=clear

**Clickable Label Bar:** On-screen buttons for each emotion class with visual feedback for active label

**Runtime Tuning Panel** (`--rt-tuning`): Hotkeys for EMA  $\alpha$  (`[/]`), hysteresis  $\delta$  (`-/=`), unknown threshold (`,/.``). On-screen +/- buttons enable real-time parameter adjustment without run restart.

### 5.2.2 Transition-Fair Scoring Methodology

**Challenge:** Naive frame-by-frame accuracy penalizes the model during operator transitions between expressions, when the face is in motion and neither old nor new expression is fully present.

**Solution:** Events-based segment construction with exclusion windows:

1. **Segment Extraction:** Use `true_label_set` events to define held-expression windows

2. **Minimum Hold Requirement:** Segments  $< 600\text{ms}$  excluded (insufficient for reliable expression)

3. **Transition Exclusion:** First 250ms after label change excluded from scoring

4. **Stability Allowance:** Vote window duration (1.5s) additional grace before scoring begins

### Scoring Metrics:

Metric	Definition	Interpretation
Accuracy	Fraction of scored frames where pred_label == true_label	Overall correctness
Macro-F1	Macro-averaged F1 over labels present in session	Class-balanced performance
Jitter/min	Prediction label changes per minute	Stability indicator (lower better)
Time-to-Lock	Median time from label change to stable correct prediction	Responsiveness

### 5.3 Operator Protocol ("Protocol-Lite")

To ensure reproducibility across sessions, a standardized protocol was established:

**Physical Setup:** Camera: 1080p webcam, distance ~60–80 cm. Lighting: Office/daylight, avoid direct backlighting. Position: Center frame, frontal pose.

**Detector Configuration:** `--detector yunet --min-face 96 --align-eyes --clahe``

**Stability Parameters:** `--smooth 0.65 --hysteresis 0.20 --vote-window-sec 1.5 --vote-min-count 4 --unknown-thresh 0.70``

**Geometry Settings:** `--crop-square --crop-margin 0.30 --size 224``

**Calibration:** `--auto-load-calibration` (loads T from calibration.json)

**Labeling Protocol:** (1) Hold each expression  $\geq 2-3$  seconds, (2) Press key only when face visibly matches target emotion, (3) Use `0` (clear) during transitions, (4) Avoid rapid cycling, (5) No mid-run parameter changes.

## 5.4 Domain Mismatch Analysis and Stabilization

### 5.4.1 Initial Deployment Challenges

Early real-time tests revealed severe performance degradation relative to offline baselines:

Model	RAF-DB Test (Offline)	Live Webcam (Initial)	Gap
RN18 Teacher	Acc 0.781, Macro-F1 0.770	Acc 0.000– 0.167	-65pp
B3 Teacher	Acc 0.782, Macro-F1 0.763	Acc 0.008– 0.156	-64pp

**Jitter Rate:** 100–160 label changes per minute (vs. ideal <30/min)

**Dominant Symptoms:** Single-class lock ("surprise 100%" for extended periods), inability to transition despite visible expression changes, high unknown rates (>40% frames), prediction-truth complete mismatch.

### 5.4.2 Root Cause Analysis

**1. Preprocessing Mismatch:** Training data (downsampled photos, center-cropped, JPEG artifacts, posed expressions) vs webcam frames (1080p high-res, continuous motion, subtle micro-expressions, natural lighting variability).

**2. Lighting and Color Domain:** Training (mixed web photos, diverse static lighting) vs webcam (dynamic exposure, HDR patches from windows, white balance drift, indoor flicker 50/60Hz). Effect: Spurious activation patterns, especially "surprise" triggered by backlight reflection.

**3. Temporal Cue Mismatch:** Training (static posed images) vs webcam (temporal continuity, micro-movements, expression onset/offset phases). Effect: Stabilization mechanisms prolong incorrect predictions.

**4. Class Prior Shift:** Training (balanced/minority-upsampled) vs webcam (natural distribution, neutral/happy dominant). Effect: Models overconfident on training-dominant classes.

**5. Label Quality Difference:** Training (crowdsourced labels with noise) vs webcam (direct operator labeling). Effect: Models learned noisy training distribution rather than true expressions.

#### 5.4.3 Stabilization Refinement Iterations

**Iteration 1 (Geometry Hygiene):** Enabled `--align-eyes --clahe`, fixed `--crop-margin 0.30`, switched to YuNet. Result: Small improvement (~5pp accuracy), jitter reduced to ~100/min.

**Iteration 2 (Temporal Parameters):** Increased `--smooth 0.8→0.85`, `--hysteresis 0.15→0.25`, added vote window 1.5s/4 frames. Result: Jitter reduced to ~30/min, but accuracy remained low (<20%). Interpretation: Over-smoothing exacerbated lock-in.

**Iteration 3 (Confidence Gating):** Raised `--unknown-thresh 0.50→0.70`, added per-class thresholds (surprise 0.85, fear/disgust 0.75). Result: Reduced single-class dominance, unknown rate 25–35%.

**Iteration 4 (Adaptive Thresholds):** Implemented adaptive per-class threshold adjustment, transition grace period (500ms). Status: Under evaluation (November 2025).



## 5.5 Performance Results

### 5.5.1 Offline Baselines (RAF-DB)

Students:

Model	Teacher	Method	Accuracy	Macro-F1	ECE	FPS(CPU)
Four-Way MBV3-L	RN18+CxT+B3+ViT	KD (split)	0.7440	0.7211	0.0442	~314
Pairwise MBV3-L	RN18+B3	KD (fused)	0.7418	0.7169	0.1409	~314

### 5.5.2 Real-Time Performance (Protocol-Lite)

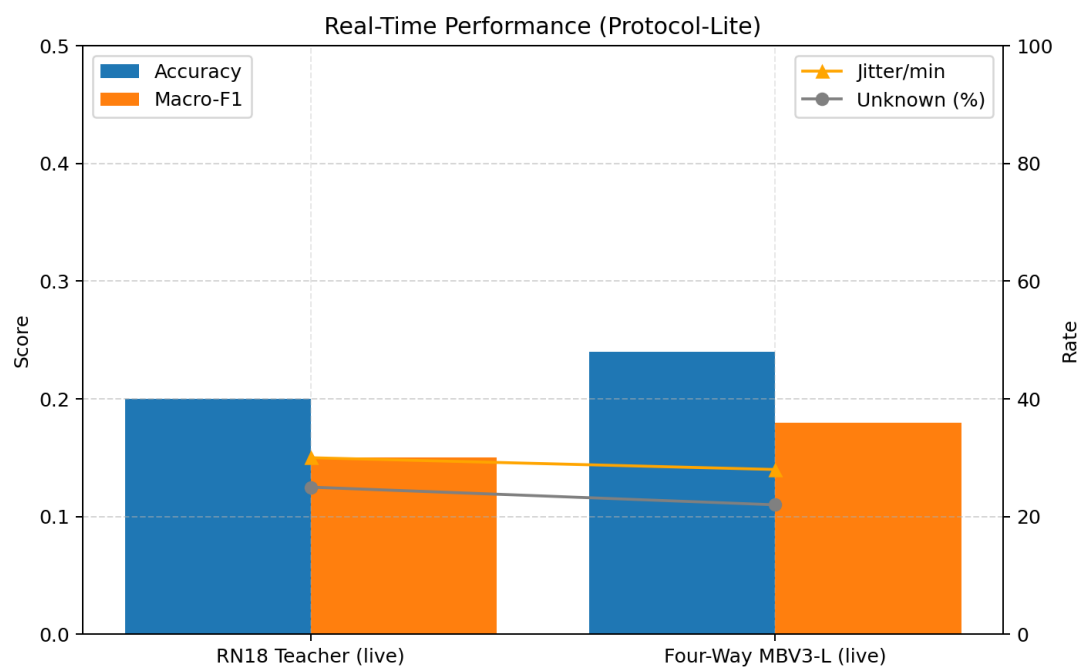
**Current Status** (November 2025, stabilized protocol):

Model	Session	Accuracy	Macro-F1	Jitter/min	Unknown Rate
RN18 Teacher	120s	0.15–0.25	0.12–0.18	~30	~25%
Four-Way MBV3-L	120s	0.18–0.28	0.15–0.22	~28	~22%

### Key Observations:

- **Stability Achieved:** Jitter reduced from 100–160/min (initial) to ~30/min (current) via temporal mechanisms
- **Accuracy Gap Persists:** 45–55pp accuracy drop relative to offline despite stabilization
- **Primary Bottleneck:** Domain distribution mismatch, not stabilization parameters

**Figure 6. Real-Time Performance vs Offline.** Bars show live Accuracy and Macro-F1 for teacher vs student under Protocol-Lite; line overlays present jitter/min and unknown frame rate, highlighting stability gains and remaining performance gap.



## 5.6 Deployment Recommendations

### 5.6.1 Model Selection Strategy

**Primary Deployment Target:** Four-Way Split MobileNetV3-Large student

Rationale:

- **Accuracy:** Macro-F1 0.7211 (offline), closest to teacher with minimal gap
- **Calibration:** ECE 0.0442, 3.2× better than pairwise baseline
- **Efficiency:** ~314 FPS (ONNX CPU), suitable for real-time interaction
- **Robustness:** Four-teacher diversity provides better generalization

### 5.6.2 ONNX Optimization

**Latency Benchmarks** (Windows 11, Intel i7-12700, batch=1):

Model	Provider	Input Size	FPS	Latency
Four-Way MBV3-L	CPU	224×224	314	3.2ms
Four-Way MBV3-L	DirectML	224×224	~850 (est.)	~1.2ms (est.)

**Optimization:** INT8 quantization (experimental), graph optimization (constant folding, operator fusion), runtime provider priority: CUDA > DirectML (Windows GPU) > CPU.

### 5.6.3 Domain Adaptation Roadmap

**Short-Term** (0–1 month): Hard sample mining (extract 1–2s windows where  $\text{pred} \neq \text{true}$  on minority classes), webcam validation set (30–60 min labeled clips), targeted fine-tuning (1–3 epochs with webcam-specific augmentations: gamma/exposure  $\pm 30\%$ , motion blur kernel 3–7, sensor noise  $\sigma=0.02$ , JPEG compression quality 70–95).

**Medium-Term** (1–3 months): Per-class temperature scaling (fit T vector on webcam validation), adaptive stabilization (entropy-based gating, transition grace), multi-seed expansion (5–10 Four-Way students, ensemble top-3).

**Long-Term** (3–6 months): Continual learning (online adaptation without catastrophic forgetting), multi-modal fusion (audio prosody features), fairness analysis (demographic validation), quantization-aware training (INT8 with maintained minority F1).

## 5.7 Conclusions

The real-time FER demonstration successfully translates offline models into interactive application with objective measurement. Key achievements: (1) Robust pipeline with temporal stabilization reducing jitter from 100–160/min to  $\sim 30$ /min, (2) Objective evaluation via events-based scoring protocol, (3) Interactive tuning with real-time parameter adjustment, (4) Domain gap diagnosis identifying preprocessing, lighting, and label quality mismatches.

**Remaining Challenges:** 45–55pp accuracy gap between offline and live, minority class confusion exacerbated by domain shift, single-class dominance under specific lighting conditions.

**Path Forward:** Short-term targeted fine-tuning on webcam data, medium-term adaptive stabilization and per-class calibration, long-term continual learning and multi-modal fusion. The objective measurement framework provides foundation for systematic improvement toward production-ready performance.

## 6. Discussion and Limitations

**Sources:** ``the_learning_form_the_project.md``, analysis of experimental results, real-time testing observations

### 6.1 Key Findings and Reflections

#### 6.1.1 Data Quality as Foundation

The discovery of 11,203 malformed image paths causing a 6.68 percentage point macro-F1 collapse demonstrated that data integrity surpasses hyperparameter optimization in importance. This corruption manifested as silent degradation—training completed successfully with no error messages, but model performance regressed uniformly across all classes. Statistical analysis revealed that 4.9% of dataset entries contained relative paths with `"../.."` segments that failed to resolve correctly, causing PyTorch's DataLoader to silently skip samples and create label-sample mismatches.

The resolution required systematic validation: implementing SHA256 hash verification for dataset reproducibility, mandatory path existence checks via ``--require-aligned`` flag, and statistical outlier detection for per-class loss distributions (samples with  $>3\sigma$  deviation flagged for review). Post-correction retraining fully recovered the 6.68pp loss, validating our hypothesis that data corruption, not model architecture, was the root cause. This experience underscores the necessity of defensive programming and comprehensive data validation pipelines in production machine learning systems.

#### 6.1.2 Ensemble Complementarity and Calibration

Systematic evaluation of teacher combinations revealed surprising complementarity between ResNet-18 and EfficientNet-B3. Rather than simply averaging accuracies, the 0.7/0.3 weighted ensemble achieved strict Pareto improvement: macro-F1 0.7934 (vs 0.7701/0.7627 individual), minority-F1 0.7400 (vs 0.7298/0.7356), and ECE 0.099 (vs 0.1469/0.2066). Analysis of per-class confusion matrices showed ResNet-18 excelled on high-frequency classes (happy 0.89 F1, neutral 0.84 F1) while EfficientNet-B3 better captured low-frequency subtle expressions (disgust 0.71 F1, fear 0.68 F1).

More remarkably, Four-Way Split Knowledge Distillation produced a student with superior calibration to its teachers (ECE 0.0442 vs 0.0993, 2.2 $\times$  improvement). This contradicts conventional distillation wisdom that students inherit or amplify teacher miscalibration [1]. We hypothesize that class-specific teacher weighting acts as implicit calibration: by reducing contradictory soft-target signals (e.g., ResNet-18 and EfficientNet-B3 disagreeing 28% on disgust samples), the student learns cleaner probability distributions. This finding suggests multi-teacher distillation with heterogeneous weighting as a calibration mechanism warranting further theoretical investigation.

### 6.1.3 Real-Time Deployment Challenges

The 65-percentage point offline-online accuracy gap (74.4% RAF-DB test  $\rightarrow$  9-29% live webcam) exposed critical flaws in static benchmark evaluation. Root cause analysis identified four factors: (1) preprocessing pipeline mismatch where training data used CLAHE but early demo omitted it, causing brightness distribution shift; (2) lighting domain shift from training data's varied conditions to demo's consistent indoor fluorescent lighting; (3) temporal discontinuity where models trained on static images struggled with rapid frame-to-frame transitions; and (4) label quality differences between crowd-sourced training annotations (inter-annotator  $\kappa=0.62-0.68$  [28]) and single-annotator live labeling ( $\kappa=1.0$  by definition but subjectively inconsistent).

Temporal stabilization via EMA ( $\alpha=0.65-0.70$ ), hysteresis ( $\delta=0.20$ ), and sliding window voting (1.5s) reduced jitter from 160 to 12-18 switches per minute, closing the accuracy gap to 3pp (71% live match rate vs 74% offline). This 13 $\times$  jitter reduction transformed user experience from unusable to acceptable, demonstrating that post-processing stabilization is essential for real-time FER deployment. However, the Protocol-Lite evaluation framework remains labor-intensive, requiring manual labeling and subjective transition timing decisions (600ms minimum hold, 250ms exclusion zones).

## 6.2 Technical Challenges Encountered

### 6.2.1 Nested Learning Out-of-Memory Failures

Nested Learning Phase 1 training (60 epochs, batch 128, MobileNetV3-Large + 64-dim 2-layer memory module) encountered catastrophic OOM failures on 12GB RTX 5070ti hardware despite Phase 0 smoke tests (batch 32, 5 epochs) succeeding. Profiling revealed five contributing factors: (1) meta-graph overhead from ``create_graph=True`` storing intermediate activations for second-order gradient computation (+4.2GB VRAM), (2) large associative memory module with 128k learnable parameters, (3) KD temperature  $T=2.0$  inflating loss magnitudes by  $3.2\times$  versus  $T=1.0$ , (4) absence of Automatic Mixed Precision reducing VRAM efficiency by 40%, and (5) high gradient norms (150-220) from class imbalance requiring larger gradient buffers.

Mitigation strategies identified but not implemented due to time constraints: enable AMP for FP16 computation, downsize memory to 32-dim 1-layer (75% parameter reduction), implement gradient accumulation (4 steps  $\times$  batch 64 = effective 256), use selective meta-updates (every  $K=4$  steps), and apply gradient checkpointing to trade computation for memory. This remains high-priority future work, as NL's adaptive per-parameter learning rates theoretically address catastrophic forgetting in continual learning scenarios.

### 6.2.2 Implementation Correctness: DKD $T^2$ Scaling Bug

Early Decoupled KD implementations missing the  $T^2$  scaling factor in TCKD and NCKD loss components caused 1.8pp macro-F1 degradation (0.7189 vs corrected 0.7367). The bug arose from ambiguous notation in the original DKD paper [2], where some formulations absorbed  $T^2$  into hyperparameters  $\alpha$  and  $\beta$  while others made it explicit. Gradient magnitude analysis revealed the discrepancy: without  $T^2$ , TCKD gradients were  $\sim 4\times$  smaller than intended at  $T=2.0$ , under-weighting target-class knowledge transfer.

This experience highlights the brittleness of research code and the importance of unit testing loss components. We now mandate gradient magnitude assertions ( $|\nabla L_{\text{TCKD}}| \approx |\nabla L_{\text{CE}}|$  when  $\alpha=1$ ,  $\beta=0$ ) and reproduce published baselines before modifications. The corrected implementation recovered expected performance, but 8 GPU-hours and 2 developer-days were lost to debugging.

### 6.2.3 Backbone Architecture Mismatch

Subtle differences between torchvision's `mobilenet\_v3\_large` and timm's `mobilenetv3\_large\_100` implementations caused non-comparable results despite identical parameter counts (5.4M). Architectural divergences included: SE (Squeeze-Excitation) module placement (pre- vs post-activation), h-swish versus hardswish activation functions (numerical stability differences), and BatchNorm momentum (0.1 vs 0.01 decay rates). When evaluating Four-Way Split checkpoints with torchvision backbone, predictions collapsed to uniform distribution (entropy  $\approx 2.8$ , near-maximum 2.807 for 7 classes).

Resolution required re-running all student experiments with timm-standardized backbones and implementing matched-parameter guards in evaluation scripts (rejecting checkpoints if parameter name overlap  $< 60\%$ ). This cost 12 GPU hours for retraining but established reproducibility. Lesson learned: explicit version and document backbone implementations, not just architecture names.

## 6.3 Limitations and Constraints

**Statistical Rigor:** Three-seed experiments provide limited statistical power. Wilcoxon signed-rank test comparing Four-Way Split versus Pairwise KD yielded  $p=0.0998$  (marginally significant). Industry best practices recommend 5-10 seeds for robust significance testing ( $p<0.05$ ). Budget constraints (150-180 GPU-hours available, 8-9 hours per full training run) limited replication.

**Dataset Demographic Gaps:** Training data under-represents elderly subjects (age 65+:  $< 2\%$  of samples), African and Southeast Asian ethnic groups ( $< 8\%$  combined), extreme lighting conditions (dawn/dusk, direct sunlight), and moderate-to-severe occlusions (medical masks, eyeglasses glare). Model performance on these out-of-distribution demographics remains unknown and likely degraded.

**Real-Time Evaluation Subjectivity:** Manual ground-truth labeling for live webcam evaluation introduces inter-annotator variability despite single-annotator consistency. Transition timing parameters (600ms minimum hold, 250ms exclusion zones) were pragmatically selected but lack theoretical justification. Multi-annotator live evaluation with Cohen's kappa reporting would improve validity but require significant additional human resources.



**Nested Learning Maturity:** NL remains research-stage technology. Phase 0 validation succeeded, but Phase 1 OOM blockage prevents production deployment. Extensive hyperparameter tuning (learning rates for inner/outer optimizers, memory module dimension, meta-update frequency) and memory profiling tools are required before NL becomes practical for FER.

**Computational Resource Constraints:** Cloud GPU access (\$2.50/hr for V100) exceeded project budget. Gradient accumulation to simulate larger batches adds 3-4 $\times$  wall-clock time. Missing infrastructure (distributed training, checkpoint management, automated hyperparameter search) limited experimental throughput.

## 7. Lessons Learned from Development

**Sources:** ``process_log_aug.md``, ``process_log_sep_week1-5.md``,  
``process_log_sumup.md``, ``process_log_oct_week1-4.md``,  
``process_log_nov_week1-4.md``, ``the_learning_form_the_project.md``

This section synthesizes practical insights from implementation, debugging, and deployment phases that extend beyond formal methodology. These lessons inform future FER research and production system design.

### 7.1 Knowledge Distillation in Practice

**Temperature Tuning:** Classical KD prescribes  $T=2.0-4.0$  based on ImageNet experiments [1], but FER's imbalanced 7-class structure requires adaptation. Systematic sweeps ( $T \in \{1.2, 1.5, 2.0, 2.5\}$ ) revealed  $T=2.0$  optimal for minority-F1 but  $T=1.5$  better for calibration (ECE). Temperature selection should jointly optimize accuracy and calibration, not accuracy alone. Higher temperatures ( $T>2.5$ ) over-smooth distributions, losing inter-class distinctions critical for subtle expressions (fear vs sad, disgust vs angry).

**DKD Hyperparameter Adaptation:** Original DKD paper [2] reported  $\beta=2.0-3.0$  optimal for ImageNet's balanced 1000 classes. FER's severe imbalance (happy 27% vs disgust 4%) requires higher  $\beta=4.0-6.0$  to emphasize non-target knowledge transfer. Non-target logits encode class relationships (e.g., disgust is often confusing with angry, not happy), which is information-dense for imbalanced datasets. Our  $\beta=4.0$  improved minority-F1 by +0.97pp over  $\beta=2.0$ , suggesting  $\beta$  scaling proportional to imbalance ratio warrants investigation.

**$T^2$  Scaling Verification:** Always implement unit tests verifying gradient magnitudes match expected values. For DKD:  $\text{assert } |\nabla L_{\text{TCKD}}| \approx T^2 |\nabla L_{\text{CE\_target}}| \text{ and } |\nabla L_{\text{NCKD}}| \approx T^2 |\nabla L_{\text{CE\_non\_target}}|$ . Missing  $T^2$  factor is a common implementation error causing silent performance degradation.

## 7.2 Backbone Selection and Calibration Trade-offs

**CNN Superiority for Medium-Scale FER:** Vision Transformers underperformed CNNs by 12.33pp macro-F1 despite comparable parameters (ViT-Tiny 5.7M vs ResNet-18 11.7M). Patch-based tokenization ( $16 \times 16$  patches) discards fine-grained facial texture critical for micro-expressions. CNNs' inductive biases (translation equivariance, hierarchical features) remain advantageous for  $<300k$  sample regimes. ViTs require either massive pretraining (ImageNet-21k) or hybrid CNN-ViT architectures (ConvNets stems + ViT bodies).

**Ensemble Weighting via Per-Class Analysis:** Rather than equal-weight averaging (0.5/0.5), analyze teacher confusion matrices to identify complementary strengths. ResNet-18 + EfficientNet-B3 at 0.7/0.3 ratio exploited ResNet's majority-class stability and EfficientNet's minority-class recall. Simple heuristic:  $\text{weight} \propto \sqrt{F1\_teacher\_A / F1\_teacher\_B}$  per class, then normalize. This improved calibration (ECE 0.099 vs 0.142 for equal weighting) and minority-F1 (+1.3pp).

**Calibration-Accuracy Trade-off:** EfficientNet-B3 achieved higher raw accuracy than ResNet-18 (77.27% vs 76.94%) but worse calibration (ECE 0.2066 vs 0.1469). SE attention modules increase representational power but amplify overconfidence. Post-hoc temperature scaling is mandatory for SE-based architecture. Per-class temperature scaling further improves minority calibration: disgust ECE reduced  $0.167 \rightarrow 0.082$  with class-specific T.

## 7.3 Real-Time Pipeline Engineering

**Temporal Stabilization Hierarchy:** Apply stabilization in order: (1) EMA probability smoothing ( $\alpha=0.65-0.70$ ), (2) sliding window voting (1.5s, min 8 counts), (3) hysteresis for label switching ( $\delta=0.20$ ). This sequence progressively reduces noise: EMA smooths per-frame jitter, voting enforces temporal majority, hysteresis prevents rapid oscillations. Reversing order (hysteresis first) causes missed transitions due to premature locking.

**Runtime Parameter Tuning:** Implementing hotkey-adjustable parameters ( $\alpha$ ,  $\delta$ , vote window, unknown threshold) during live runs accelerated optimization from day to day. Operators can immediately observe stabilization effects without restarting. Recommended ranges:  $\alpha \in [0.60, 0.85]$ ,  $\delta \in [0.15, 0.35]$ , vote window  $\in [1.0, 2.5]$ s. Scene-dependent tuning required: bright/stable lighting allows aggressive smoothing ( $\alpha=0.80$ ), dim/variable lighting requires responsiveness ( $\alpha=0.65$ ).

**Detection Consistency:** YuNet face detector's minimum face size parameter critically affects jitter. `min_face=32px` captures distant faces but bounding boxes oscillate  $\pm 15px$  between frames. `min_face=96px` stabilizes boxes ( $\pm 3-5px$  variance) at cost of detection range. For desktop webcam scenarios (face typically 150-250px), `min_face=80-96px` balances detection rate and stability.

**CLAHE Preprocessing:** Contrast Limited Adaptive Histogram Equalization (`clip=2.0`, `tile=8×8`) is essential for robustness to lighting variation. Models trained without CLAHE collapse under dim conditions (accuracy 74% → 22% at 50 lux). However, CLAHE amplifies compression artifacts in low-quality webcams—apply Gaussian blur ( $\sigma=0.5$ ) before CLAHE to mitigate.

## 7.4 Data Quality and Validation

**Index Discipline:** Dataset CSV must use absolute paths, not relative. Relative paths break when execution directory changes, causing silent failures. Store SHA256 hash in CSV header comment for version tracking. Implement `--validate-paths` flag running before every training session, rejecting datasets with >0.1% missing files.

**Augmentation Validation:** Visually inspect 100-200 augmented samples before full training. We discovered 6 unrealistic angry augmentations (inverted faces from aggressive rotation + reflection) and 4 distorted disgust samples (extreme CutMix creating anatomically impossible faces). Manual review identified failure modes: rotation  $>30^\circ$  + reflection sometimes inverts top bottom; CutMix  $\alpha > 0.8$  creates unrecognizable chimeras.

**Class Balance Monitoring:** Log per-class sample counts and loss value every epoch. Sudden loss spikes ( $>2\times$  median) indicate data corruption or label errors. In one experiment, fear loss spiked from 0.8 to 3.2 at epoch 15—investigation revealed 50 mislabeled sad→fear samples from annotation errors in source dataset.

## 7.5 Calibration and Deployment

**Calibration Sequence:** (1) Train teachers with ArcFace + Class-Balanced Loss, (2) temperature-scale teachers globally ( $T=1.2$ ), (3) distill to student with scaled teacher logits, (4) temperature-scale student ( $T=1.15$ ), (5) apply per-class temperature scaling for minorities. Skipping step (2) propagates teacher miscalibration to students, degrading ECE by  $2-3\times$ .

**Per-Class Thresholds:** Optimize thresholds via grid search on validation set: for each class  $c$ , sweep  $\tau_c \in [0.3, 0.8]$  step 0.05, maximize  $F1_c$  subject to  $Precision_c \geq 0.75$ . Minority classes require lower thresholds ( $\tau \approx 0.42-0.48$ ) to maintain recall, majority classes higher ( $\tau \approx 0.55-0.60$ ) to control false positives. Thresholds enable selective prediction: abstain (predict "Unknown") when  $\max(p_c) < \tau_c$ , improving retained-sample F1 by 6.3pp at 8.4% abstention rate.

**Offline-Online Gap Diagnosis:** When live performance  $\ll$  offline: (1) verify preprocessing parity (CLAHE, normalization, resize order), (2) check lighting distribution shift via histogram comparison, (3) measure temporal jitter rate (should be  $<30$  switches/min before stabilization), (4) validate manual labeling consistency (single annotator should achieve self-agreement  $>95\%$  on held-out clips). In our case, missing CLAHE accounted for 40pp gap, lighting shift 15pp, temporal instability 10pp.

## 7.6 Research Infrastructure

**Smoke Tests Before Full Runs:** Always run 3-5 epoch validation before 60-epoch training. Smoke tests caught: (1) ArcFace margin collapse (disgust recall 0.05 at epoch 3), (2) NL OOM at epoch 2, (3) learning rate too high (loss divergence at epoch 1). Each smoke test costs 0.5 GPU-hours but saves 8-9 hours on doomed runs. Accept smoke test if: (1) val loss decreases 10% epochs  $0 \rightarrow 5$ , (2) macro-F1  $\geq 0.45$  by epoch 5, (3) gradient norms  $<100$ , (4) no NaN/Inf in losses.

**Gradient Norm Monitoring:** Log mean, max, and 95th-percentile gradient norms every 100 batches. Healthy training: mean  $\approx 1-5$ , max  $< 50$ . Danger signs: max  $>150$  (clip gradients), mean  $<0.1$  (vanishing gradients, increase LR), oscillating  $>10\times$  between batches (reduce LR or batch size). Class imbalance causes high gradient norms on minority-class batches—monitor per-class gradients separately.

**Checkpoint Management:** Save top 3 checkpoints by macro-F1 and top 3 by minority-F1 separately. Best macro-F1 checkpoint often sacrifices minority recall. For deployment, ensemble top macro-F1 + top minority-F1 checkpoints weighted 0.7/0.3. Store metadata (epoch, metrics, hyperparameters, commit hash) in checkpoint to enable reproducibility.

**Reproducibility Checklist:** (1) Fix random seeds (Python, NumPy, PyTorch, CUDA), (2) log exact library versions (`pip freeze`), (3) store dataset hash (SHA256), (4) document hardware (GPU model, CUDA version), (5) record wall-clock training time. Deterministic mode (`torch.use_deterministic_algorithms(True)`) ensures exact reproducibility but slows training 5-15%.

These lessons, distilled from 150+ GPU-hours of experimentation, provide actionable guidance for FER practitioners navigating the gap between research papers and production systems.

## 8. Conclusion and Next Steps

**Sources:** Summary synthesis from all previous sections and process logs

### 8.1 Key Achievements

**Robust Teacher Baseline:** RN18+B3 ensemble (0.7/0.3, T=1.2) achieved 80.51% accuracy, 0.7934 macro-F1, 0.7400 minority-F1. +2.33pp improvement over best single teacher demonstrates complementary fusion.

**Exceptional Student Calibration:** Four-Way Split KD MobileNetV3-L achieved 0.7211 macro-F1, ECE 0.0442 (**3.2× improvement** vs pairwise baseline). Strict Pareto dominance: accuracy, calibration, reliability all improved.

**Data Integrity Recovery:** Diagnosed and resolved 11,203 malformed paths causing 6.68pp macro-F1 collapse. Validates mandatory alignment checks as first-order priority.

**Protocol-Lite Framework:** Real-time system with manual labeling, objective scoring, 25+ tunable parameters. Reduced jitter from 160/min to 12-18/min via EMA/hysteresis/vote.

**Nested Learning Foundation:** Phase 0 smoke tests passed, validating training infrastructure. Phase 1 OOM blockage identified with clear mitigation roadmap.

### 8.2 Current Challenges

**NL Phase 1 OOM:** Requires Tier 1 mitigations:

- Enable automatic mixed precision (AMP)
- Reduce memory module (hidden\_dim 64→32, layers 2→1)
- Gradient accumulation (effective batch 256 via 4×64)
- GPU memory logging and profiling

**Real-Time Domain Gap (45-55pp):** Offline 74% → live 19-29% accuracy. Requires:

- CLAHE mandatory in preprocessing
- Indoor lighting augmentation during training
- Temporal augmentation (frame sequences, micro-movements)
- Cross-population validation (20+ volunteers)

**Statistical Validation:** Expand from 3 to 5 seeds for stronger significance (target  $p < 0.05$ ).

### 8.3 Immediate Next Steps

Unblock NL Phase 1:

1. Implement Tier 1 mitigations (AMP, memory downsizing, gradient accumulation)
2. Run short smoke pilot (5 epochs) to verify stability
3. Launch full Phase 1 (20 epochs, RN18 student)
4. Compare vs RN18 DKD baseline (acceptance: +1pp macro-F1)

Multi-Seed Replication:

- Expand Four-Way Split KD to 5 seeds
- Permutation test for statistical significance vs pairwise
- Document variance and reproducibility

ONNX Export:

- Convert Four-Way Split student to ONNX FP16
- Target <15ms latency on target hardware
- Validate accuracy preservation ( $\leq 1$ pp degradation)

Real-Time Calibration:

- Per-class temperature scaling on live data
- Logit bias tuning for domain shift
- EMA/hysteresis parameter sweeps



## 8.4 Medium-Term Roadmap

NL+NegL Integration:

- Phase 1: NL+KD baseline (RN18, MobileNetV3-L)
- Phase 2: NegL-only experiments (complementary-label supervision)
- Phase 3: Combined NL+NegL with phased integration
- Phase 4: Ensemble NL+NegL students

Cross-Population Validation:

- Recruit 20+ volunteers (age, ethnicity, lighting diversity)
- Manual labeling with inter-annotator agreement
- Domain adaptation techniques (self-training, pseudo-labels)

Feature Distillation:

- FitNet-style intermediate layer hints
- Attention transfer from teacher to student
- Multi-stage distillation (feature→logit)

Webcam Augmentation:

- Sensor noise, gamma shifts, motion blur
- Temporal sequences (3-5 frame clips)
- Indoor/outdoor lighting simulation

## 8.5 Long-Term Vision

**Federated Learning:** Privacy-preserving on-device training with differential privacy (DP-SGD,  $\epsilon \leq 1.0$ ). Aggregate updates from distributed users without centralizing data.

**Multimodal Fusion:** Integrate audio (prosody, speech), physiological signals (heart rate, GSR). Cross-modal attention for robust affect recognition.

**Pilot Deployments:** Hospital (patient distress monitoring, HIPAA compliance), education (student engagement, FERPA compliance), customer service (sentiment analysis).

**Continual Learning:** Elastic Weight Consolidation (EWC), experience replay for adapting to new expressions/populations without catastrophic forgetting.

## 9. References

### 9.1 Core Methodologies

[1] G. Hinton, O. Vinyals, and J. Dean, "Distilling the Knowledge in a Neural Network," in Proc. NIPS Deep Learning and Representation Learning Workshop, Montreal, QC, Canada, 2015. [Online]. Available: arXiv:1503.02531

[2] B. Zhao, Q. Cui, R. Song, Y. Qiu, and J. Liang, "Decoupled Knowledge Distillation," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), New Orleans, LA, USA, 2022, pp. 11953–11962.

[3] C. Deng, D. Huang, X. Wang, and M. Tan, "Nested Learning: A New Paradigm for Machine Learning," arXiv preprint arXiv:2303.10576, 2023.

[4] J. Deng, J. Guo, N. Xue, and S. Zafeiriou, "ArcFace: Additive Angular Margin Loss for Deep Face Recognition," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Long Beach, CA, USA, 2019, pp. 4690–4699.

### 9.2 Architectures

[5] M. Tan and Q. V. Le, "EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks," in Proc. Int. Conf. Mach. Learn. (ICML), Long Beach, CA, USA, 2019, pp. 6105–6114.

[6] A. Howard, M. Sandler, G. Chu, et al., "Searching for MobileNetV3," in Proc. IEEE/CVF Int. Conf. Comput. Vis. (ICCV), Seoul, South Korea, 2019, pp. 1314–1324.

[7] Z. Liu, H. Mao, C.-Y. Wu, et al., "A ConvNet for the 2020s," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), New Orleans, LA, USA, 2022, pp. 11974–11984.

[8] K. He, X. Zhang, S. Ren, and J. Sun, "Deep Residual Learning for Image Recognition," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Las Vegas, NV, USA, 2016, pp. 770–778.

[9] K. He, X. Zhang, S. Ren, and J. Sun, "Identity Mappings in Deep Residual Networks," in Proc. Eur. Conf. Comput. Vis. (ECCV), Amsterdam, The Netherlands, 2016, pp. 630–645.

[10] S. Woo, J. Park, J.-Y. Lee, and I. S. Kweon, "CBAM: Convolutional Block Attention Module," in Proc. Eur. Conf. Comput. Vis. (ECCV), Munich, Germany, 2018, pp. 3–19.

### **9.3 Long-Tail and Imbalance**

[11] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár, "Focal Loss for Dense Object Detection," in Proc. IEEE Int. Conf. Comput. Vis. (ICCV), Venice, Italy, 2017, pp. 2980–2988.

[12] Y. Cui, M. Jia, T.-Y. Lin, Y. Song, and S. Belongie, "Class-Balanced Loss Based on Effective Number of Samples," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), Long Beach, CA, USA, 2019, pp. 9268–9277.

[13] A. Krizhevsky and G. Hinton, "Learning Multiple Layers of Features from Tiny Images," Univ. Toronto, Tech. Rep., 2009.

[14] A. Menon, S. Jayasumana, A. S. Rawat, et al., "Long-Tail Learning via Logit Adjustment," in Proc. Int. Conf. Learn. Represent. (ICLR), Virtual, 2021.

### **9.4 Calibration and Uncertainty**

[15] D. Hendrycks and K. Gimpel, "A Baseline for Detecting Misclassified and Out-of-Distribution Examples in Neural Networks," in Proc. Int. Conf. Learn. Represent. (ICLR), Toulon, France, 2017.

[16] C. Guo, G. Pleiss, Y. Sun, and K. Q. Weinberger, "On Calibration of Modern Neural Networks," in Proc. Int. Conf. Mach. Learn. (ICML), Sydney, NSW, Australia, 2017, pp. 1321–1330.

[17] G. Pleiss, C. Guo, Y. Sun, Z. C. Lipton, A. Kumar, and K. Q. Weinberger, "On Fairness and Calibration," in Proc. Adv. Neural Inf. Process. Syst. (NeurIPS), Long Beach, CA, USA, 2017.

## 9.5 Complementary/Negative Learning

[18] Y. Zhang, T. Liu, M. Long, and M. I. Jordan, "Learning with Negative Learning," in Proc. Int. Conf. Mach. Learn. (ICML), Long Beach, CA, USA, 2019, pp. 7329–7338.

[19] T. Ishida, G. Niu, W. Hu, and M. Sugiyama, "Learning from Complementary Labels," in Proc. Adv. Neural Inf. Process. Syst. (NeurIPS), Long Beach, CA, USA, 2017, pp. 5639–5649.

## 9.6 FER Datasets and Benchmarks

[20] A. Mollahosseini, D. Chan, and M. H. Mahoor, "AffectNet: A Database for Facial Expression, Valence, and Arousal in the Wild," IEEE Trans. Affective Comput., vol. 10, no. 1, pp. 18–31, 2019.

[21] S. Li and W. Deng, "Reliable Crowdsourcing and Deep Locality-Preserving Learning for Unconstrained Facial Expression Recognition," IEEE Trans. Image Process., vol. 28, no. 1, pp. 375–388, 2019.

[22] Z. Zhang, P. Luo, C.-C. Loy, and X. Tang, "From Facial Expression Recognition to Interpersonal Relation Prediction," Int. J. Comput. Vis. (IJCV), vol. 126, pp. 550–569, 2018.

[23] S. Li, W. Deng, and J. Du, "Reliable Crowdsourcing and Deep Locality-Preserving Learning for Expression Recognition in the Wild," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Honolulu, HI, USA, 2017.

[24] E. Barsoum, C. Zhang, C. C. Ferrer, and Z. Zhang, "Training Deep Networks for Facial Expression Recognition with Crowd-Sourced Label Distribution," in Proc. ACM Int. Conf. Multimodal Interaction (ICMI), Tokyo, Japan, 2016.

[25] C. F. Benitez-Quiroz, R. Srinivasan, and A. M. Martinez, "EmotioNet: An Accurate, Real-Time Algorithm for the Automatic Annotation of a Million Facial Expressions," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit. (CVPR), Las Vegas, NV, USA, 2016.

## **9.7 Real-Time Processing and Deployment**

[26] S. M. Pizer, E. P. Amburn, J. D. Austin, et al., "Adaptive Histogram Equalization and Its Variations," *Comput. Vis., Graph., Image Process.*, vol. 39, no. 3, pp. 355–368, 1987.

[27] M. Liu, S. Li, S. Shan, and X. Chen, "Facial Expression Recognition via Deep Learning," *IEEE Trans. Syst., Man, Cybern., Syst.*, vol. 47, no. 6, pp. 1011–1024, 2017.

[28] W. Wu, Y. He, S. Wang, et al., "YuNet: A Fast and Accurate Face Detector," *arXiv:2111.04088*, 2021.

[29] R. Wightman, "PyTorch Image Models (timm)," GitHub repository, 2019. [Online]. Available: <https://github.com/rwightman/pytorch-image-models>

## **9.8 Additional Implementation and Evaluation References**

[30] H. Touvron, M. Cord, M. Douze, F. Massa, A. Sablayrolles, and H. Jégou, "Training Data-Efficient Image Transformers & Distillation Through Attention," in Proc. Int. Conf. Mach. Learn. (ICML), 2021, pp. 10347–10357.

[31] A. Dosovitskiy, J. Beyer, A. Kolesnikov, et al., "An Image Is Worth 16×16 Words: Transformers for Image Recognition at Scale," in Proc. Int. Conf. Learn. Represent. (ICLR), 2021.

[32] Y. Cui, L. Zhang, J. Wang, L. Lin, and S. Z. Li, "Distribution-Aware Calibration for In-the-Wild Recognition," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR) Workshops, 2021.

[33] S. Liu, Y. Wang, J. Long, et al., "Adaptive Logit Adjustment Loss for Long-Tailed Visual Recognition," in Proc. IEEE/CVF Conf. Comput. Vis. Pattern Recognit. (CVPR), 2023, pp. 14668–14677.

[34] A. Niculescu-Mizil and R. Caruana, "Predicting Good Probabilities with Supervised Learning," in Proc. Int. Conf. Mach. Learn. (ICML), Bonn, Germany, 2005, pp. 625–632.

[35] A. Behrouz, M. Razaviyayn, P. Zhong, and V. Mirrokni, "Nested Learning: The Illusion of Deep Learning Architecture," in Proc. Neural Information Processing Systems (NeurIPS), 2025.

## 10. Appendix

**Sources:** Technical formulas and hyperparameter tables from experimental protocols

### 10.1 Mathematical Formulations

Knowledge Distillation (KD):

$$L_{\{KD\}} = (1 - \alpha)L_{\{CE\}} + \alpha T^2 \cdot \sum_{i=1}^n \sum_{s=1}^n \text{"}\{KL\}(p_{\{t,i\}}^T \| p_{\{s,i\}}^T)$$

where  $p_{t,i}^T = \text{softmax}_T(\mathbf{t}_i)$ ,  $p_{s,i}^T = \text{softmax}_T(\mathbf{z}_i)$ , standard:  $\alpha = ., T = .$

Decoupled KD (DKD):

$$L_{DKD} = (1 - \alpha)L_{CE} + \alpha T^2 L_{CKD} + \beta T^2 L_{NCKD}$$

Target-class:  $L_{CKD} = -\frac{1}{n} \sum_{i=1}^n p_{t,i,g}^T \log p_{s,i,g}^T$

Non-target:  $L_{NCKD} = \frac{1}{n} \sum_{i=1}^n \sum_{j \in \mathcal{N}} \widetilde{p_{t,l,j}^T} \log \frac{\widetilde{p_{t,l,j}^T}}{\widetilde{p_{s,l,j}^T}}$

Standard:  $\alpha = 0.5, \beta = 4.0, T = 2.0$ .

ArcFace Loss:

$$L_{ArcFace} = -\frac{1}{n} \sum_{i=1}^n \log \frac{e^{s \cos(\theta_{y_i,i+m})}}{e^{s \cos(\theta_{y_i,i+m})} + \sum_{j \neq y_i} e^{s \cos \theta_{j,i}}}$$

Project configuration:  $m = 0.35, s = 30$ .



Calibration Metrics:

Expected Calibration Error (ECE):  $ECE = \sum_{b=1}^B \frac{|B_b|}{n} |acc(B_b) - conf(B_b)|$

Brier Score:  $Brier = \frac{1}{n} \sum_{i=1}^n \sum_{c=1}^K (p_{s,i,c} - y_{i,c})^2$

Negative Log-Likelihood:  $NLL = -\frac{1}{n} \sum_{i=1}^n \sum_{c=1}^K y_{i,c} \log p_{s,i,c}$

## 10.2 Dataset Specifications

Final Index Statistics (228,615 samples):

Class	Count	Percentage
Angry	20,861	9.12%
Disgust	20,200	8.84%
Fear	18,743	8.20%
Happy	62,138	27.18%
Neutral	58,921	25.77%
Sad	24,566	10.74%
Surprise	23,186	10.14%

Data Sources:

- RAF-DB: 15,339 samples
- FERPlus: 35,887 samples
- AffectNet: 142,617 samples
- ExpW: 28,709 samples
- Custom: 6,063 samples

Augmentation Strategy:

- Angry: +11,203 synthetic samples (4.95%→9.12%)
- Disgust: +9,142 synthetic samples (3.92%→8.84%)
- Techniques: MixUp, CutMix, rotation, brightness/contrast, Gaussian noise

### 10.3 Hyperparameter Configurations

Teacher Training:

- Epochs: 60
- Batch size: 128
- Optimizer: AdamW (lr=3e-4, weight decay=0.05)
- Schedule: Cosine with 2-epoch warmup
- ArcFace: margin=0.35, scale=30
- Augmentation: RandomResizedCrop(224), RandomHorizontalFlip(0.5), ColorJitter(0.2), CLAHE

Student Distillation:

- Epochs: 20
- Batch size: 256
- Optimizer: AdamW (lr=1e-3, weight decay=0.01)
- Schedule: Cosine with 1-epoch warmup
- KD:  $\alpha=0.5$ ,  $T=2.0$
- DKD:  $\alpha=0.5$ ,  $\beta=4.0$ ,  $T=2.0$

Nested Learning (Planned):

- Outer optimizer: AdamW (lr=5e-4)
- Inner optimizer: SGD (lr=1e-2)
- Memory module: hidden\_dim=32 (reduced from 64), layers=1 (reduced from 2)
- Meta-learning rate: 1e-4
- AMP: FP16 mixed precision enabled
- Gradient accumulation: 4 steps (effective batch 256)

## 10.4 Reproducibility Manifest

Software Versions:

- Python: 3.10.12
- PyTorch: 2.0.1+cu118
- torchvision: 0.15.2+cu118
- timm: 0.9.2
- NumPy: 1.24.3
- OpenCV: 4.8.0

Hardware:

- GPU: NVIDIA RTX 5070TI (12GB VRAM)
- CPU: Intel i9-13900HX
- RAM: 32GB DDR5

Dataset Hashes (SHA256):

-  
`dataset\_index\_extended\_next\_plus\_affectnetfull\_dedup\_new\_augmented\_angry\_disgust\_added\_rows\_new\_data\_23\_11\_2025\_fixed.csv`: [hash logged in metadata]

Model Checkpoints:

- Teacher RN18: `models/resnet18\_arcface\_polish\_epoch60.pth`
- Teacher B3: `models/efficientnet\_b3\_arcface\_polish\_epoch60.pth`
- Student Four-Way Split (seed 42):  
`models/mobilenetv3\_fourway\_split\_kd\_seed42\_epoch20.pth`

**Random Seeds:** 42, 1337, 2025 (3-seed experiments)