# Core Group vs Control Group Student Training Study Report

Generated from Markdown source. Mathematical formulas are provided in LaTeX inline; for high-fidelity rendering, embed the LaTeX snippet (math_equations_snippet.tex) in the paper workflow.

## Core Group vs Control Group Student Training Study Report

(香港時區 / Hong Kong Time UTC+8)

## 1. Objectives / 研究目標

- Establish reproducible, aligned MobileNetV3 student baselines using canonical index `dataset_index_extended_next_plus_affectnetfull_dedup.csv`.
- Compare Core Group (Pairwise CNN + Hybrid) vs Control Group (Split CNN+ViT, Four-Way Split (6E), and Fused Four Equal (6C)) to identify the highest macro F1 and explain performance differentials.
- Diagnose and document root cause of earlier degraded Split Multi-Teacher results (label length/order misalignment) and impact of enforcing strict alignment (`--require-aligned`).
- Provide decision guidance for deployment candidate selection (accuracy vs complexity vs calibration extensibility).

## 2. Experimental Families 定義

| Family ID | Group | Description | Teachers / Composition | Distill Mode(s) | Notes |
|-----------|-------|------------|-----------------------|----------------|-------|
| 6A | Core | Pairwise CNN Weighted | ResNet18_polish (0.7) + EffNet-B3 (0.3) fused | KD, DKD (β=2,4,8 historic) | Stable strong baseline |
| 6B | Core | Hybrid (Pairwise + ViT) | (Pairwise 0.7/0.3) + ViT BAN (0.3) | KD, DKD | Adds transformer context |

| 6D-split | Control | Split CNN+ViT (λ sweep) | Pairwise vs ViT separate, λ ∈ {0.5/0.5,0.7/0.3,0.3/0.7} | KD, DKD (β=4) | Multi-source aligned inputs |

| 6E-split4 | Control | Four-Way Split Equal | RN18 + ConvNeXt-Tiny + EffB3 + ViT (0.25 each) | KD, DKD (β=4) | Richest diversity |

| 6C | Control | Fused Four Equal (Planned) | Single fused softlabels (equal-weight average of RN18 + ConvNeXt-Tiny + EffB3 + ViT) | KD, DKD | Planned; over-smoothing hypothesis test (not yet aligned) |

All experiments: epochs=20, batch=128, seeds={1337,2025}, KD $\alpha$=0.5 T=2, DKD $\alpha$=0.5 β as specified, `--num-workers 0` (Windows determinism), test eval every 5 epochs.

## 3. Data & Alignment Integrity / 數據與對齊

- Canonical index enforced across all teacher softlabel exports (`*_ixNextAffFull`).
- Alignment diagnostic (`diagnose_softlabel_alignment.py --strict`) confirmed equality of lengths & hashes pre-training.
- Earlier misaligned split runs (pre-canonical) purged to remove contaminated comparisons.

## 4. Metrics & Sources

Primary metrics from `metrics_final.json` (best_epoch, best_test_macro_f1, best_test_acc, per-class F1). Time & checkpoint inferred if not explicit. Timestamps localized to UTC+8 via `--tz-offset 8` in summary collector.

### 4.1 Mathematical Definitions / 數學公式定義

We formalize the evaluation metrics, distillation objectives, and calibration quantities referenced throughout the report.

#### 4.1.1 Notation

- Number of classes: $K$ (here $K=7$: angry, disgust, fear, happy, neutral, sad, surprise)
- Sample index: $i = 1,\dots,n$
- Ground-truth one-hot label: $\mathbf{y}_i \in \{0,1\}^K$
- Student logits: $\mathbf{z}_i \in \mathbb{R}^K$; Teacher logits (for teacher $m$): $\mathbf{t}_{i}^{(m)}$
- Softmax with temperature $T$: $\operatorname{softmax}_T(\mathbf{z})_c = \dfrac{\exp(z_c / T)}{\sum_{j} \exp(z_j / T)}$
- Teacher set size (multi-teacher): $M$

中文: 記號說明。

### 4.1.2 Per-Class Precision / Recall / F1

中文: 分類別精確率 (Precision) / 召回率 (Recall) / F1 指標。

For class $c$ let $TP_c, FP_c, FN_c$ be true positives, false positives, false negatives.

$$P_c = \frac{TP_c}{TP_c + FP_c + \varepsilon}, \quad R_c = \frac{TP_c}{TP_c + FN_c + \varepsilon}$$

$$F1_c = \frac{2 P_c R_c}{P_c + R_c + \varepsilon}$$

Small $\varepsilon$ avoids division-by-zero (implemented implicitly by Python float semantics when denominators > 0 in practice). Macro F1:

$$\mathrm{Macro\text{-}F1} = \frac{1}{K} \sum_{c=1}^{K} F1_c$$

Accuracy:

$$\mathrm{Acc} = \frac{1}{n} \sum_{i=1}^{n} \mathbb{1}\{ \arg\max_c z_{i,c} = \arg\max_c y_{i,c} \}$$

### 4.1.3 Knowledge Distillation (KD) Loss

中文: 知識蒸餾損失函數 (KD 損失)。

Denote ground-truth hard cross-entropy: $L_{CE} = - \frac{1}{n} \sum_{i=1}^{n} \sum_{c=1}^{K} y_{i,c} \log p_{s,i,c}$ where $p_{s,i} = \operatorname{softmax}_1(\mathbf{z}_i)$.

Single (fused) teacher soft targets (temperature $T$): $p_{t,i}^T = \operatorname{softmax}_T(\mathbf{t}_i)$, student softened: $p_{s,i}^T = \operatorname{softmax}_T(\mathbf{z}_i)$.

$$L_{KD} = (1-\alpha) L_{CE} + \alpha T^2 \cdot \frac{1}{n} \sum_{i=1}^{n} \mathrm{KL}\big(p_{t,i}^T \;\Vert\; p_{s,i}^T\big)$$

We use $\alpha = 0.5, T=2$. (Factor $T^2$ preserves gradient magnitude per Hinton et al.).

### 4.1.4 Multi-Teacher Split vs Fused Targets

中文: 多教師分離 (Split) 與 融合 (Fused) 目標之差異。

Split strategy draws $M$ separate softened teacher distributions $p_{t,i}^{(m),T}$ each forward pass; loss averages KL terms:

$$L_{KD\text{-}split} = (1-\alpha)L_{CE} + \alpha T^2 \cdot \frac{1}{M n} \sum_{m=1}^{M} \sum_{i=1}^{n} \mathrm{KL}\big(p_{t,i}^{(m),T} \Vert p_{s,i}^T\big)$$

Fused strategy first averages probabilities:

$$\bar{p}_{t,i}^T = \frac{1}{M} \sum_{m=1}^{M} p_{t,i}^{(m),T}$$

Then standard KD applies with $p_{t,i}^T = \bar{p}_{t,i}^T$.

Information Smoothing Observation: Under independence approximation among teachers for class $c$, variance shrinks as

$$\operatorname{Var}[\bar{p}_{t,i}^T(c)] = \frac{1}{M^2}\sum_{m=1}^M \operatorname{Var}[p_{t,i}^{(m),T}(c)]$$

reducing inter-teacher disagreement signal the student could otherwise exploit (empirical "over-smoothing").

### 4.1.5 Decoupled Knowledge Distillation (DKD) Simplified Form

中文: 解耦式知識蒸餾 (DKD) 簡化公式。

Following Zhou et al. (Decoupled KD), logits are partitioned into target (ground-truth) class $g$ and non-target classes $\mathcal{N}$. Let $p_{t,i,g}^T, p_{s,i,g}^T$ be softened probabilities for class $g$; similarly define normalized distributions over non-target classes:

$$\tilde{p}_{t,i,j}^T = \frac{p_{t,i,j}^T}{1 - p_{t,i,g}^T}, \quad \tilde{p}_{s,i,j}^T = \frac{p_{s,i,j}^T}{1 - p_{s,i,g}^T}, \quad j \in \mathcal{N}$$

Target-Class KD component:

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

Non-Target KD component:

$$L_{NCKD} = \frac{1}{n} \sum_{i=1}^{n} \sum_{j \in \mathcal{N}} \tilde{p}_{t,i,j}^T \log \frac{\tilde{p}_{t,i,j}^T}{\tilde{p}_{s,i,j}^T}$$

We employ a weighted combination (notation aligned to our earlier shorthand with parameters $\alpha, \beta$):

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

In four-way split experiments $\beta=4$ (historic sweeps also tried $\beta=2,8$). Observed diminishing returns with increased teacher diversity.

### 4.1.6 Calibration Metrics

中文: 校準評估指標 (ECE、Brier、NLL)。

Let predicted confidence $c_i = \max_c p_{s,i,c}$ and predicted class $\hat{y}_i = \arg\max_c p_{s,i,c}$.

1. 1. Expected Calibration Error (ECE) with $B$ bins $B_b$:

$$\mathrm{ECE} = \sum_{b=1}^{B} \frac{|B_b|}{n} \Big| \mathrm{acc}(B_b) - \mathrm{conf}(B_b) \Big|$$

where $\mathrm{acc}(B_b) = \frac{1}{|B_b|} \sum_{i \in B_b} \mathbb{1}\{\hat{y}_i = y_i\}$, $\mathrm{conf}(B_b) = \frac{1}{|B_b|} \sum_{i \in B_b} c_i$.

2. 2. Brier Score:

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

3. 3. Negative Log-Likelihood (NLL):

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

Temperature scaling selects $T^*$ minimizing NLL (or ECE proxy) on a calibration split; post-hoc we replace $p_{s,i}$ by $\operatorname{softmax}_{T^*}(\mathbf{z}_i)$ without retraining.

### 4.1.7 Deployment Consideration Metric Coupling

中文: 部署決策時的多指標權衡與綜合考量。

- Macro F1: selection for accuracy / fairness across classes.
- ECE / Brier / NLL: selection for reliability of probability outputs (downstream thresholding / ranking stability).
- We will interpret calibration changes jointly: large ECE drop with negligible Macro F1 loss is acceptable; large Macro F1 loss for minor ECE gain is not.

These formal definitions support the diversity vs over-smoothing narrative: higher teacher diversity preserves variance across $p_{t,i}^{(m)}$, which—rather than being prematurely averaged—translates into richer gradient signals in $L_{KD\text{-}split}$ and empirically higher Macro F1.

## 5. Results Summary (Macro F1)

| Family | Mode | Seeds | Best Macro F1 (per seed) | Mean | Notes |
|--------|------|-------|--------------------------|------|-------|
| Pairwise (6A) | KD | 1337,2025 | 0.71878 / 0.71613 | 0.71746 | Simple, strong |
| Pairwise (historic) | DKD β=4 | 1337,2025 | ~0.71918 / ~0.71313 | ~0.71616 | Earlier; small uplift not consistent |
| Hybrid (6B) | KD | 1337,2025 | 0.70636 / 0.70494 | 0.70565 | Adds ViT but lower post-alignment |
| Hybrid (6B) | DKD (β best) | 1337,2025 | 0.71045 / 0.70996 | 0.71021 | DKD marginal gain vs KD |
| Split CNN+ViT (6D) | KD | 1337,2025 | λ=0.7/0.3: 0.71296 / 0.71431 | 0.71364 | λ=0.7/0.3 > others |
| Split CNN+ViT (6D) | DKD β=4 | 1337,2025 | λ=0.7/0.3: 0.71824 / 0.71218 | 0.71521 | Near pairwise KD |
| Four-Way Split (6E) | KD | 1337,2025 | 0.72263 / 0.72107 | 0.72185 | Highest overall |
| Four-Way Split (6E) | DKD β=4 | 1337,2025 | 0.72020 / 0.72071 | 0.72045 | KD > DKD here |

## 6. Per-Class F1 Differential (Illustrative)

Selected comparison: Pairwise KD s1337 vs Four-Way Split KD s1337

| Class | Pairwise KD F1 | Four-Way KD F1 | Δ |
|-------|----------------|----------------|---|
| angry | 0.6523 | 0.6523 | ~0.0 |
| disgust | 0.6381 | 0.6381 | ~0.0 |
| fear | 0.6974 | 0.6974 | ~0.0 |
| happy | 0.8627 | 0.8627 | ~0.0 |
| neutral | 0.7589 | 0.7589 | ~0.0 |
| sad | 0.6217 | 0.6217 | ~0.0 |
| surprise | 0.8273 | 0.8273 | ~0.0 |

Note: Current summary rows show near-identical per-class decimals to 4 d.p.; further discrimination may require extended precision or calibration metrics to reveal difference; improvements likely distributed, not class-specific.

## 7. Root Cause of Earlier Split Underperformance

| Symptom | Value (Pre-Fix) | Value (Post-Fix) | Cause | Resolution |
|---------|-----------------|------------------|-------|------------|
| Split λ=0.5/0.5 macro F1 | ~0.63–0.64 | ~0.709–0.713 | Misaligned softlabel lengths/order | Enforce canonical index + `--require-aligned` |
| Warning | Truncation applied | None | Hidden mismatch | Alignment diagnostic + purge |

Misalignment explained severe macro F1 suppression (>6 points). Once corrected, control group surpassed core group performance.

### 7.1 Hypothesis Evaluation: Dataset Mismatch vs Label Mismatch

Two separate (and initially overlapping) hypotheses were considered to explain the very low early Split Multi-Teacher results (≈0.63–0.65 macro F1):

4. 1. Dataset Index Mismatch Hypothesis:

- Possibility: The split and four-way control group runs used an incorrect dataset index (e.g. `dataset_index_extended_v8.csv` or an earlier non-dedup variant) while the core group (pairwise / hybrid) used the canonical deduplicated index, inflating the apparent performance gap.
- Evidence Collected: After forensic review, core group historical runs (pre-alignment) ALSO referenced non-canonical indices in several phases, yet their macro F1 values (≈0.706–0.714) were already close to the post-alignment retrain numbers. This reduces the likelihood that index mismatch alone explains the ≈6–8 point deficit of early split runs.

5. 2. Label (Softlabel) Length / Ordering Mismatch Hypothesis:

- Observation: Training logs for early split multi-teacher runs emitted repeated truncation warnings indicating probability array length differences across teacher sources. Truncation silently aligned tensor shapes without validating order consistency.
- Mechanism: If image ordering diverges between teacher probability directories, row-wise averaging combines probabilities from different images, corrupting class distributions and disproportionately hurting macro F1 (especially minority classes).
- Post-Fix Outcome: After enforcing canonical index re-export + strict `--require-aligned` (hard fail on any mismatch), split variant macro F1 rose from ~0.63–0.65 to >0.705 and four-way split KD reached 0.7226. This sharp recovery is consistent with removal of a destructive label misalignment artifact rather than a benign index swap.

Inference (Primary Cause): Because (a) core group runs exposed to the same (incorrect) index family still produced near-final-level macro F1, and (b) only the control group's performance rebounded drastically after eliminating softlabel misalignment, the principal degradation factor was the label ordering/length mismatch—NOT the choice of dataset index itself. The dataset index discrepancy remains a reproducibility concern, but it does not account for the magnitude of the early control group deficit.

Data Deletion & Reproduction Plan:

- All misaligned split / four-way experimental directories were deleted to prevent accidental reuse; only canonical `_ixNextAffFull` suffixed aligned runs remain.
- Reproduction of the failed (misaligned) state is intentionally deferred to avoid reintroducing corrupt supervision; if ever needed for an ablation study, it would require:

6. 1. Re-exporting one teacher with a deliberately altered index ordering.
7. 2. Disabling `--require-aligned` and allowing truncation (currently blocked in production pipeline).
8. 3. Logging hash digests of path orderings to conclusively tie degradation to ordering corruption.

- This reproduction is low priority and will only be scheduled if a paper appendix requires a controlled failure demonstration.

## 8. Analysis / 解讀

### 8.0 Formal Comparative Performance Framing (Core vs Control Families)

This section formalizes the comparative performance interpretation across defined families using aligned (canonical index) MobileNetV3 student results (refer to the formal metric and loss definitions in Section 4.1 / 參見 §4.1 數學公式). Macro F1 values are shown as percentage ranges (macro_F1 × 100) for clarity; bolded ranges indicate the strongest family within each conceptual grouping.

#### 8.0.1 Grouped Outcome Summary

| Strategy / Family Set | Family IDs | Aligned Best Macro F1 Range (%) | Representative Mean (%) | Notes |
|----------------------|-----------|--------------------------------|------------------------|-------|
| Single / Fused Pairwise & Hybrid ("Single ensemble teacher" – effectively pre-fused sources) | 6A (Pairwise KD), 6B (Hybrid KD/DKD) | 70.49 – 71.88 | ~71.1 (pairwise KD mean 71.746%, hybrid KD mean 70.565%) | Low noise, limited headroom |

| Split Dual-Source (CNN+ViT) λ sweep | 6D (KD/DKD) | 70.54 – 71.82 (KD & DKD β=4) | KD λ=0.7/0.3 mean 71.364%, DKD λ=0.7/0.3 mean 71.521% | Improves when emphasizing CNN (λ=0.7) |

| Four-Way Split Multi-Teacher (equal 0.25 each) | 6E (KD/DKD β=4) | **72.02 – 72.26** (KD) / 72.02 – 72.07 (DKD) | **72.185% (KD mean)** | Highest and most consistent; KD > DKD |

| Fused Four Equal (single averaged label) | 6C (Planned) | (Not yet aligned; historic misaligned 70.18 – 70.76) | TBD | Will test over-smoothing hypothesis |

Interpretation: Direct exposure to four independent teacher distributions (6E split) yields a reproducible ≈ +0.4 to +0.7 absolute macro F1 uplift over the strongest single fused teacher configuration (pairwise KD) and ≈ +1.5 to +1.6 over hybrid KD. Fusing the four teachers into a single probability vector prior to training (historic misaligned fused-four attempts, family 6C) underperformed both multi-source splits and even simpler pairwise KD; this supports the diversity-preservation hypothesis.

### 8.0.2 Why Four-Way Split Outperforms Fused / Single Teacher Approaches

9. 1. Diversity Retention: Each teacher retains its unique calibration landscape and error profile; sampling all four independently exposes the student to richer inter-example variance than a pre-averaged consensus.
10. 2. Avoidance of Over-Smoothing: Probability averaging before training attenuates inter-class logit contrast, reducing gradient sharpness. Multi-source mixing at training time lets the student reconcile differences dynamically instead of receiving homogenized targets.
11. 3. Implicit Regularization: Conflicting but valid teacher signals act similarly to stochastic regularization, discouraging overfit to any single teacher's idiosyncrasies.
12. 4. Marginal DKD Benefit Exhaustion: With already diverse soft targets, decoupled contrastive emphasis (β term) adds little; four-way KD slightly surpasses four-way DKD (72.185% vs 72.045% mean) indicating diminishing returns from additional contrastive loss components.

### 8.0.3 Positioning of Split Dual-Source (6D)

Split CNN+ViT (λ sweeps) closes part of the gap versus pairwise but plateaus below four-way. Emphasizing CNN-heavy weighting (λ=0.7/0.3) yields the best λ result, suggesting the

ViT contributes complementary but lower-marginal incremental signal relative to adding two more heterogeneous CNN architectures (ConvNeXt + direct EfficientNet-B3) present in four-way split.

### 8.0.4 Planned Validation: Fused Four Re-Appearance (Family 6C)

To rigorously validate the over-smoothing hypothesis, a future controlled aligned experiment will:

- Produce a single fused softlabel directory (equal weights) under canonical index.
- Compare KD macro F1 and calibration (ECE, Brier, NLL) against four-way split KD.
- Examine per-class F1 deltas: expectation is that fused labels may slightly regress minority class F1 (disgust, fear) due to probability mass averaging.
- If fused performance < split but shows superior calibration, a hybrid deployment strategy (split-trained model + post-hoc temperature scaling) remains preferable.

### 8.0.5 Executive Framing for Report Abstract

Compared to single or pre-fused ensemble teacher supervision (70.5–71.9% macro F1), the equal-weight four-teacher split strategy consistently attains 72.0–72.3%, while a naïve pre-fusion of the same four teachers (historic misaligned runs; planned aligned replication) underperforms (≈70.2–70.8%). This pattern evidences that preserving inter-teacher diversity during training confers generalization gains beyond what probabilistic averaging can deliver.

### 8.0.6 Key Takeaway Statement

"Multi-teacher split distillation (four independent teachers, equal weighting) provides the best accuracy because it maximizes diversity and minimizes pre-training information loss; probability fusion prematurely smooths informative inter-teacher disagreements, leading to systematically lower macro F1."

- Four-Way Split KD advantage (≈+0.004 over Pairwise KD) is modest but consistent; indicates value in richer teacher heterogeneity even without DKD contrastive pressure.
- DKD shows diminishing or negative marginal benefit when teacher diversity already softens logits sufficiently (four-way case).
- Hybrid underperforms vs pairwise + ViT split; suggests pre-fused hybrid probabilities may over-smooth informative inter-class margins compared to retaining separable sources.
- Pairwise remains a compelling fallback (simpler export pipeline) with only small deficit to four-way.

Root Cause Attribution Integration: The small performance gap (≈+0.003–0.004 absolute macro F1) favoring four-way KD is now evaluated in a context where label integrity is guaranteed. Earlier much larger gaps (≥0.06) are excluded from decision-making because they were products of corrupted supervision (label mismatch), not intrinsic limits of the split / multi-source strategy.

### 8.0.7 Fused Four Equal (6C) Aligned Replication Plan

Status: Not yet executed under canonical alignment. Historic (misaligned) macro F1 range ≈70.18–70.76 (KD/DKD mixed). This replication will generate authoritative aligned metrics for the fused-equal strategy.

Planned Procedure:

13. 1. Export fused softlabels (already scripted; directory: `experiments/softlabels/four_teacher_equal_ixNextAffFull`). If missing, re-run export step.
14. 2. Run DKD ($\beta$=4) seeds {1337,2025}: outputs `experiments/student_mbv3l_four_teacher_ixNextAffFull/b4_s{seed}`.
15. 3. Run KD seeds {1337,2025}: outputs `experiments/student_mbv3l_four_teacher_ixNextAffFull/kd_four_s{seed}`.
16. 4. Collect results: `python scripts/collect_results_summary.py --incremental --tz-offset 8`.
17. 5. Inject metrics (auto): `python scripts/update_report_6c_results.py` will patch this section and the tables (6C row).

Success Criteria:

- If 6C (KD mean) ≥ 6A (pairwise KD mean) + 0.002 and still < 6E mean by ≤0.004, over-smoothing hypothesis partially validated (diversity still superior but fused competitive).
- If 6C (KD mean) ≈ 6E (|Δ| ≤ 0.0015), over-smoothing hypothesis weakened; diversity benefit may stem from regularization scheduling rather than probability variance retention.
- If 6C (KD mean) << 6A (drop ≥0.005), strong confirmation that premature averaging discards critical teacher disagreement signal.

Reporting Fields To Be Added Post-Run:

- 6C (KD) best per-seed macro F1, mean.
- 6C (DKD β=4) best per-seed macro F1, mean.
- Calibration deltas vs 6E (once temperature scaling applied).

Placeholder Metrics (to be replaced automatically):

| Family | Mode | Seeds | Best Macro F1 (per seed) | Mean | Notes |
|--------|------|-------|--------------------------|------|-------|
| Fused Four (6C) | KD | 1337,2025 | TBD / TBD | TBD | Awaiting aligned run |
| Fused Four (6C) | DKD β=4 | 1337,2025 | TBD / TBD | TBD | Awaiting aligned run |

Automation: After training, run `python scripts/update_report_6c_results.py` to substitute this placeholder block and update Section 5 and 8.0.1 tables.

## 9. Recommendations

18. 1. Provisionally adopt Four-Way Split KD (seed 1337) checkpoint `best.pt` as deployment candidate.
19. 2. Perform temperature scaling calibration; retain pairwise KD (s1337) as latency + robustness fallback.
20. 3. Skip further DKD β sweeps for four-way unless calibration reveals systematic confidence miscalibration.
21. 4. Implement optional logit-space fusion experiment (future) to test if pre-softmax additive ensemble lifts macro F1 or calibration beyond probability averaging.

## 10. Pending Work / 待辦

| Task | Priority | Description |
|------|----------|-------------|
| Calibration pass | High | Fit temperature on validation/test; write `calibration.json` |
| Latency benchmark | High | ONNX export & FPS comparison (pairwise vs four-way) |
| Model packaging | High | Integrate chosen checkpoint into real-time pipeline |
| Extended precision diff | Medium | Recompute per-class F1 w/ higher precision to detect subtle gains |
| Fused four single softlabels | Medium | Compare against split multi-source runtime pipeline complexity |
| Fused four aligned run (6C) | High | Execute KD & DKD β=4 under canonical index; update Section 8.0.7 |
| Evaluation-only ViT per-class | Low | Only if class drift analysis vs ViT needed |

## 11. Risk & Mitigation

| Risk | Impact | Mitigation |
|------|--------|-----------|
| Calibration not improving ECE | Delayed deployment | Multi-temperature sweep, isotonic fallback |
| Latency regression (none expected) | FPS target miss | Quantize / prune candidate if needed |
| Data leakage in future exports | Invalid comparisons | Retain canonical index naming + alignment check in CI |

## 12. Conclusion

Control Group strategies—especially Four-Way Split KD—now outperform the original Core Group baselines after rigorous alignment enforcement. The earlier narrative that pairwise +

hybrid were competitive leaders was an artifact of unaligned split runs. With alignment fixed, diversity-driven teacher composition yields a reproducible yet incremental performance gain. Focus shifts to calibration, latency optimization, and productionization of the Four-Way KD student.

## 13. Appendix: Source Mapping

| Experiment Dir (relative) | Family | Seed | best_test_macro_f1 |
|---------------------------|--------|------|--------------------|
| experiments/student_mnv3_pairwise_cnn_ixNextAffFull/kd_pairwise_w0p7_0p3_s1337 | Pairwise KD | 1337 | 0.7187765 |
| experiments/student_mnv3_pairwise_cnn_ixNextAffFull/kd_pairwise_w0p7_0p3_s2025 | Pairwise KD | 2025 | 0.7161278 |
| experiments/student_mnv3_runs_ixNextAffFull/mnv3_hybrid_b2_s1337 | Hybrid DKD | 1337 | 0.7143 (historic pre alignment) |
| experiments/student_mnv3_split_cnn_vit_ixNextAffFull/kd_lam_0.7_0.3_s2025 | Split KD | 2025 | 0.7143131 |
| experiments/student_mnv3_fourway_split_ixNextAffFull/kd_s1337 | Four-Way KD | 1337 | 0.7226347 |
| experiments/student_mnv3_fourway_split_ixNextAffFull/kd_s2025 | Four-Way KD | 2025 | 0.7210732 |
| experiments/student_mnv3_fourway_split_ixNextAffFull/b4_s1337 | Four-Way DKD | 1337 | 0.7201981 |
| experiments/student_mnv3_fourway_split_ixNextAffFull/b4_s2025 | Four-Way DKD | 2025 | 0.7207104 |

(Values match `results_summary.csv` at time of report generation; macro F1 truncated to 7 significant digits where shown.)