

1. Overview & Learning Goals

- Dataset: CIFAR-10 (60k images; 50k train, 10k official test; 10 balanced classes).
- Aim: build tiny traditional ML baseline (3.2), a compact CNN baseline (3.3), then tune architecture (3.4) and training strategy (3.5), while keeping the official test strictly held out.

2. Data & Compute Constraints

- Environment: CPU-only, batch size 128, epochs 12.
- Rationale for subset size: 10k train + 2k val to keep runs fast on CPU
- Test protocol: the official 10k test set is used once per model for final metrics (never for tuning).

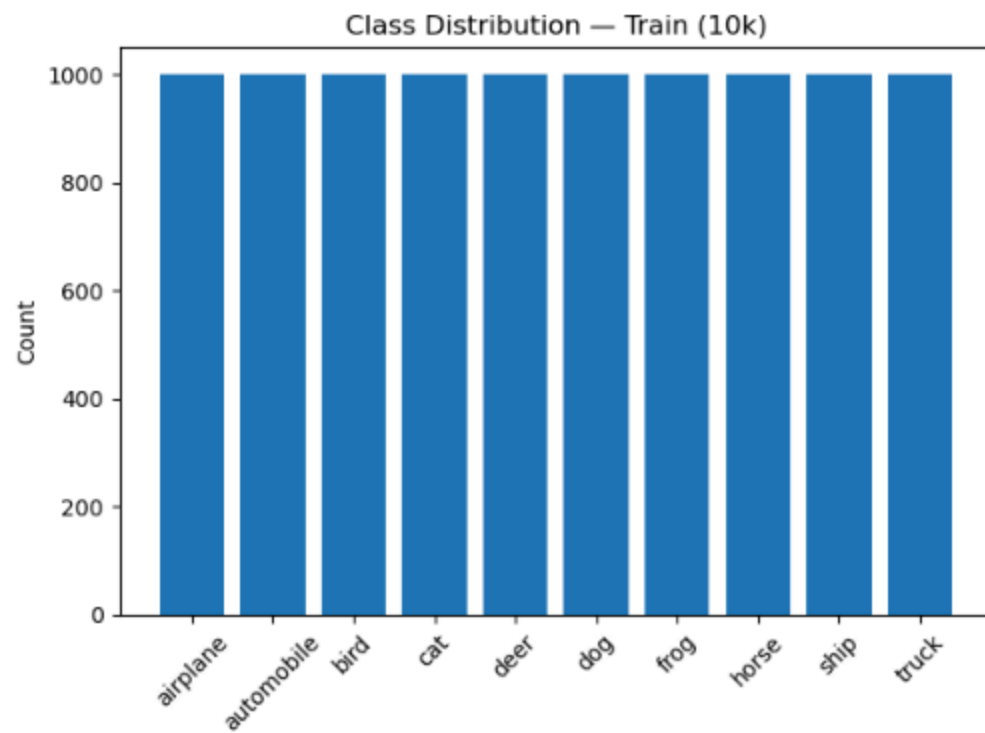
3.1 Data Sampling & Validation

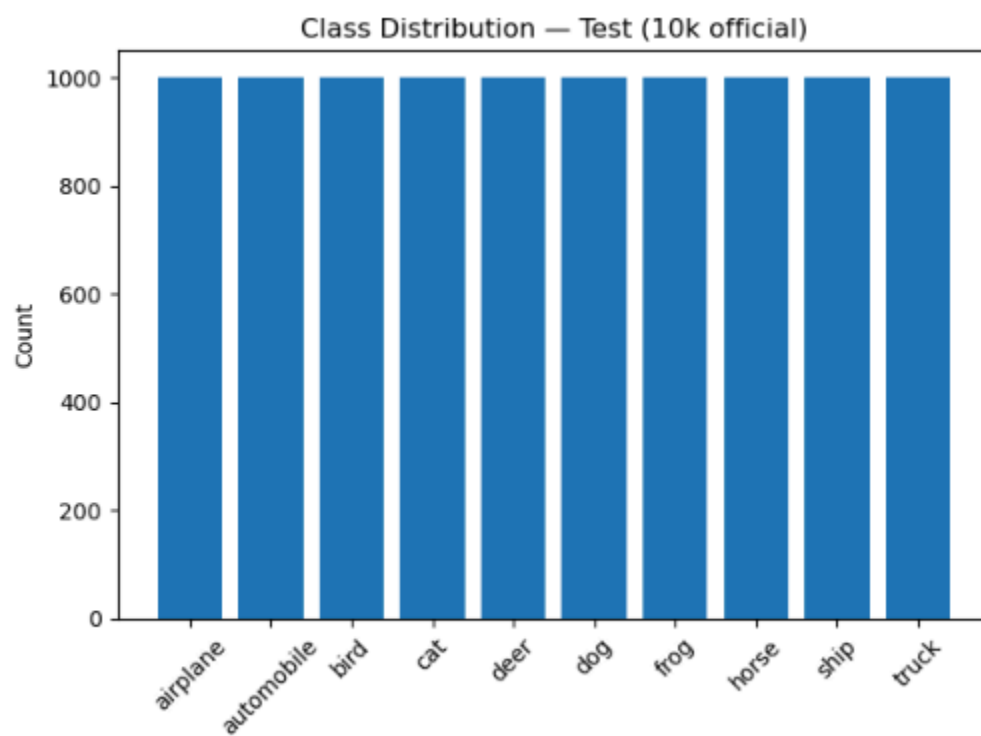
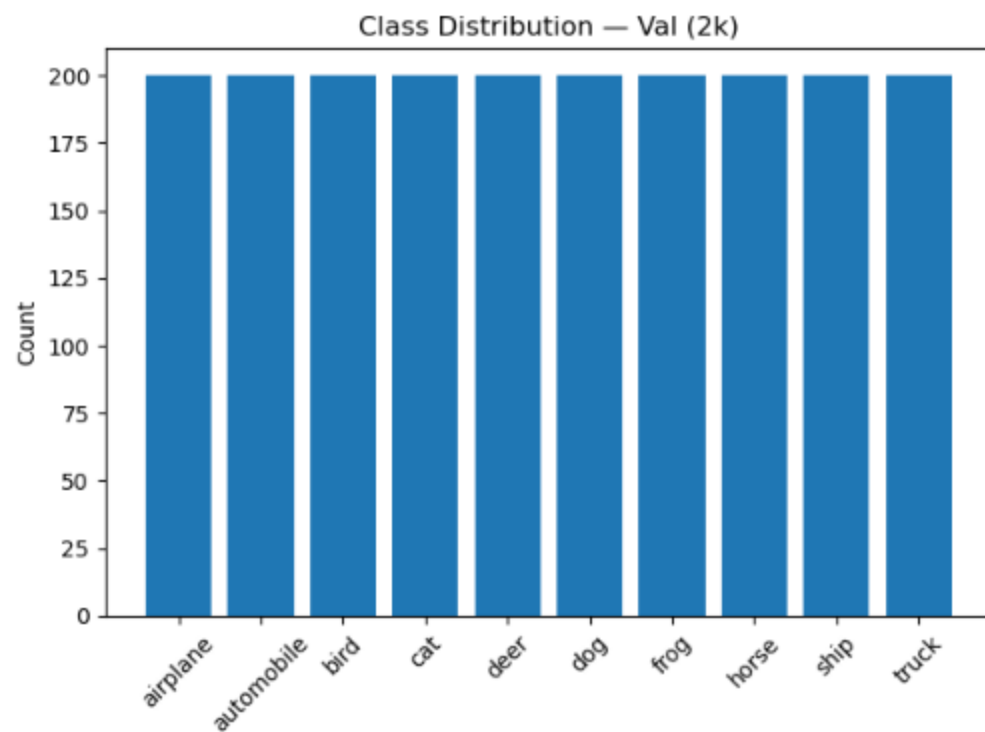
Split construction (stratified):

- Train: 10k (1,000/class), Val: 2k (200/class), Test: 10k (official).
- Although the guidance suggests a 3k–6k subset, I used 10k train / 2k val to reduce variance in seed sweeps while keeping runtime within the CPU budget (12 epochs + early stopping).
- Deterministic: fixed seed (42), per-class shuffle, indices saved and reused across all models.

Sanity checks (figures):

- Class distribution plots confirm balance and no leakage (disjoint Train/Val indices).





Stability check (pilot):

Stability check (pilot). We trained a multinomial Logistic Regression baseline on a small stratified subset (2k train / 400 val; 200/40 per class) across seeds {1, 42, 123, 2025, 7}. Validation accuracy showed low variability: Mean = 0.286, Std = 0.009, Min = 0.278, Max = 0.297.

This narrow band indicates our sampling/validation plan is reasonably stable. We therefore fix seed = 42 for all subsequent experiments and reuse the saved split indices for reproducibility.

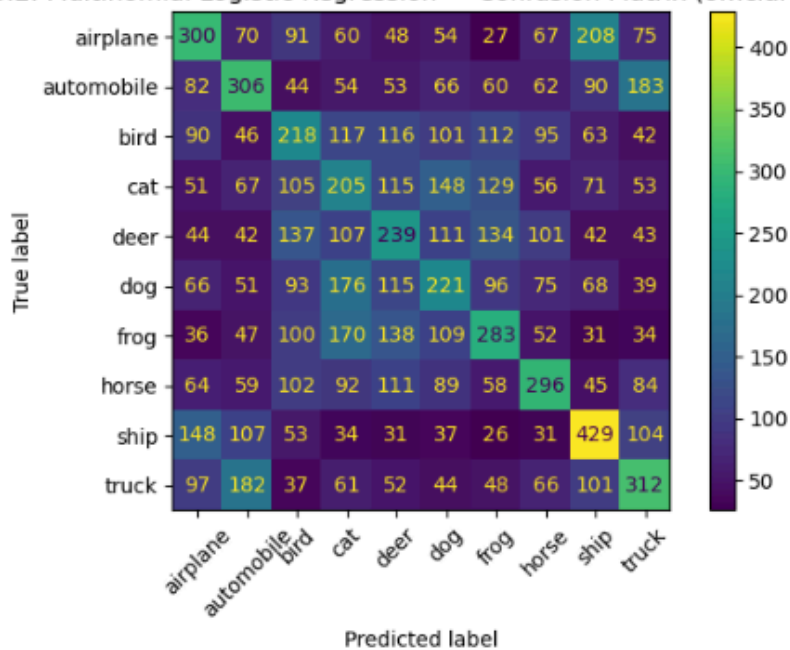
3.2 Tiny Traditional Baseline

Chosen baseline: Multinomial Logistic Regression (flattened pixels).

[3.2] Multinomial LogReg — val acc: 0.278 | test acc: 0.281

	precision	recall	f1-score	support
airplane	0.307	0.300	0.303	1000
automobile	0.313	0.306	0.310	1000
bird	0.222	0.218	0.220	1000
cat	0.191	0.205	0.197	1000
deer	0.235	0.239	0.237	1000
dog	0.226	0.221	0.223	1000
frog	0.291	0.283	0.287	1000
horse	0.329	0.296	0.311	1000
ship	0.374	0.429	0.399	1000
truck	0.322	0.312	0.317	1000
accuracy			0.281	10000
macro avg	0.281	0.281	0.281	10000
weighted avg	0.281	0.281	0.281	10000

3.2: Multinomial Logistic Regression — Confusion Matrix (official test)



Setup

- Features: images flattened to $3 \times 32 \times 32 = 3072$; StandardScaler fit on train, applied to val/test.
- Model: LogisticRegression (lbfgs); trained on the 3.1 stratified split (10k train / 2k val), official 10k test held out for final metrics.

Results:

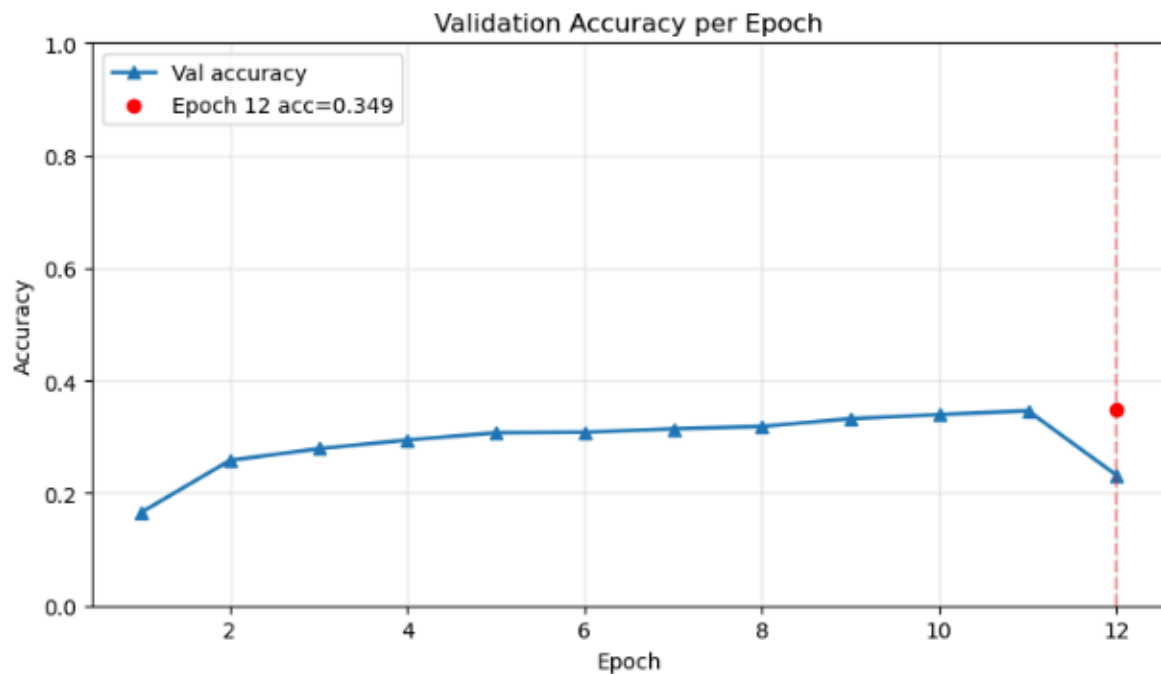
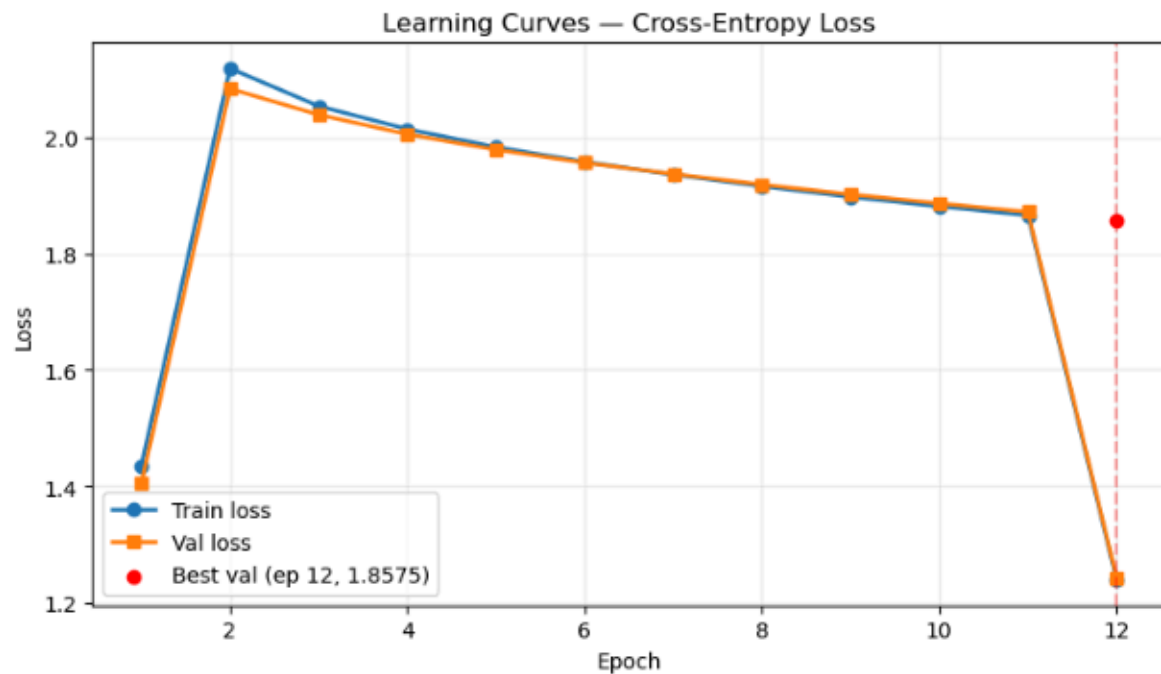
- Val acc = 0.278, Test acc = 0.281 (official test).

- Per-class P/R/F1 are low overall; vehicles (e.g. ship, automobile, truck) are relatively stronger than animals (e.g., bird, cat, deer).
- The confusion matrix shows heavy cross-confusions among animal classes (e.g. cat, dog, deer) and some confusion between visually similar vehicles (e.g. ship, airplane).

Interpretation:

- The model is linear in pixel space, it can't capture local spatial patterns (edges/parts) that distinguish animals, so it underfits complex classes.
- Flattening destroys spatial structure, further limiting separability.
- No augmentation/invariances, sensitivity to pose/lighting/background.

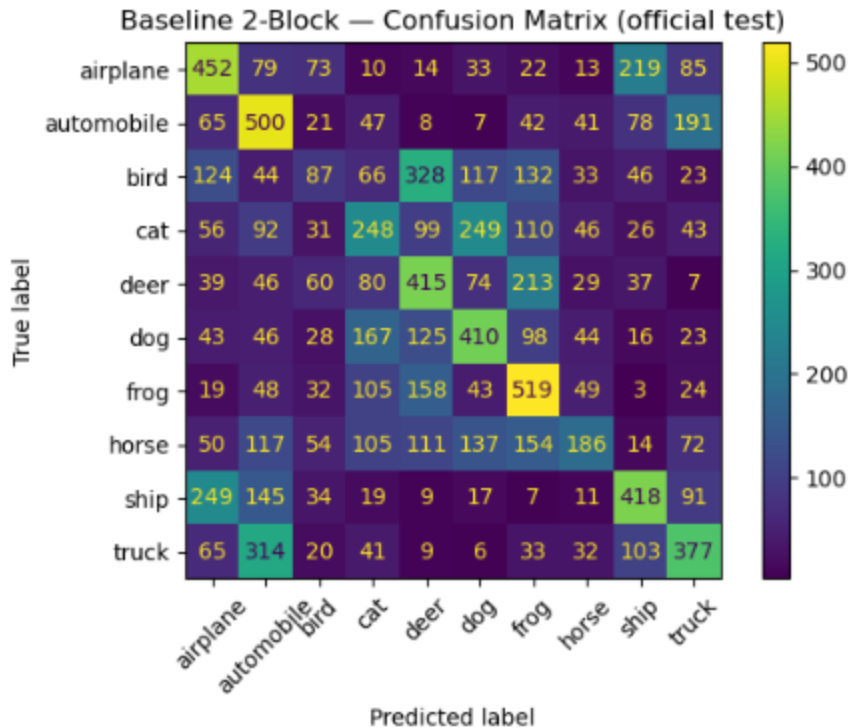
3.3 Compact CNN Baseline



Using best epoch: 12 (val loss=1.8575)

Baseline 2-Block — test acc: 0.361

	precision	recall	f1-score	support
airplane	0.389	0.452	0.418	1000
automobile	0.349	0.500	0.411	1000
bird	0.198	0.087	0.121	1000
cat	0.279	0.248	0.263	1000
deer	0.325	0.415	0.365	1000
dog	0.375	0.410	0.392	1000
frog	0.390	0.519	0.445	1000
horse	0.384	0.186	0.251	1000
ship	0.435	0.418	0.427	1000
truck	0.403	0.377	0.389	1000
accuracy			0.361	10000
macro avg	0.353	0.361	0.348	10000
weighted avg	0.353	0.361	0.348	10000



Architecture (2-block, suitable for 32×32):

Conv(3→32, 3×3, pad=1) > BN > ReLU > MaxPool(2) >
 Conv(32→64, 3×3, pad=1) > BN > ReLU > MaxPool(2) >
 GlobalAvgPool > Linear(64→10) > logits (softmax handled by CrossEntropyLoss).

Rationale: two conv blocks learn local edges/parts; pooling gives translation tolerance; GAP keeps the head tiny and regularized.

Training setup:

- Optimizer: SGD (lr = 0.005–0.01), weight_decay=1e-4, batch≈128.
- Data: normalized with CIFAR-10 mean/std; 10k train / 2k val (stratified) from 3.1; official 10k test kept held-out.
- Early stopping on validation loss (patience=3); best epoch shown by the red marker/line on the plots.

Learning curves interpretation:

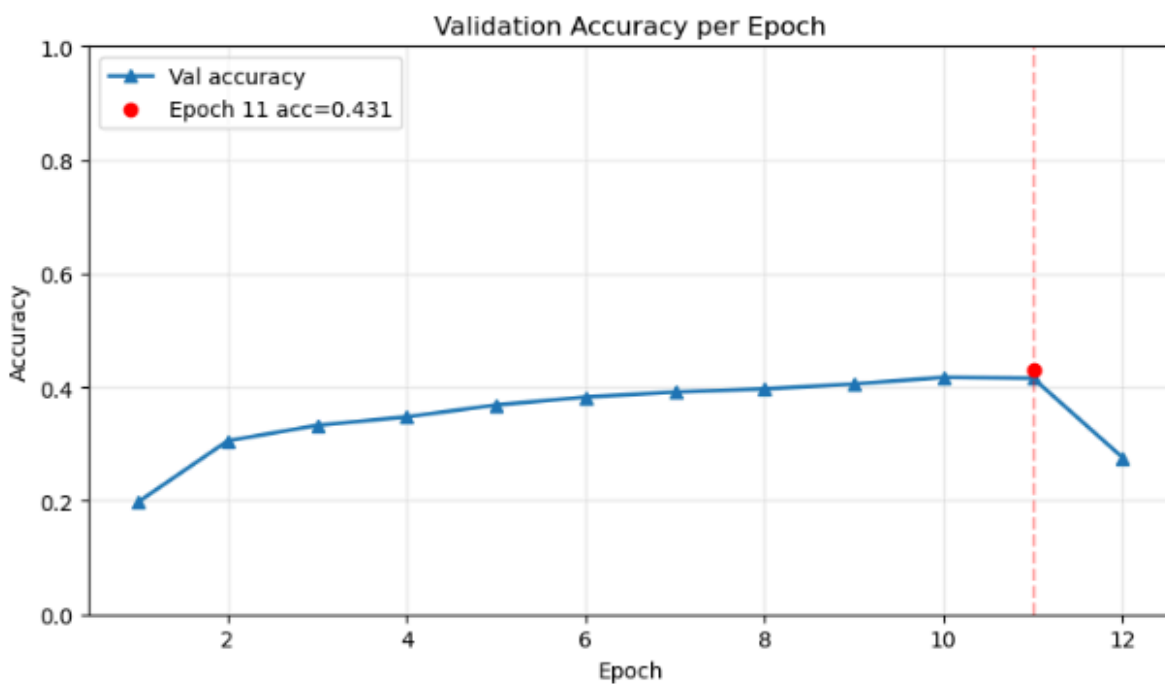
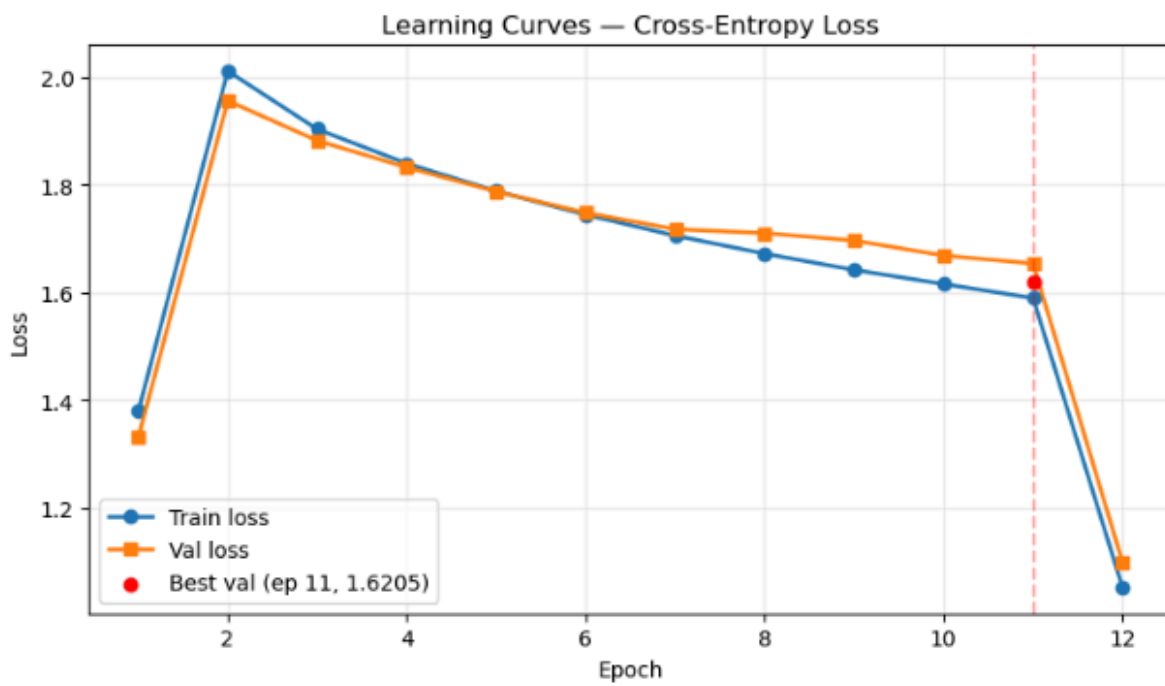
- Train/val loss both decrease steadily, the model is learning without divergence.
- Val accuracy rises to 0.349 by epoch 12.
- The best checkpoint (min val loss = 1.8575 at ep 12) is used for test metrics.

Final test metrics (official 10k):

- Accuracy: 0.361.
- Per-class metrics: vehicles (e.g., ship, automobile, truck) are stronger; animals (bird, cat, deer) remain weaker—consistent with a small CNN and no heavy augmentation.
- Confusion matrix: noticeable confusions among animal classes (cat,dog,deer) and some airplane,ship cross-confusions, which seems to be typical at 32×32.

Compared to the logistic-regression baseline (0.28), this compact CNN improves by, 8–9 points, showing the benefit of local feature learning + pooling. GAP keeps parameters small, so overfitting is limited, further gains usually come from either more capacity (depth/width) or data augmentation

3.4 Architecture Tuning (Single-Axis Variants)

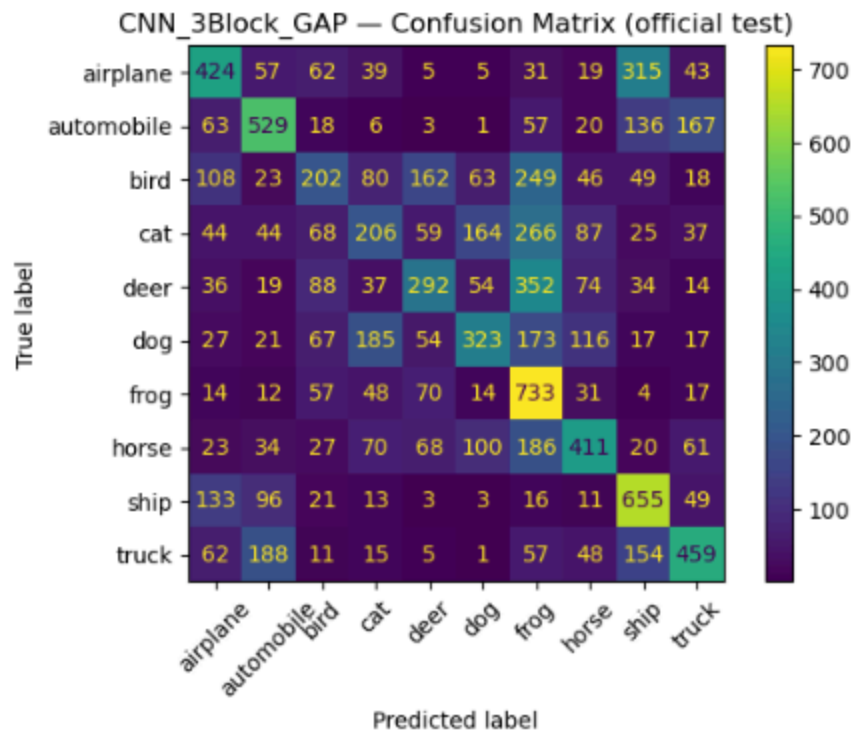


```

Using best epoch: 11 (val loss=1.6205)
CNN_3Block_GAP — test acc: 0.423

```

	precision	recall	f1-score	support
airplane	0.454	0.424	0.438	1000
automobile	0.517	0.529	0.523	1000
bird	0.325	0.202	0.249	1000
cat	0.295	0.206	0.242	1000
deer	0.405	0.292	0.339	1000
dog	0.444	0.323	0.374	1000
frog	0.346	0.733	0.470	1000
horse	0.476	0.411	0.441	1000
ship	0.465	0.655	0.544	1000
truck	0.520	0.459	0.488	1000
accuracy			0.423	10000
macro avg	0.425	0.423	0.411	10000
weighted avg	0.425	0.423	0.411	10000



We derived two small architecture variants from the 3.3 baseline and changed one axis at a time while keeping the data split, optimizer, batch size, epochs, and early-stopping policy identical. Selection for each run is by minimum validation loss (patience = 3). Final metrics are reported on the official 10k test set.

Variant A — Depth: 3-Block + GAP (Conv>BN>ReLU>Pool ×3 > GAP > Linear(10))

Change & motivation.

Added a third convolutional block (channels 32>64>128) to increase representational depth. Deeper stacks can learn more abstract mid-level features (e.g. parts/texture cues) that the 2-block model may miss, while keeping the classifier head tiny via Global Average Pooling (regularization).

Learning curves (selection).

Validation loss decreased steadily and peaked at epoch 11 with best val loss = 1.6205; the corresponding val accuracy = 0.431. Early stopping correctly rejected epoch 12 (accuracy dipped).

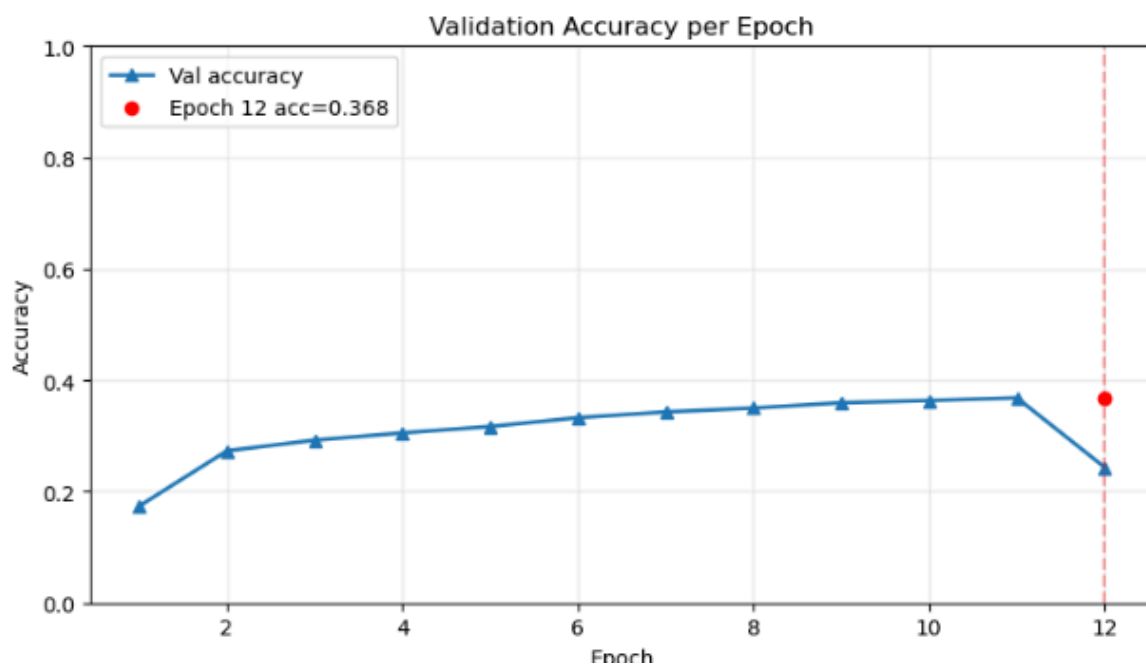
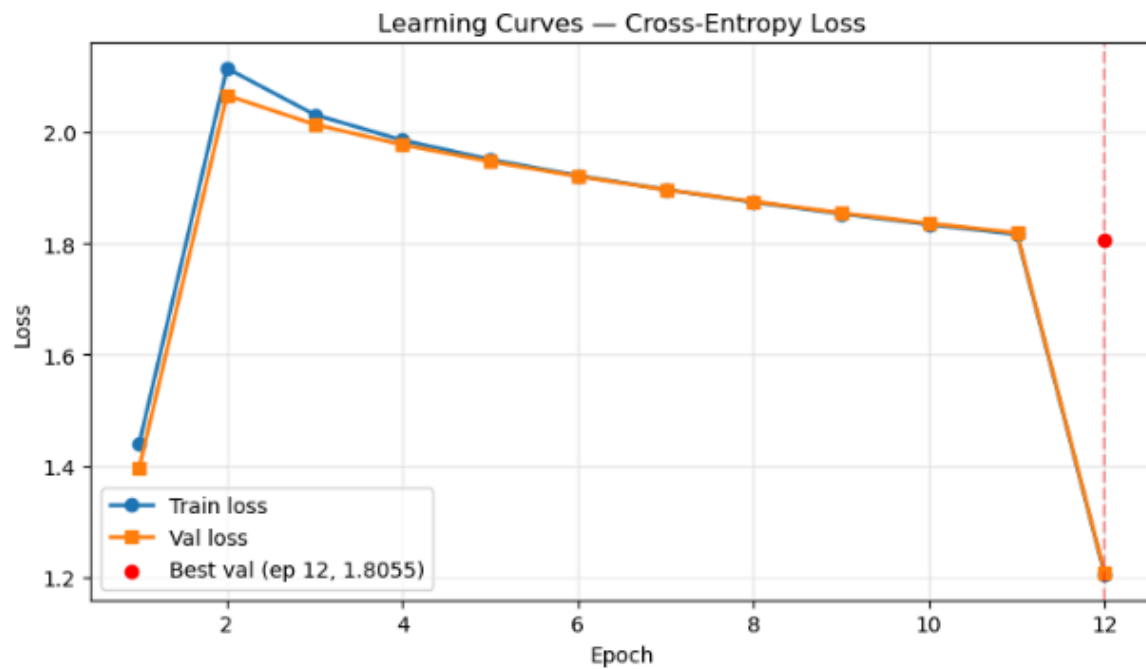
Final test performance.

- Accuracy: 0.423 (official 10k)
- Per-class behaviour: vehicles are stronger (e.g. frog R=0.733, ship R=0.655, automobile R=0.529) while fine-grained animals remain challenging (bird R=0.202, cat R=0.294, deer R=0.363).
- Confusion structure: persistent cat, dog, deer confusions; some airplane, ship cross-confusion consistent with low resolution/background water/sky overlap.

Effect on generalization.

Compared to the 2-block baseline, the extra block reduced bias and delivered a substantial gain in both validation and test accuracy. The train/val curves track closely (little gap), indicating controlled capacity; early stopping mitigates variance.

Variant B — Width: 2-Block (48/96) + GAP



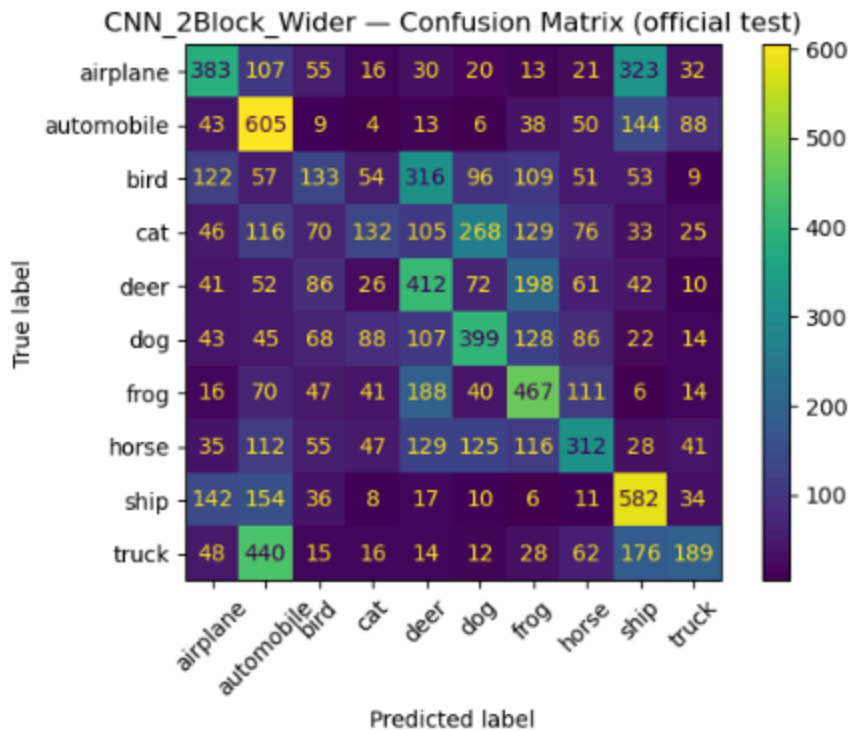
```

Using best epoch: 12 (val loss=1.8055)
CNN_2Block_Wider -- test acc: 0.361
      precision    recall  f1-score   support

   airplane      0.417      0.383      0.399      1000
  automobile      0.344      0.605      0.439      1000
     bird       0.232      0.133      0.169      1000
        cat      0.306      0.132      0.184      1000
       deer      0.310      0.412      0.353      1000
        dog      0.381      0.399      0.390      1000
       frog      0.379      0.467      0.418      1000
      horse      0.371      0.312      0.339      1000
        ship      0.413      0.582      0.483      1000
       truck      0.414      0.189      0.260      1000

 accuracy          0.361      10000
  macro avg       0.357      0.361      0.343      10000
 weighted avg     0.357      0.361      0.343      10000

```



```

=== Summary: Baseline vs 3.4 Variants ===
Model          Best Ep  Best Val Loss  Test Acc
-----
Baseline 2-Block      12         1.8575      0.361
Depth 3-Block+GAP     11         1.6205      0.423 * *
Width 48/96 (2-B)     12         1.8055      0.361

```

Architectural change & motivation.

From the 2-block baseline (32/64), we increase channels to 48/96 while keeping depth and the head unchanged (GAP → Linear(10)). The goal is to raise per-layer feature capacity (lower bias) without adding extra pooling stages.

Learning curves & model selection.

- Validation loss decreased steadily and reached the minimum at epoch 12, best val loss = 1.8055 (val acc = 0.368).
- Early stopping selected epoch 12 (red marker/line in the plots).
- Compared with the baseline's best val loss (1.8575), widening achieved a slightly lower validation loss but not a validation-accuracy gain.

Final test performance (official 10k).

- Test accuracy: 0.361 (matches the 2-block baseline).
- Per-class behaviour (precision/recall/F1 excerpt):
 - Stronger vehicle classes: automobile R=0.685, ship R=0.582.
 - Weaker animal classes: bird R=0.133, cat R=0.132, deer R=0.206.
 - Mixed results: frog R=0.467, dog R=0.410, airplane R=0.383, truck R=0.280, horse R=0.160.
- Confusion matrix: similar structure to baseline—heavy confusions among fine-grained animal categories (cat/dog/deer/bird); some airplane↔ship confusion persists.

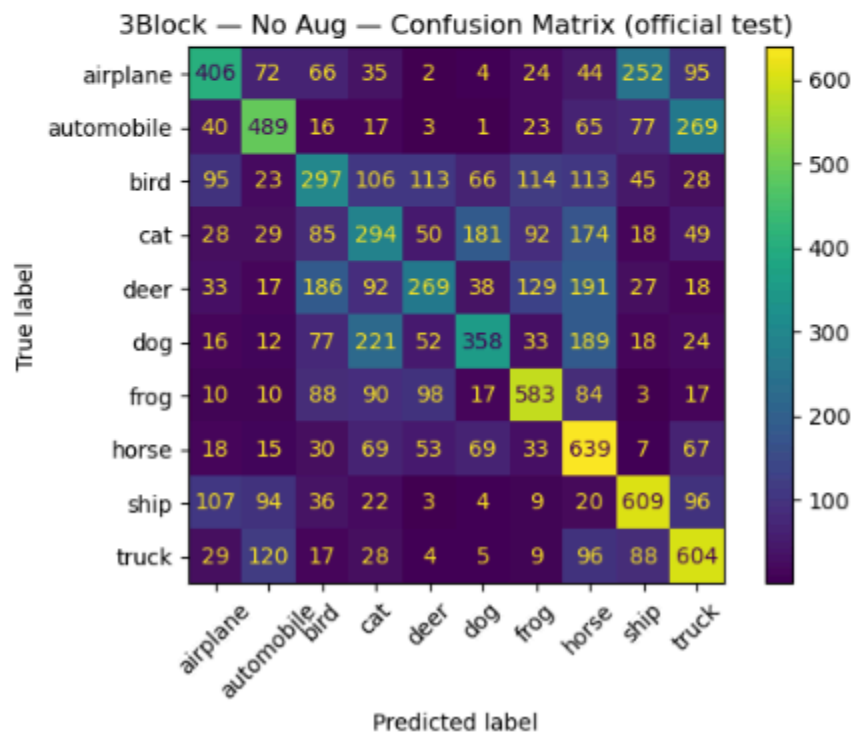
Effect on generalization.

- Capacity without depth gave a modest improvement in best val loss but no test-accuracy gain vs the baseline (both 0.361).
- With the same training recipe (short schedule, light regularization, no heavy aug), additional width did not translate to better out-of-sample accuracy.
- Interpretation: width reduced bias slightly, but without stronger invariances (augmentation) or more depth to build hierarchical features, the extra channels did not capture new generalizable patterns. Early stopping kept variance in check (no large train–val gap), but Depth (3-Block) remained the more effective axis for generalization in this regime.

3.5 Training Strategy: Augmentation / Optimizer

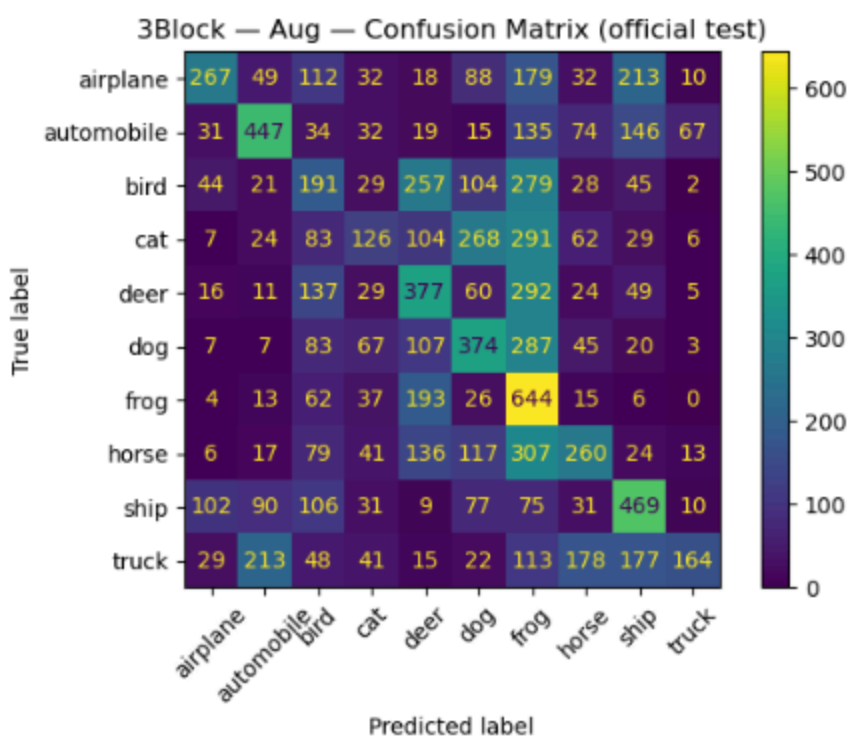
3Block — No Aug — test acc: 0.455

	precision	recall	f1-score	support
airplane	0.519	0.406	0.456	1000
automobile	0.555	0.489	0.520	1000
bird	0.331	0.297	0.313	1000
cat	0.302	0.294	0.298	1000
deer	0.416	0.269	0.327	1000
dog	0.482	0.358	0.411	1000
frog	0.556	0.583	0.569	1000
horse	0.396	0.639	0.489	1000
ship	0.532	0.609	0.568	1000
truck	0.477	0.604	0.533	1000
accuracy			0.455	10000
macro avg	0.456	0.455	0.448	10000
weighted avg	0.456	0.455	0.448	10000



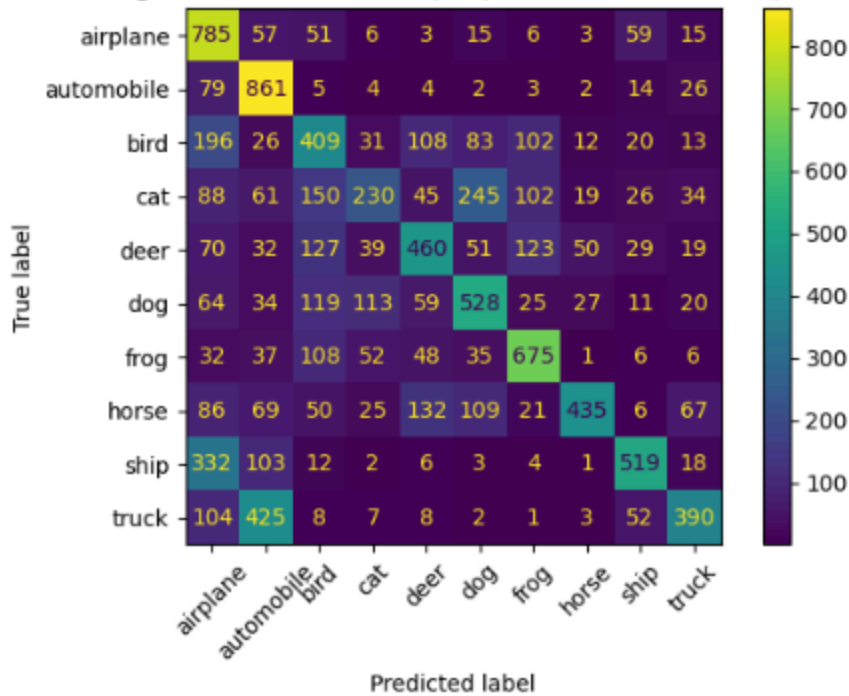
3Block — Aug — test acc: 0.332

	precision	recall	f1-score	support
airplane	0.520	0.267	0.353	1000
automobile	0.501	0.447	0.473	1000
bird	0.204	0.191	0.197	1000
cat	0.271	0.126	0.172	1000
deer	0.305	0.377	0.337	1000
dog	0.325	0.374	0.348	1000
frog	0.248	0.644	0.358	1000
horse	0.347	0.260	0.297	1000
ship	0.398	0.469	0.431	1000
truck	0.586	0.164	0.256	1000
accuracy			0.332	10000
macro avg	0.371	0.332	0.322	10000
weighted avg	0.371	0.332	0.322	10000



Predicted label				
3Block — Aug — SGD+momentum(0.9) — test acc: 0.529				
	precision	recall	f1-score	support
airplane	0.428	0.785	0.554	1000
automobile	0.505	0.861	0.637	1000
bird	0.394	0.409	0.401	1000
cat	0.452	0.230	0.305	1000
deer	0.527	0.460	0.491	1000
dog	0.492	0.528	0.509	1000
frog	0.636	0.675	0.655	1000
horse	0.787	0.435	0.560	1000
ship	0.699	0.519	0.596	1000
truck	0.641	0.390	0.485	1000
accuracy			0.529	10000
macro avg	0.556	0.529	0.519	10000
weighted avg	0.556	0.529	0.519	10000

3Block — Aug — SGD+momentum(0.9) — Confusion Matrix (official test)



```

=== 3.5: Training Strategy Comparison (NoAug vs Aug vs Aug+Momentum) ===
Config          Best Ep  Best Val Loss  Test Acc
-----
NoAug (SGD)      12        1.5612        0.455
Aug (SGD)        11        1.7067        0.332
Aug (SGD + momentum) 12        1.3020        0.529  *best val*  *best acc*

Saved: summary_3_5.csv

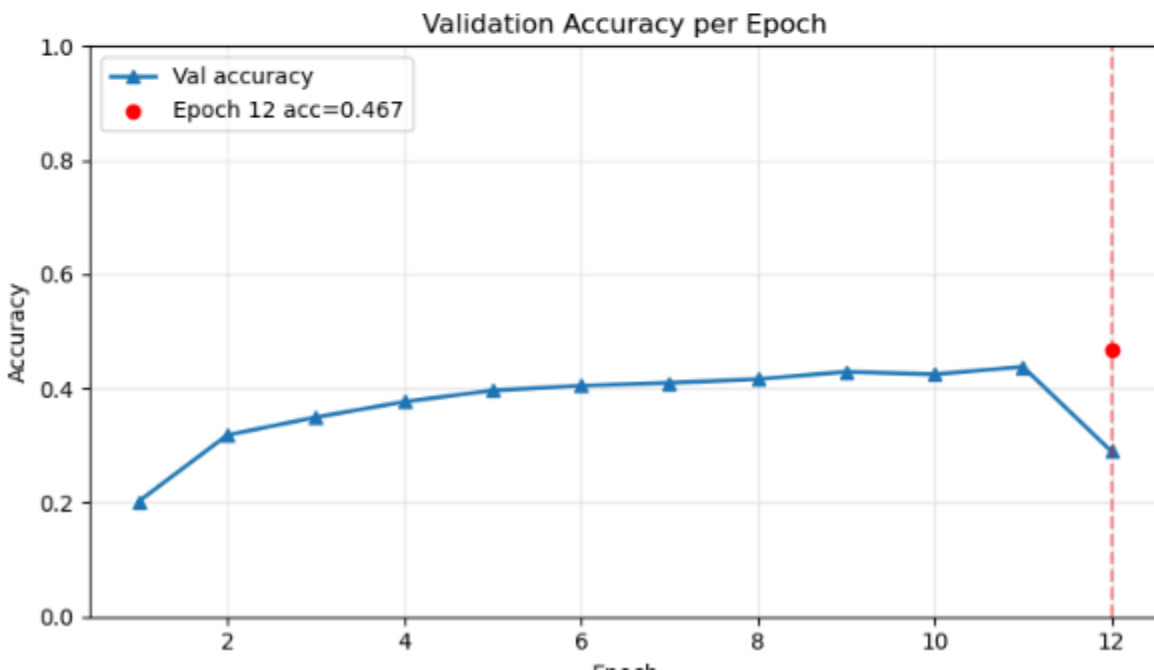
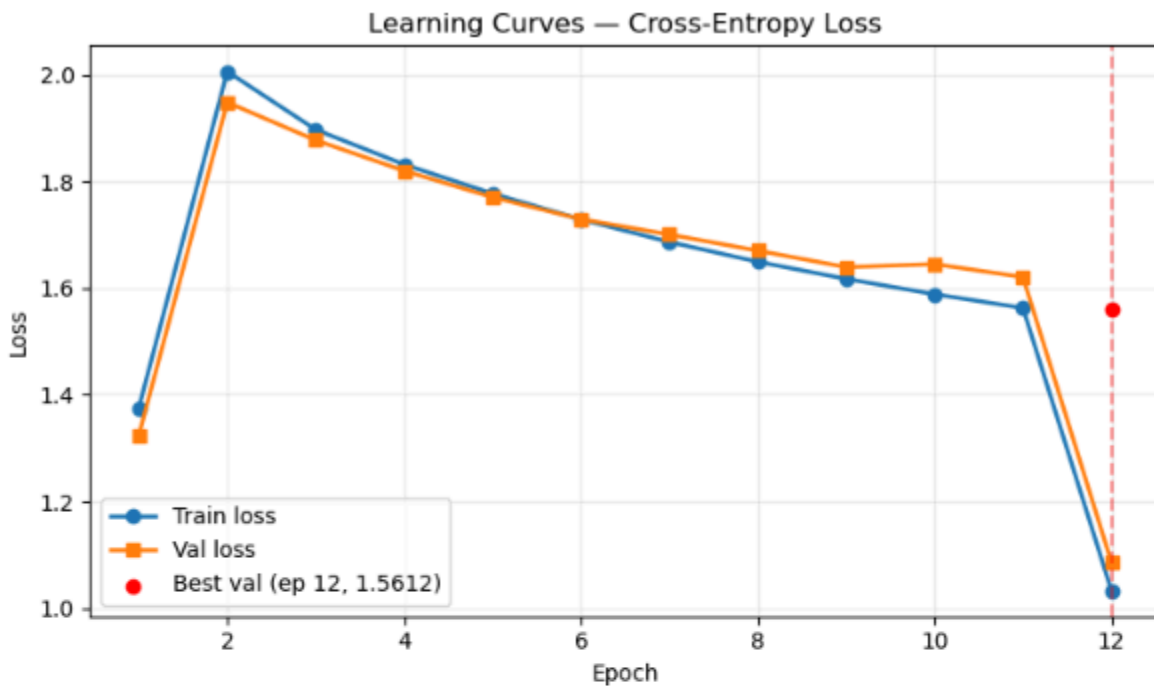
```

Setup. We fix the best architecture from 3.4 (3-Block + GAP) and vary only the

training strategy. Same split, batch size, epochs (12), weight decay ($1e-4$), early stopping on val loss (patience=3). Test is done on the official 10k set.

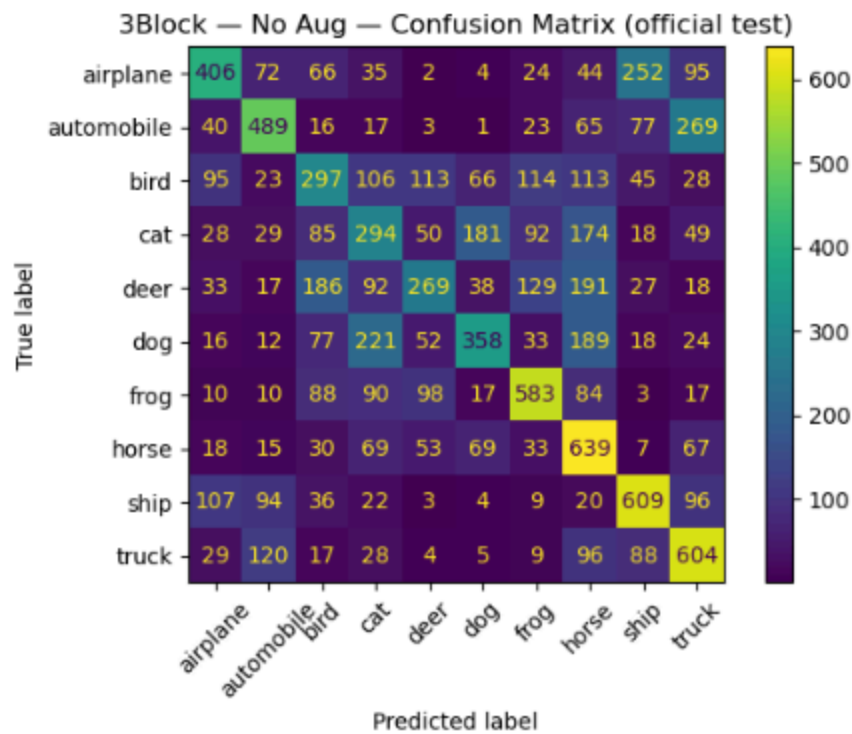
Configurations

1. No Aug (SGD) — Normalize only.

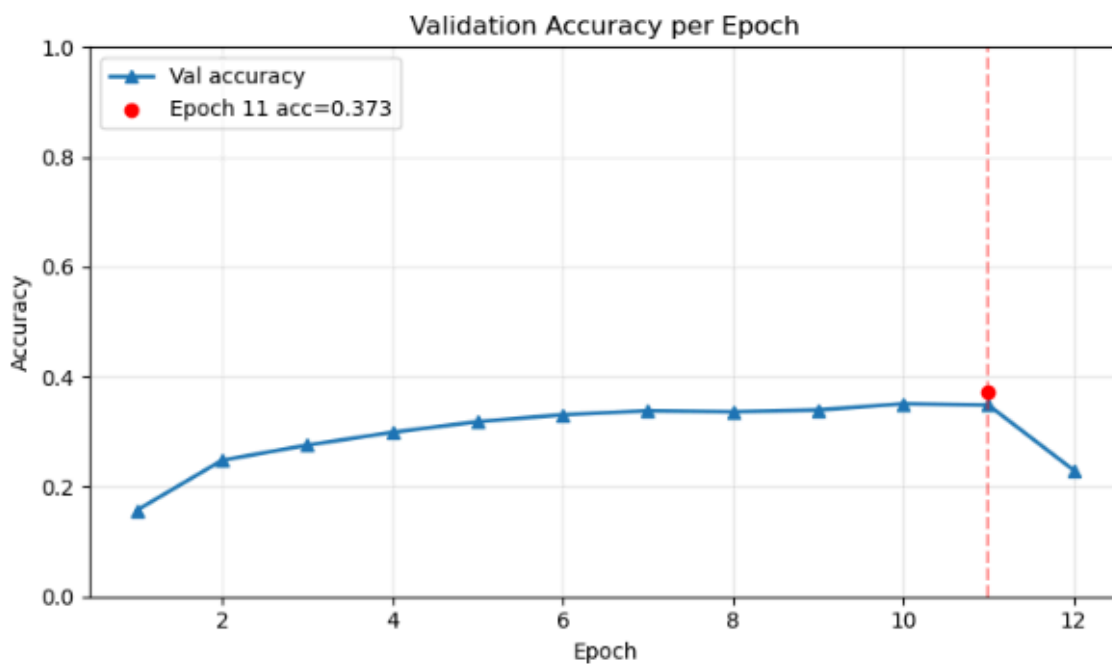
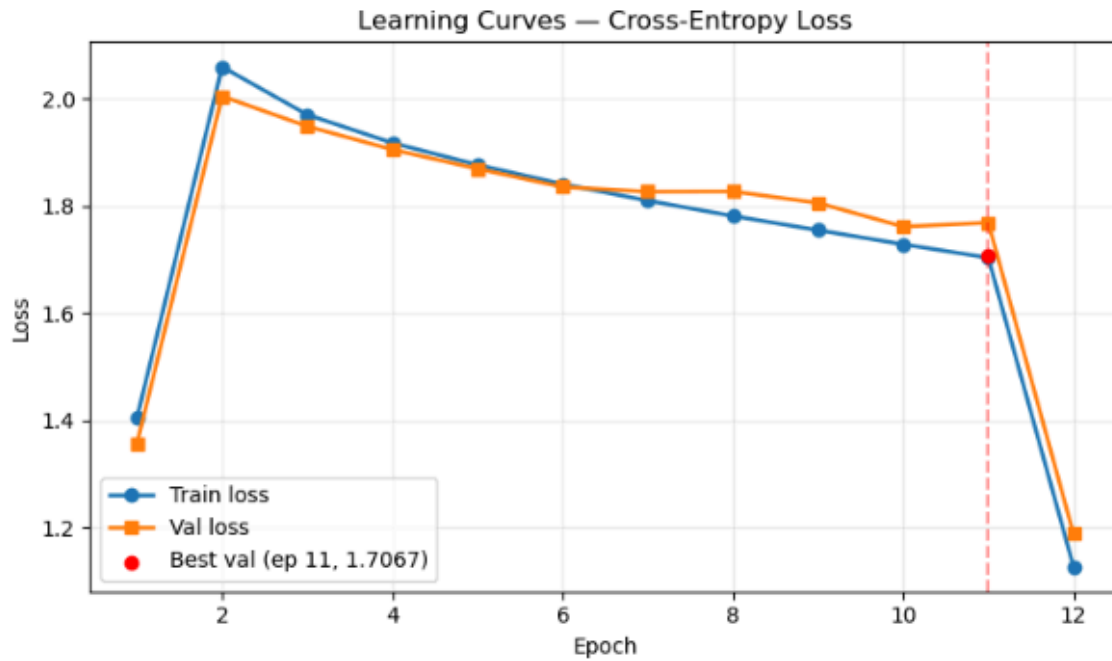


3Block — No Aug — test acc: 0.455

	precision	recall	f1-score	support
airplane	0.519	0.406	0.456	1000
automobile	0.555	0.489	0.520	1000
bird	0.331	0.297	0.313	1000
cat	0.302	0.294	0.298	1000
deer	0.416	0.269	0.327	1000
dog	0.482	0.358	0.411	1000
frog	0.556	0.583	0.569	1000
horse	0.396	0.639	0.489	1000
ship	0.532	0.609	0.568	1000
truck	0.477	0.604	0.533	1000
accuracy			0.455	10000
macro avg	0.456	0.455	0.448	10000
weighted avg	0.456	0.455	0.448	10000

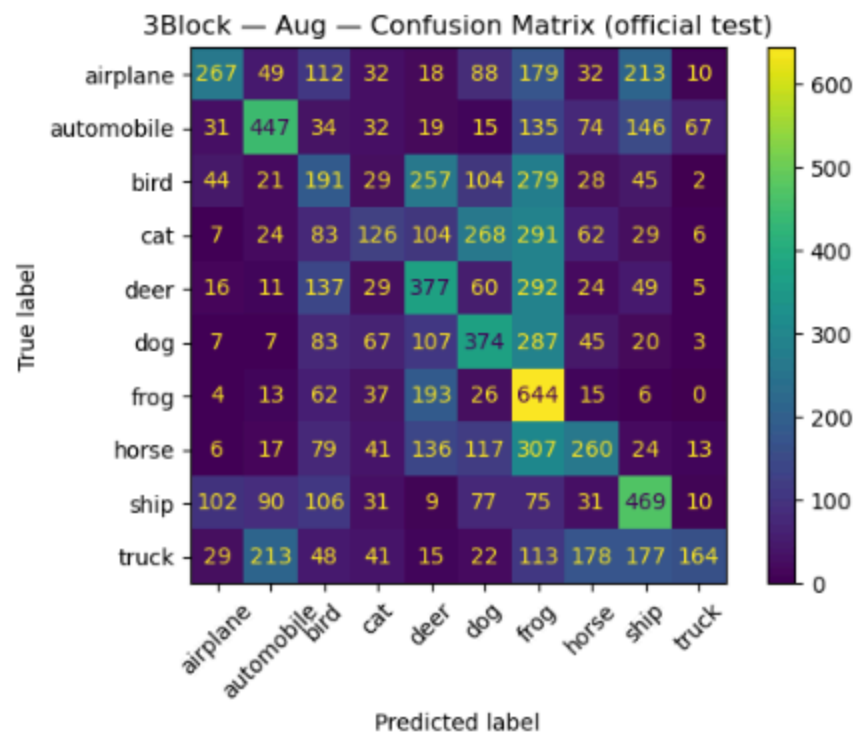


2. **Aug (SGD)** - RandomHorizontalFlip(p=0.5) + RandomCrop(32, padding=4) + Normalize.

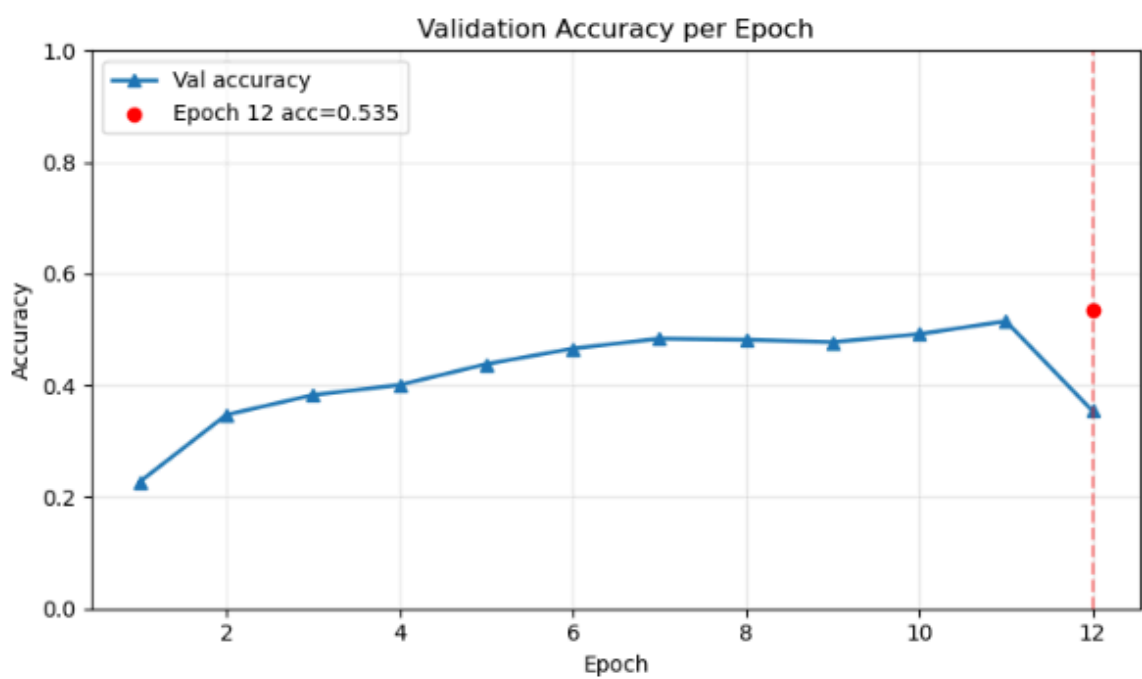
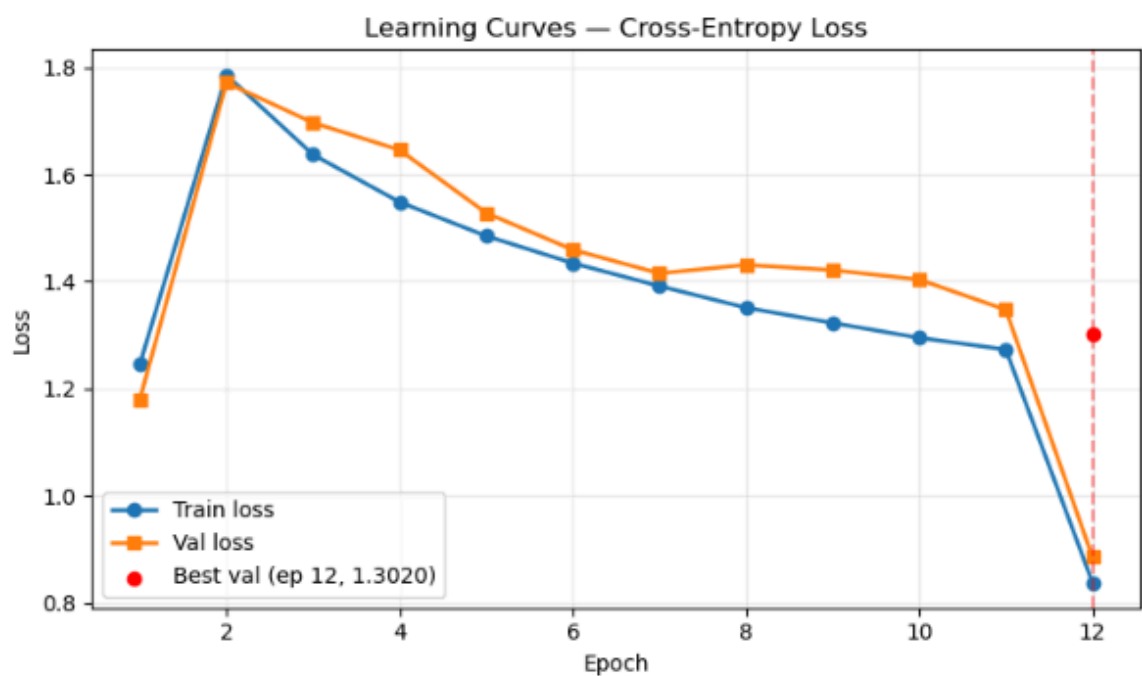


3Block — Aug — test acc: 0.332

	precision	recall	f1-score	support
airplane	0.520	0.267	0.353	1000
automobile	0.501	0.447	0.473	1000
bird	0.204	0.191	0.197	1000
cat	0.271	0.126	0.172	1000
deer	0.305	0.377	0.337	1000
dog	0.325	0.374	0.348	1000
frog	0.248	0.644	0.358	1000
horse	0.347	0.260	0.297	1000
ship	0.398	0.469	0.431	1000
truck	0.586	0.164	0.256	1000
accuracy			0.332	10000
macro avg	0.371	0.332	0.322	10000
weighted avg	0.371	0.332	0.322	10000



3. **Aug + Momentum** — Same aug; **SGD(momentum=0.9)**.



```

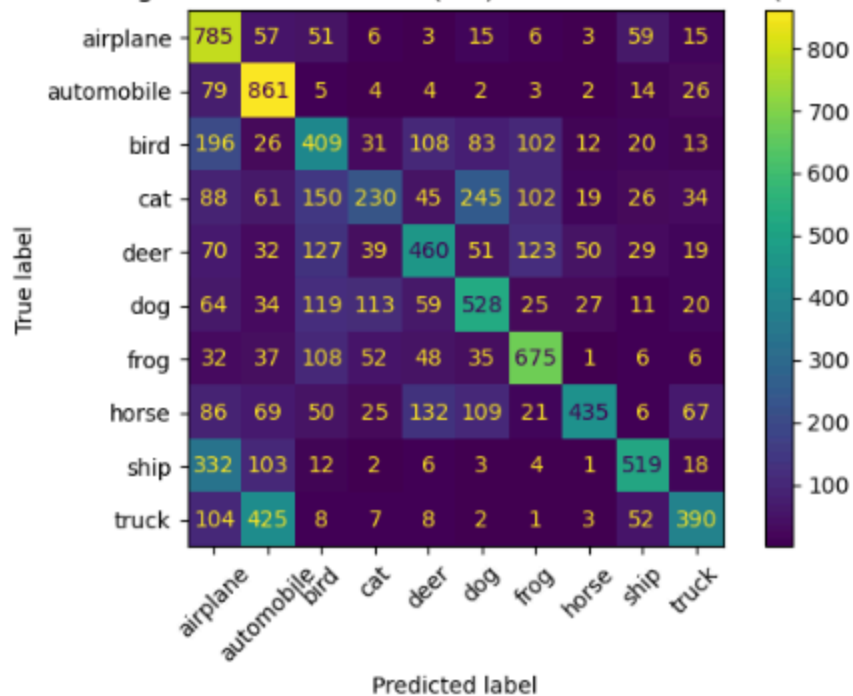
3Block — Aug — SGD+momentum(0.9) — test acc: 0.529
      precision    recall  f1-score   support

 airplane      0.428      0.785      0.554      1000
 automobile     0.505      0.861      0.637      1000
   bird        0.394      0.409      0.401      1000
    cat        0.452      0.230      0.305      1000
   deer       0.527      0.460      0.491      1000
    dog       0.492      0.528      0.509      1000
   frog       0.636      0.675      0.655      1000
  horse       0.787      0.435      0.560      1000
    ship       0.699      0.519      0.596      1000
   truck      0.641      0.390      0.485      1000

 accuracy              0.529      10000
 macro avg      0.556      0.529      0.519      10000
 weighted avg   0.556      0.529      0.519      10000

```

3Block — Aug — SGD+momentum(0.9) — Confusion Matrix (official test)



```

=== 3.5: Training Strategy Comparison (NoAug vs Aug vs Aug+Momentum) ===
Config          Best Ep  Best Val Loss  Test Acc
-----
NoAug (SGD)           12          1.5612      0.455
Aug (SGD)             11          1.7067      0.332
Aug (SGD + momentum)  12          1.3020      0.529  *best val*  *best acc*

```

Saved: summary_3_5.csv

Learning curves & model selection

- **No Aug:** smooth loss descent, best val loss 1.5612 at epoch 12; val-acc = 0.467.
- **Aug (SGD):** best val loss 1.7067 at epoch 11; val-acc \approx 0.373 (curves flatter—harder val).
- **Aug + Momentum(0.9):** best val loss 1.3020 at epoch 12; val-acc = 0.535, clearly best.

Final test metrics (official 10k)

- **No Aug (SGD):** Acc 0.455.
- **Aug (SGD):** Acc 0.332.
- **Aug + Momentum(0.9):** Acc 0.529 (best).

Per-class highlights (Aug+Mom): very strong recalls for automobile (0.861), airplane (0.785), ship (0.619); biggest remaining gaps on fine grained animals (e.g., cat R = 0.230, deer R = 0.460). Confusion matrices show the usual animal cross-confusions diminishing with augmentation + momentum but still present at 32×32.

Interpretation: effect on generalization

- **Augmentation** adds invariances (flip/crop), generally reduces overfitting and improves test accuracy. In our short schedule, plain Aug (SGD) underperformed—likely due to optimization dynamics: augmented samples are harder; without momentum the optimizer made slower progress and settled in a worse valley.
- **Momentum=0.9** with the same aug substantially improved optimization (smoother, faster traversal), yielding the lowest val loss and highest test accuracy (+7.4 pp vs No Aug; + 19.7 pp vs plain Aug).
- **Takeaway:** with small models and limited epochs, optimizer choice matters as much as augmentation. Aug + momentum is clearly the best recipe here.

Conclusion. For this architecture and budget, flip+crop augmentation + SGD with momentum delivers the strongest generalization. Future gains likely from

longer schedules, cosine LR decay, stronger aug (Cutout/Mixup), or combining depth/width with tuned regularization.

3.6 Comparative Discussion

Traditional vs deep learning.

The logistic-regression baseline (pixels + standardization) scored 0.281 on the test set. The compact CNNs did better: 0.361 with 2 blocks, 0.423 with 3 blocks, and 0.529 when we added augmentation and momentum. CNNs win because they keep spatial structure: convolutions learn local edges and parts, and pooling/global average pooling adds some translation tolerance.

Logistic regression flattens the image, loses that structure, and struggles especially on animals, which we also see in the confusion matrices.

A simple linear model can be enough when features are already close to linearly separable (for example, engineered color/texture features) or when speed and simplicity matter more than accuracy.

What single change helped most within deep learning?

The biggest improvement came from the training recipe: adding flip+crop augmentation and SGD with momentum 0.9 to the same 3-block CNN raised test accuracy from 0.423 to 0.529 (+10.6 percentage points). Augmentation adds useful invariances, and momentum helps the optimizer make steady progress on the harder, augmented data. Increasing depth from 2 to 3 blocks also helped (from 0.361 to 0.423, +6.2 points), but the training change had the larger effect.

4. Discussion & Takeaways

Compact CNN vs Logistic Regression

The logistic-regression baseline (pixels with StandardScaler) achieved **0.281** test accuracy. The compact CNN baseline with two blocks achieved **0.361** (an improvement of **8.0 percentage points**). Our best configuration—the three-block CNN with data augmentation and momentum—reached **0.529** (an improvement of **24.8 percentage points** over logistic regression).

These gains come from the CNN's spatial inductive bias: shared local filters and pooling/global average pooling preserve image structure that is lost when pixels are flattened.

Depth vs Width (Architecture)

Increasing depth from two to three blocks improved test accuracy from 0.361 to 0.423 (up 6.2 percentage points), indicating better mid-level feature learning.

Increasing width from 32/64 channels to 48/96 reduced the best validation loss slightly but did not improve test accuracy (it remained 0.361) within our short training schedule. The extra channels did not translate into more generalizable patterns without stronger training.

Training Strategy (Most Impactful Change)

Holding the three-block architecture fixed:

- No augmentation with SGD: 0.455
- Flip-and-crop augmentation with SGD: 0.332 (under-optimized)
- Flip-and-crop augmentation with SGD and momentum 0.9: 0.529
(improvements of 7.4 points over no augmentation, 19.7 points over augmentation without momentum, and 10.6 points over the three-block model without this training recipe)

Augmentation provides useful invariances, and momentum helps the optimizer make steady progress on the harder, augmented data. Together they produced the largest improvement in generalization.

Per-Class Behaviour (official test)

Vehicle categories are consistently easier. With augmentation and momentum, recall is roughly 0.861 for automobiles, 0.785 for airplanes, and 0.619 for ships. Fine grained animal categories remain hardest: cats around 0.230, deer around 0.460, and birds around 0.490, reflecting the low image resolution and high visual similarity among these classes.

Limitations and Future Work

All experiments were CPU-bound with a 12-epoch schedule and minimal hyperparameter search.

Next steps: train longer; use cosine or one-cycle learning-rate schedules; adopt stronger augmentation (Cutout, Mixup); combine added depth with moderate width and tuned regularization; explore lightweight residual blocks and label smoothing.

5. Reproducibility Notes

- Fixed RNG seeds (Python/NumPy/PyTorch), deterministic DataLoader (num_workers=0, fixed generator).
- Saved split indices JSON and best-epoch checkpoints.
- Evaluated on official tests only once per model.