

The CIFAR-10 Challenge: A Tale of Two Models

Comparing Classic Machine Learning and
Modern Deep Learning for Image Recognition



The Challenge: Recognizing Objects in a Low-Resolution World

The Dataset

CIFAR-10 is a classic benchmark for image classification, consisting of 60,000 labeled images.



The Difficulty

The images are tiny (32x32 pixels), making it hard to distinguish fine details.

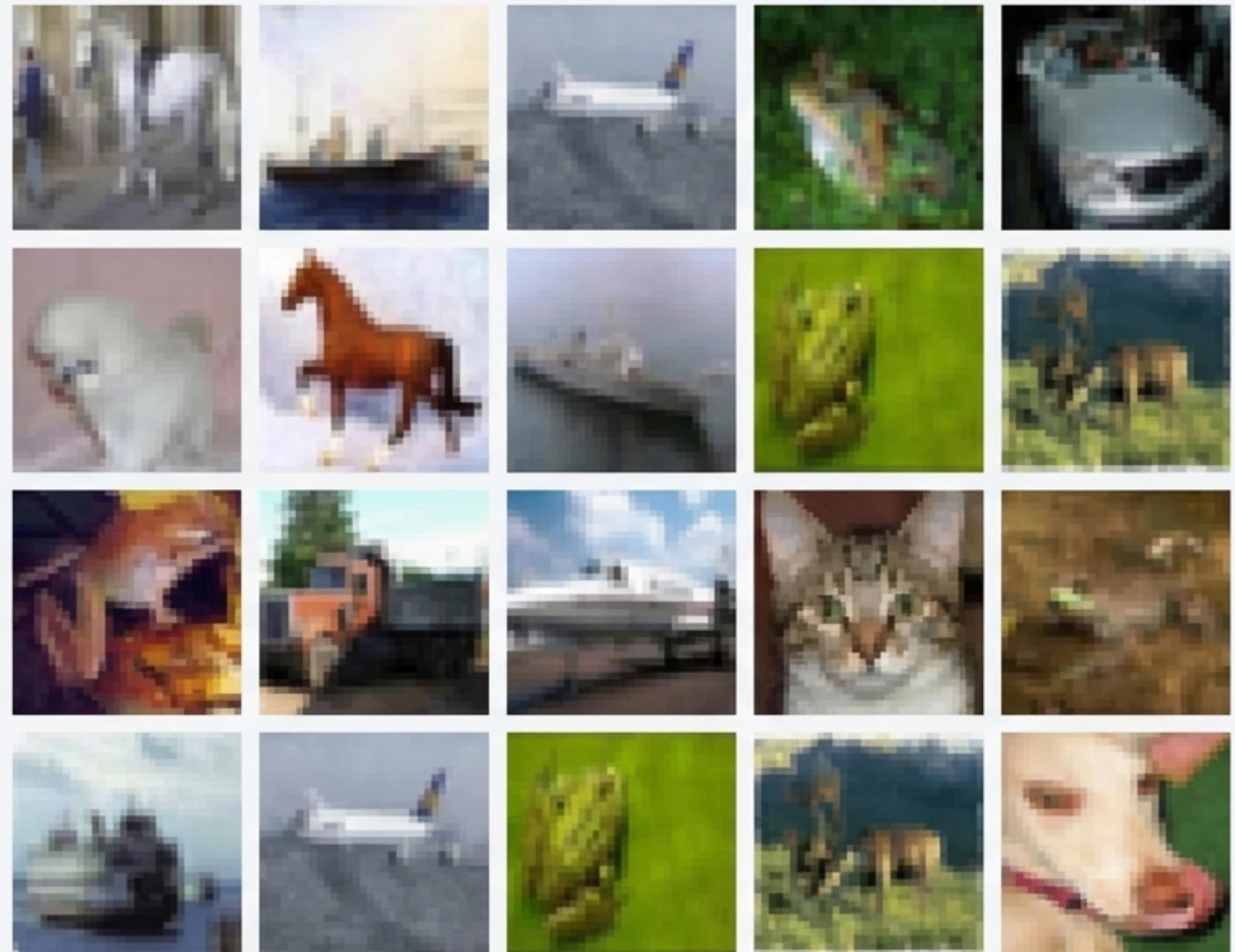
The Diversity

The dataset contains 10 balanced classes (5,000 training images each), ranging from animals ('frog', 'cat', 'dog') to vehicles ('airplane', 'ship', 'truck').

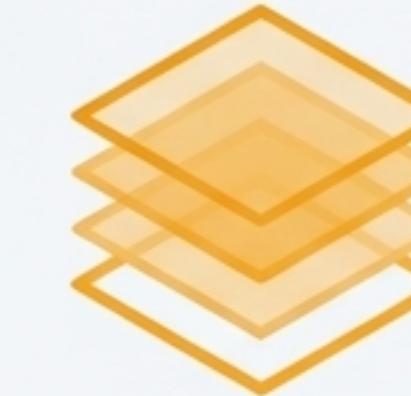


The Ambiguity

Many classes share visual characteristics, creating significant classification challenges.

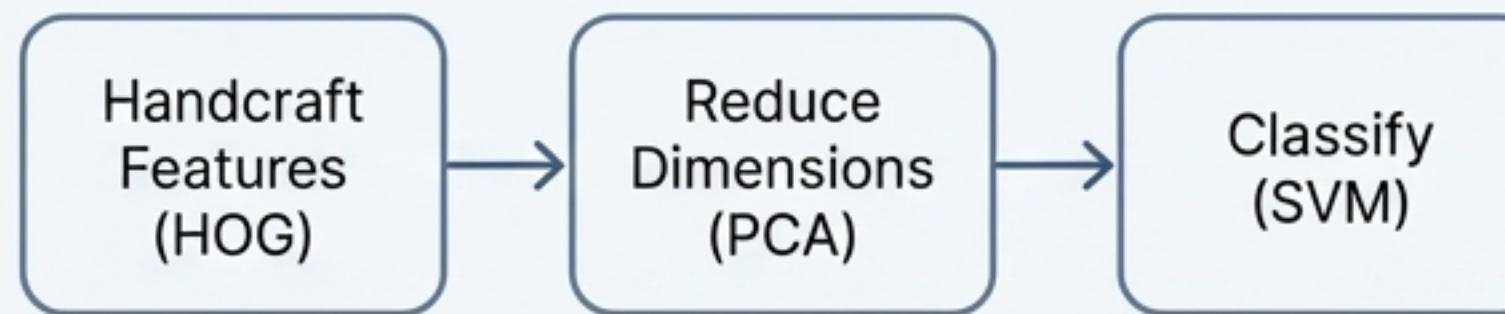


Meet the Contenders: Two Philosophies for Image Recognition



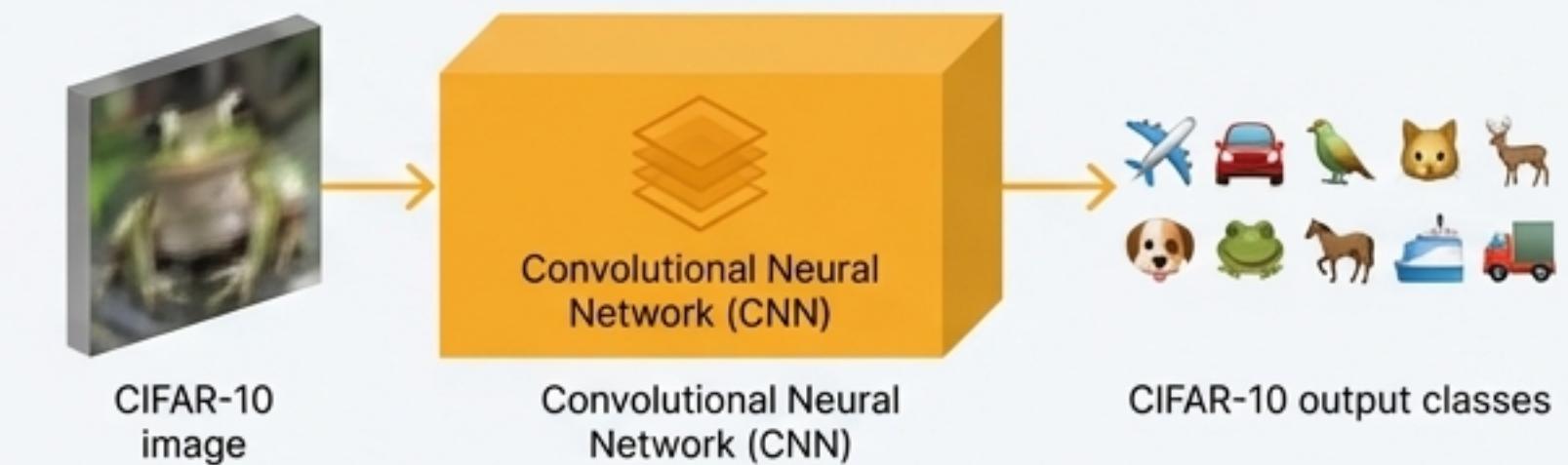
The Craftsman: SVM Pipeline

A meticulous, multi-step process reliant on human expertise to design and extract features.



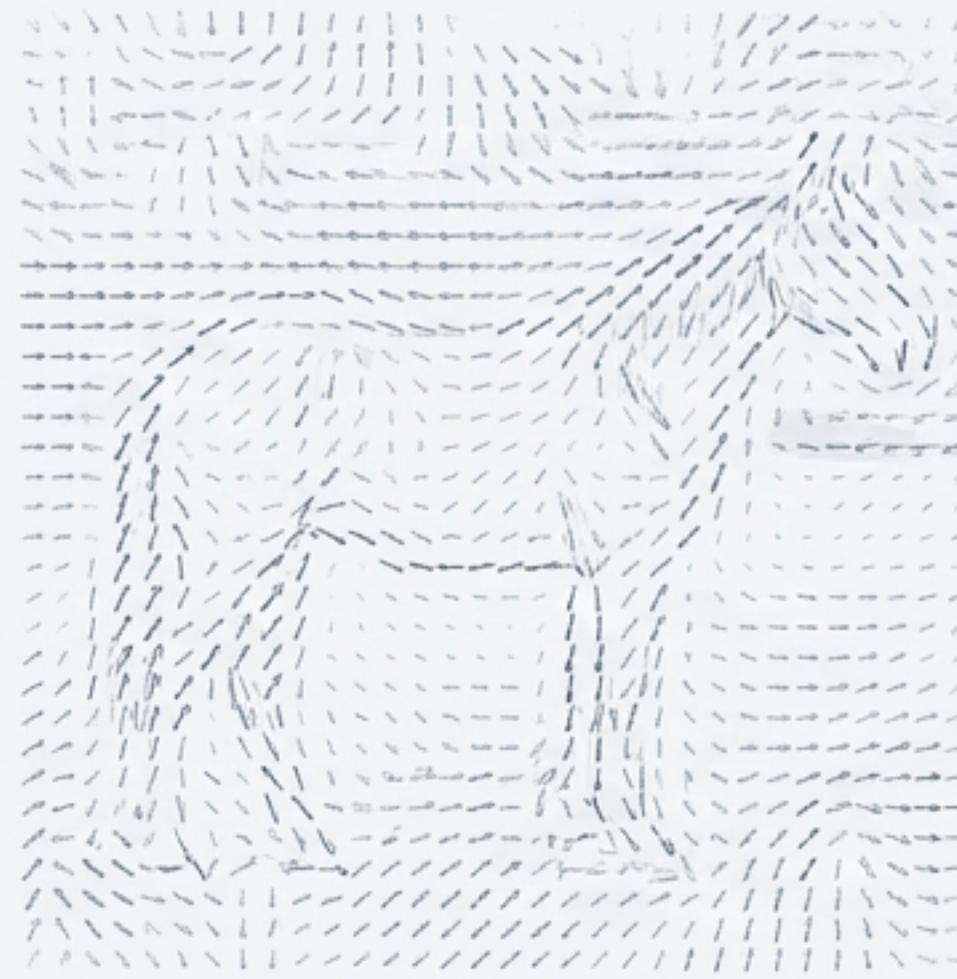
The Learner: CNN

An integrated, end-to-end approach where the model learns the optimal features directly from the raw pixel data.



Step 1: Translating Pixels into a World of Shapes with HOG

Histogram of Oriented Gradients (HOG) is a feature descriptor that captures an object's structure by counting gradient orientations in an image. It translates complex visuals into a numerical vector representing edges and shapes. To reduce computational load, RGB images are first converted to grayscale.



CIFAR-10 image

HOG Visualization



324-Dimension Feature Vector

Original Image: 32x32 pixels

HOG Feature Vector: **324 features** per image

Training Set Matrix: **(50000, 324)**

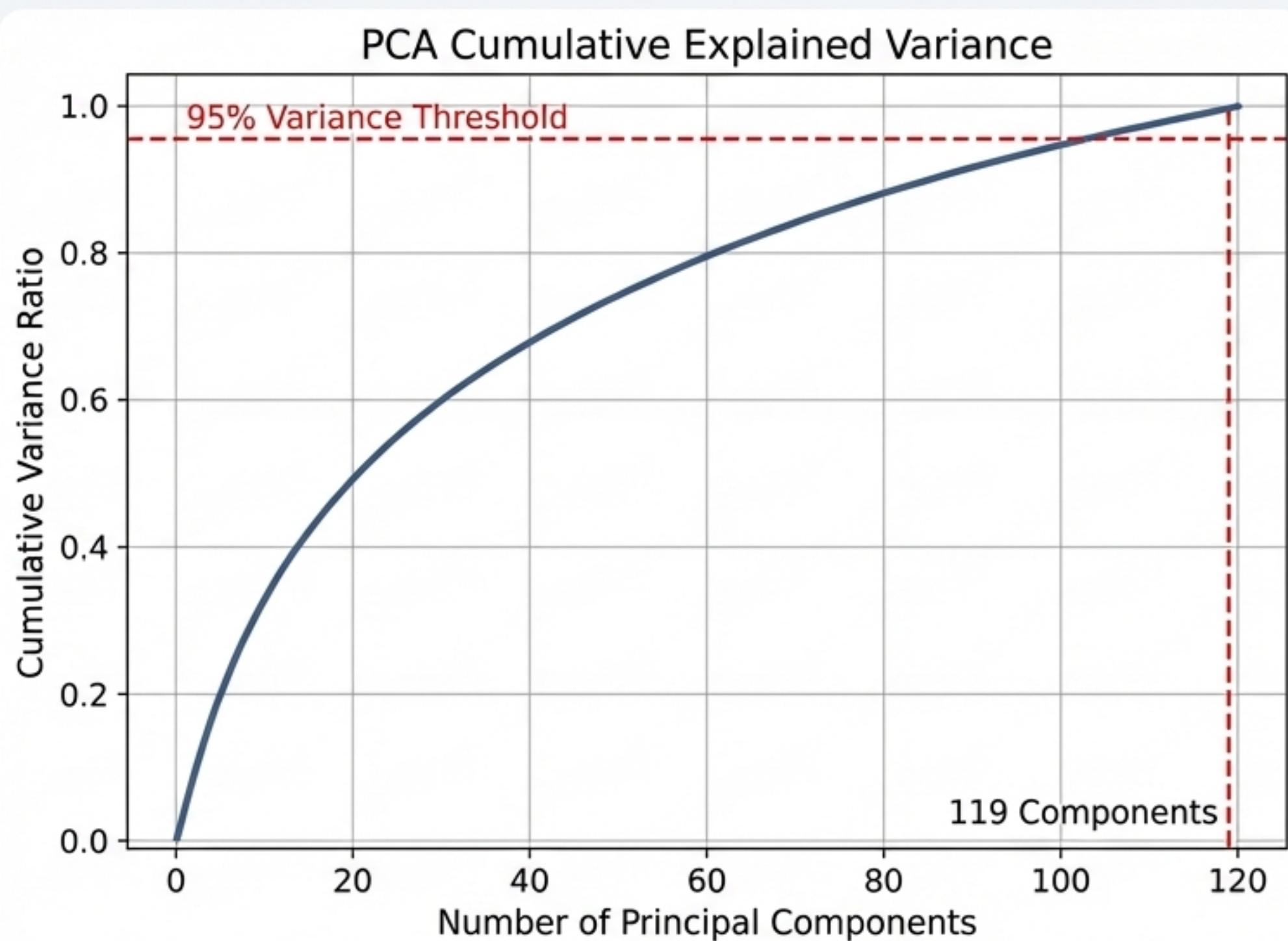
Step 2: Taming 324 Dimensions with Principal Component Analysis (PCA)

The 324-dimensional HOG feature space is too large and noisy for an SVM to handle efficiently. PCA is a linear technique used to reduce dimensionality while preserving as much of the data's variance (information) as possible.

To retain 95% of variance...

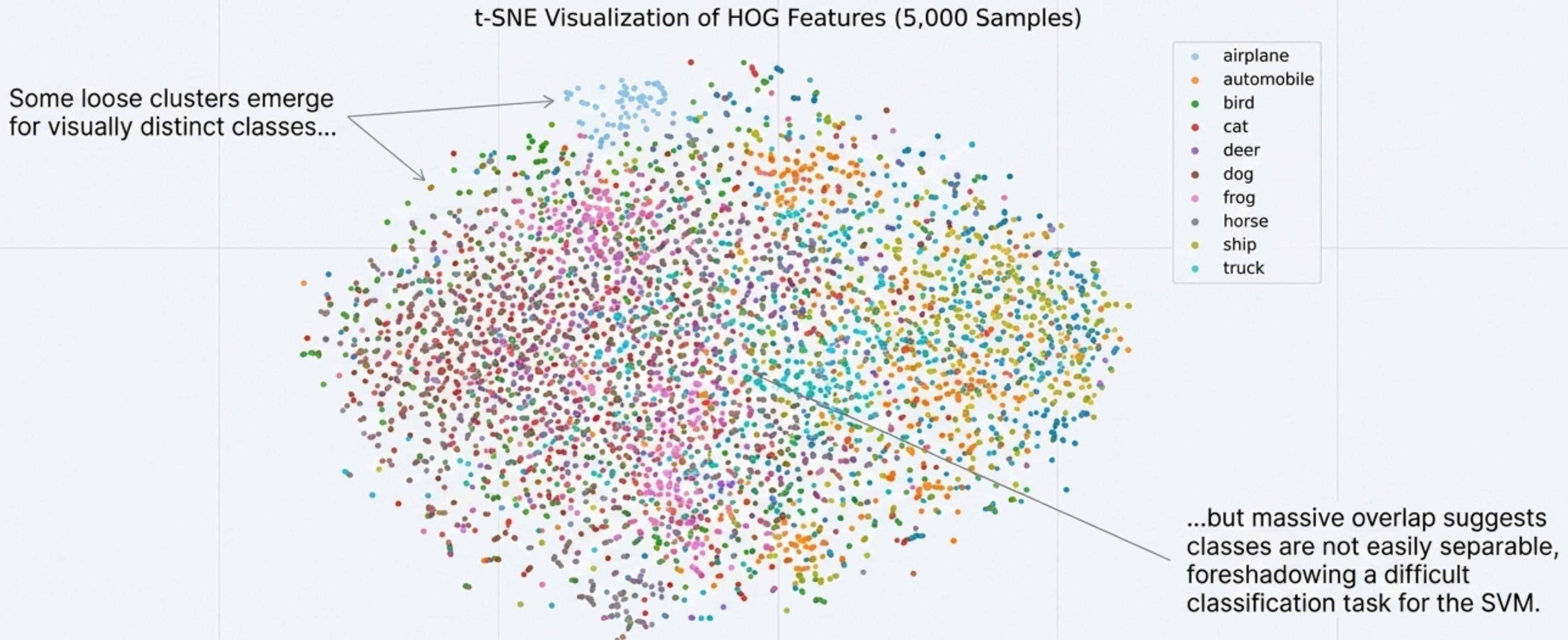
we must reduce 324 dimensions to **119** principal components.

This is a **36.7%** compression, making the data manageable for the SVM.



A Glimpse into the Feature Space: Can We See the Classes?

t-SNE is a non-linear technique that visualizes high-dimensional data in 2D, revealing underlying cluster structures. We applied it to 5,000 samples to map out the 119-dimensional feature space.

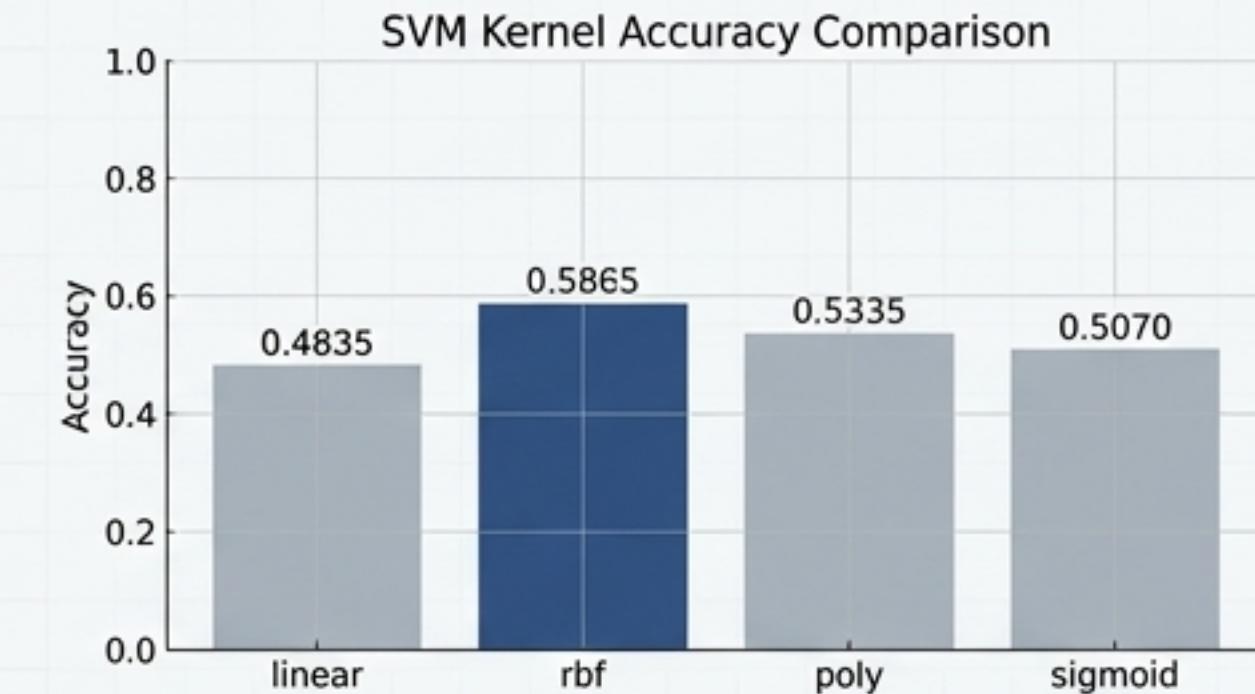


Step 3: Forging the SVM Decision Engine

Part 1: Choosing the Right Tool (Kernel Selection)

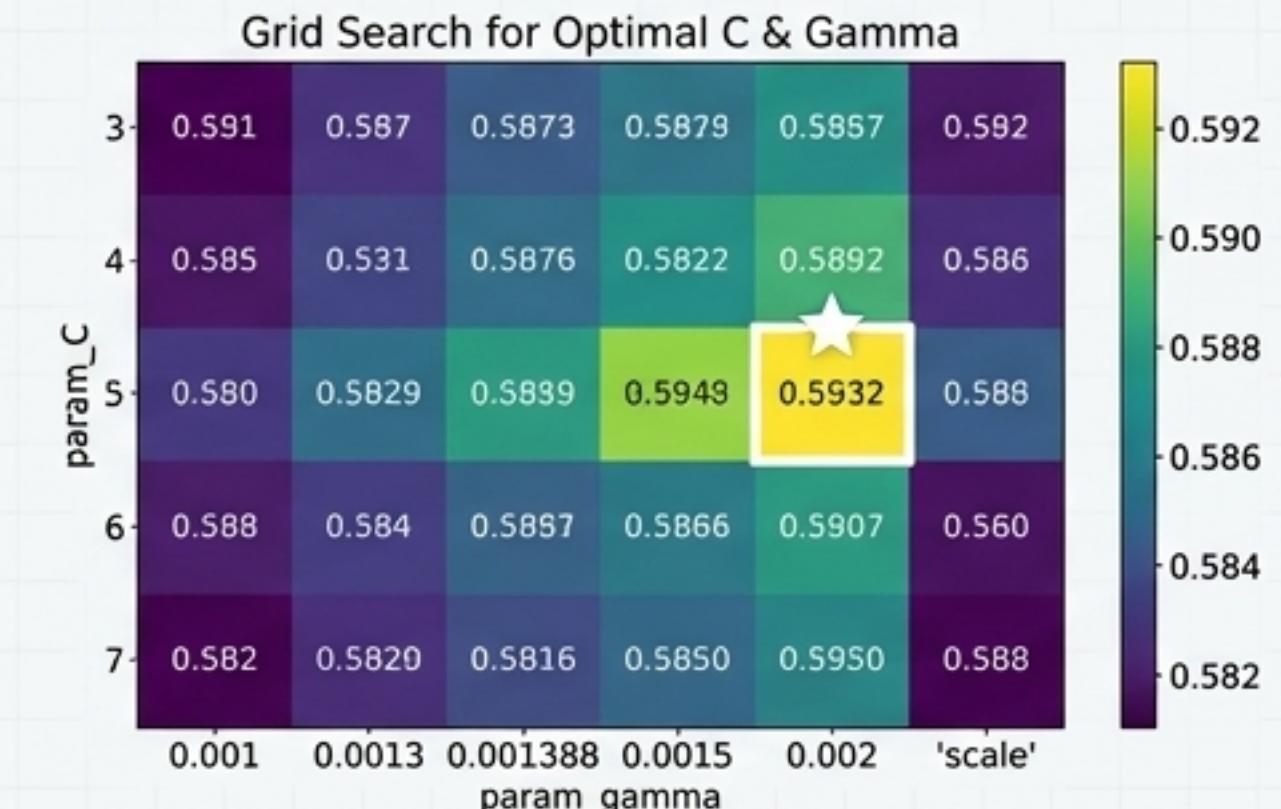
Four SVM kernels were tested: linear, rbf, poly, and sigmoid.

The **RBF (Radial Basis Function) kernel** was the clear winner, achieving the highest accuracy (**58.65%**) with a fast training time (**8.71 seconds**). The linear kernel was extremely slow (263 seconds) and less accurate.

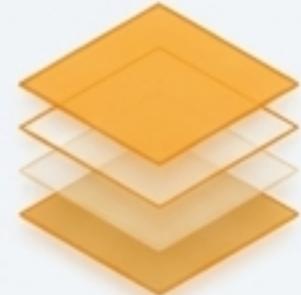


Part 2: Fine-Tuning for Peak Performance (Grid Search)

A two-stage grid search was performed to find the optimal hyperparameters ('C' and 'gamma') for the RBF kernel. The optimal parameters were found to be **C=4** and **gamma=0.002**, which pushed the cross-validation accuracy to **59.32%**.

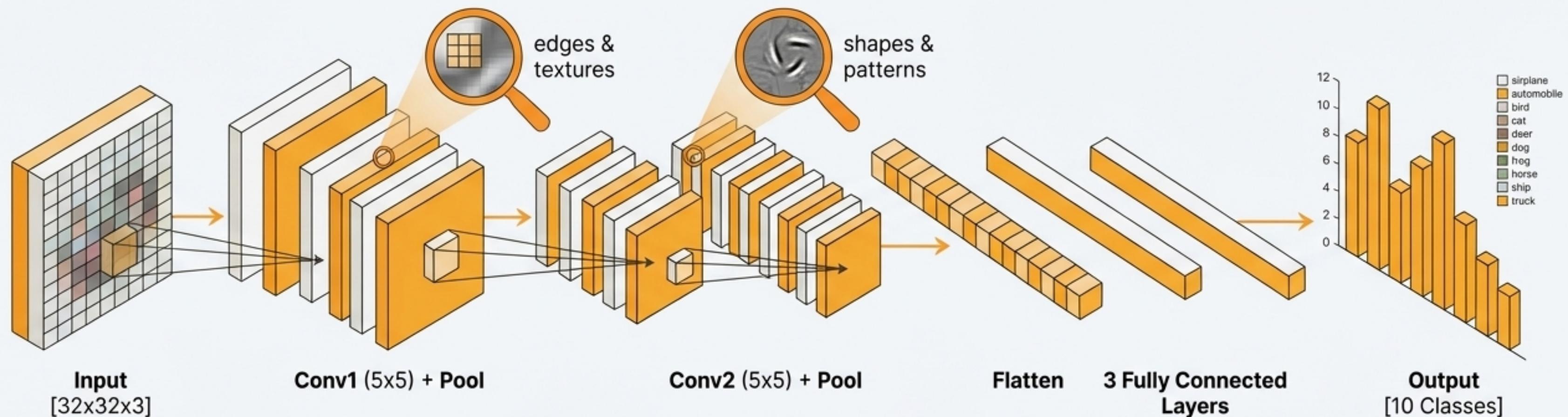


The Learner's Blueprint: A Self-Taught Feature Factory



The Learner

Unlike the SVM's manual process, a **Convolutional Neural Network (CNN)** learns a hierarchy of features directly from image pixels. Early layers learn simple features like edges and colors, while deeper layers **combine these to recognize complex shapes and object parts.**



Total Trainable Parameters: **62,006**
This simple architecture efficiently learns to classify images without any pre-defined feature extractors.

The Learning Process: Forging Intelligence Through Iteration

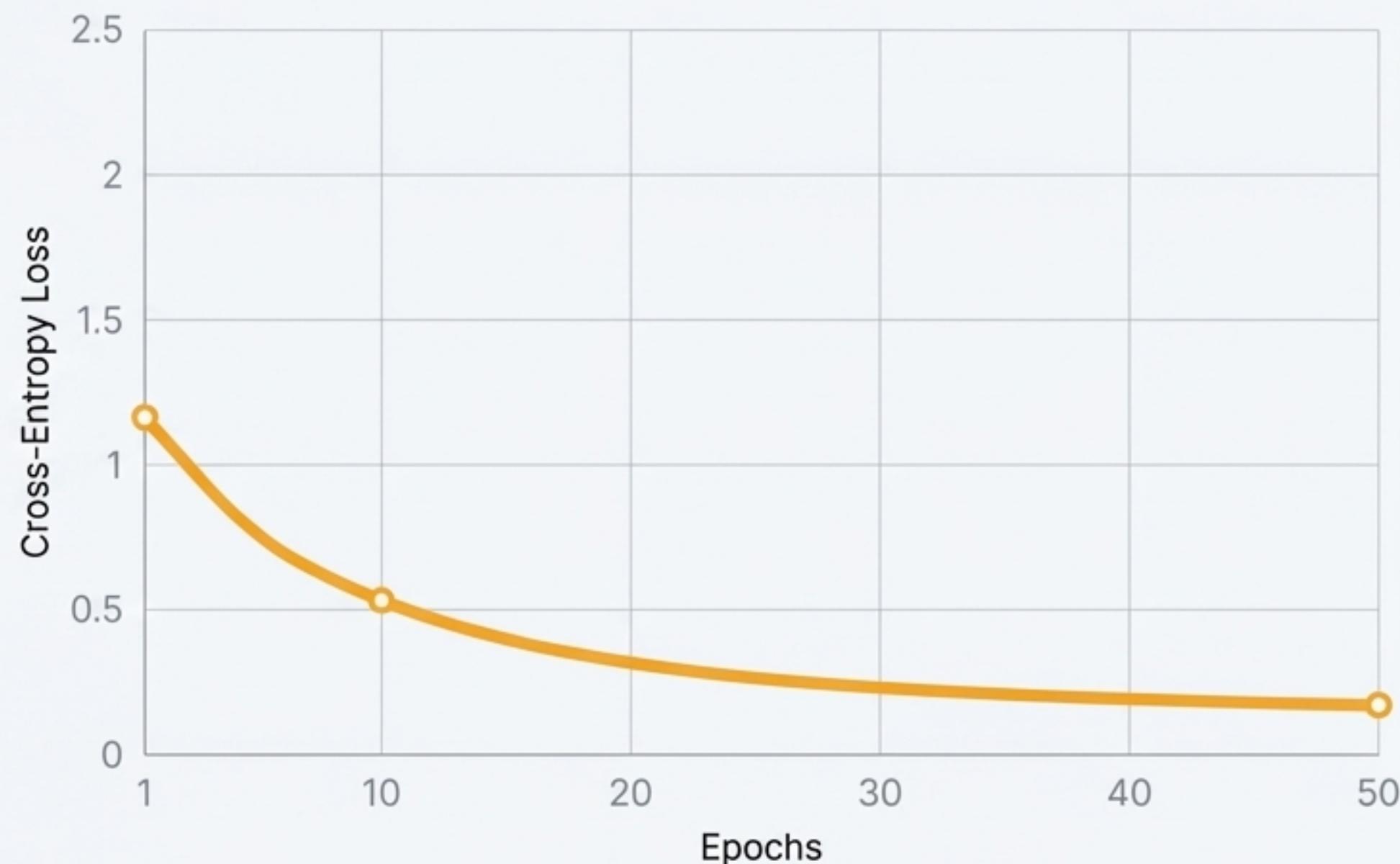
Training Loop Explained

The model was trained on the 50,000 CIFAR-10 images for **50 epochs** (50 full passes through the dataset). In each pass, the model makes predictions, calculates its error using a **Cross-Entropy Loss** function, and adjusts its internal parameters using an **SGD optimizer** to reduce that error.

Visualizing Improvement

The chart shows the training loss steadily decreasing over 50 epochs. The initial high loss indicates the model starts with random guesses. The rapid decline shows it is quickly learning the underlying patterns in the data, becoming more accurate with each iteration.

CNN Training Loss Over 50 Epochs



The Final Scorecard: A Clear Victor Emerges



The Craftsman (SVM)

59.5%

Final Test Accuracy

Achieved after extensive feature engineering, PCA, kernel selection, and hyperparameter tuning.



The Learner (CNN)

62.1%

Final Test Accuracy

Achieved through end-to-end learning directly from pixels, with no handcrafted features.

The Heat of Battle: Uncovering Confusion Between Classes

A confusion matrix shows what classes a model confuses with one another. The diagonal represents correct classifications; off-diagonal cells represent errors.

SVM Confusion Matrix

	airplane	automobile	bird	cat	deer	dog	frog	horse	ship	truck	
airplane	134	7	6	2	5	3	1	3	35	4	
automobile	5	152	2	2	5	1	5	0	16	12	
bird	16	4	93	15	17	23	17	5	7	3	
cat	7	6	11	60	16	58	17	16	3	6	
deer	10	4	17	19	101	7	16	18	4	6	
dog	0	3	15	34	8	113	6	15	2	4	
frog	3	3	16	19	7	9	137	2	2	2	
horse	3	3	11	15	19	15	5	117	0	10	
ship	29	9	5	6	2	5	4	5	128	7	
truck	5	15	2	9	6	0	2	6	14	141	

Achilles' Heel:
Out of 200 true
cats, **58 (29%)**
were misclassified
as dogs.

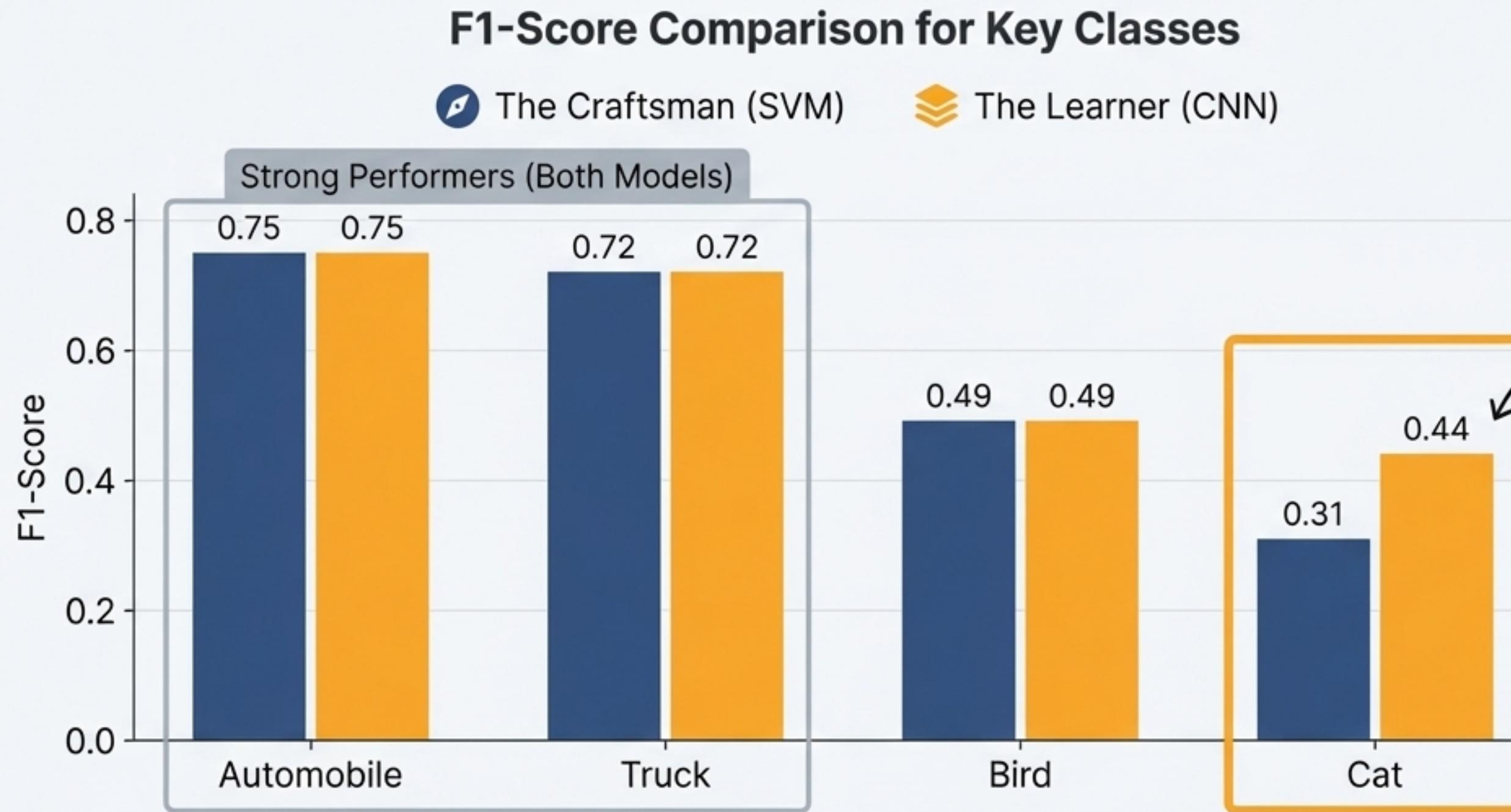
CNN Confusion Matrix

	airplane	automobile	bird	cat	deer	dog	frog	horse	ship	truck	
airplane	629	22	72	37	50	24	10	18	79	59	
automobile	17	740	14	26	12	14	7	13	29	128	
bird	57	11	457	111	127	103	64	40	13	17	
cat	9	21	48	469	81	217	75	45	10	22	
deer	19	8	66	65	619	76	60	75	8	6	
dog	5	6	60	229	68	530	28	55	7	12	
frog	4	7	46	36	66	55	699	11	14	11	
horse	11	6	53	58	108	91	13	633	2	20	
ship	88	69	22	30	19	17	8	5	683	58	
truck	32	87	13	26	19	18	13	21	23	747	

Visible Improvement:
While still the
largest source of
error, the overall
confusion
between
challenging
classes like cat,
dog, and bird is
significantly
reduced.

A Tale of Two Tiers: How Performance Varies by Class

The F1-score is a balanced measure of a model's precision and recall. A direct comparison highlights where each model excels and struggles.



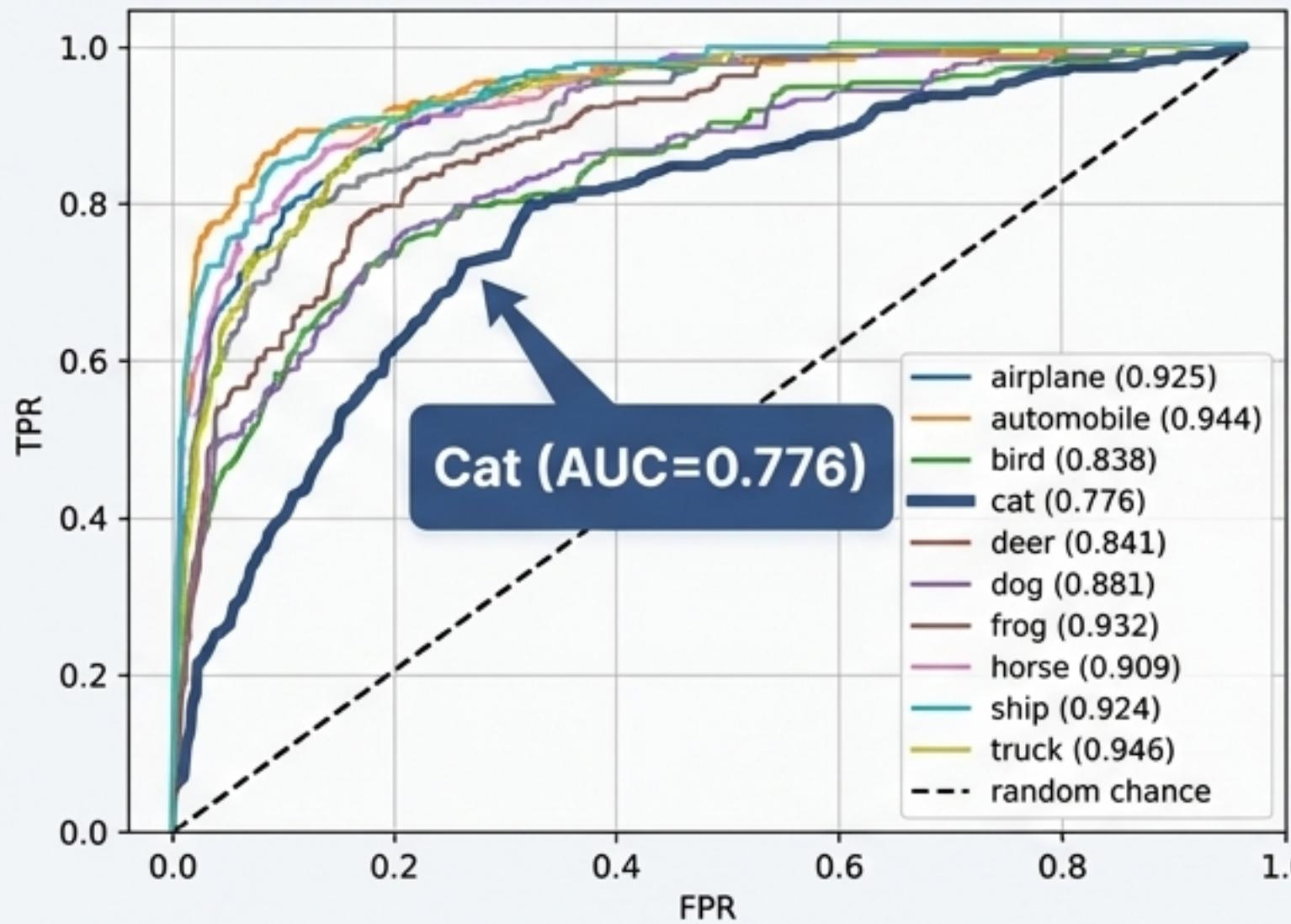
The Deciding Factor:
The CNN's victory is cemented by a **42% relative improvement** on the most difficult class.

The Diagnostic View: Measuring True Discriminative Power

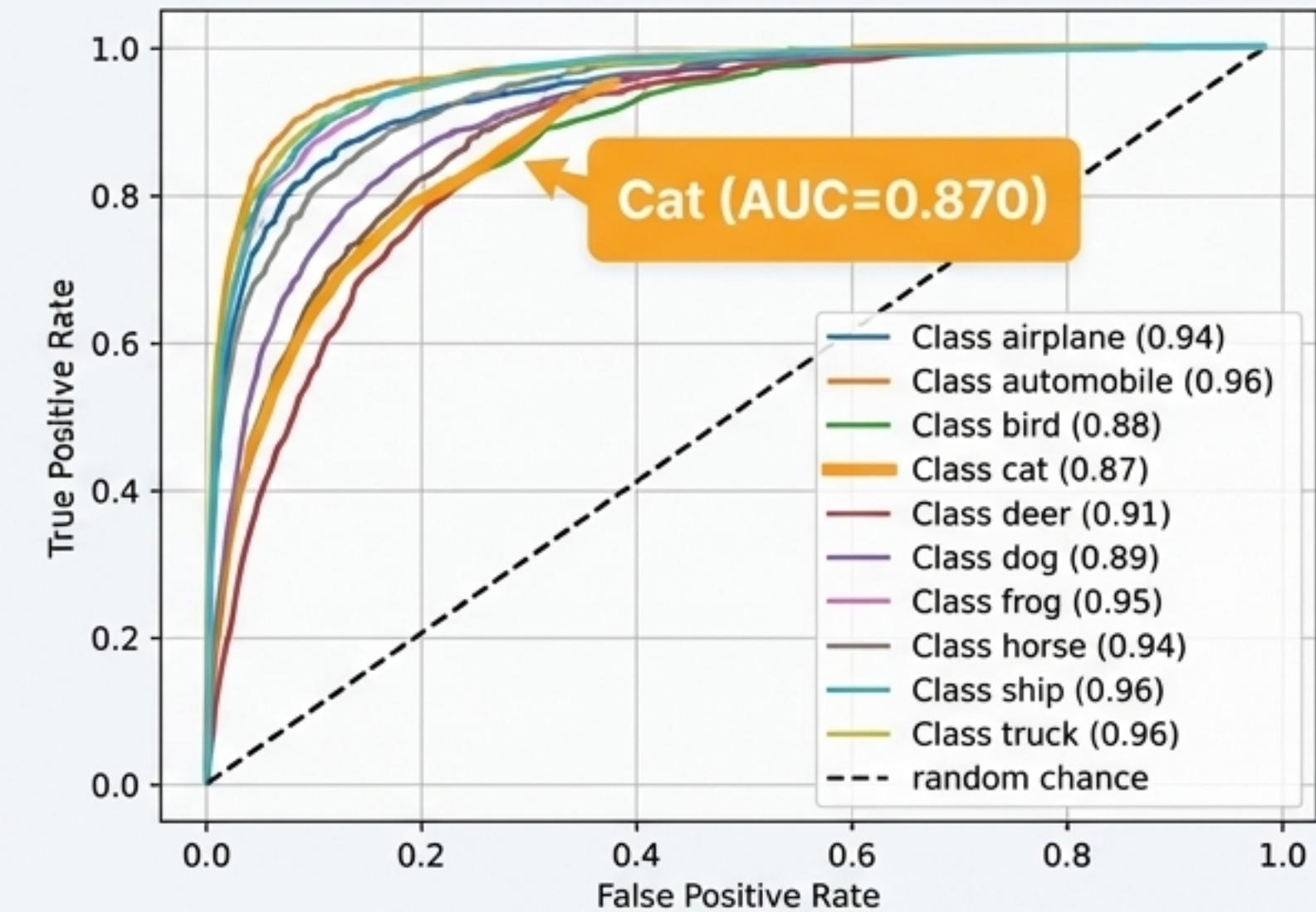
The Area Under the ROC Curve (AUC) measures a model's ability to distinguish between classes. An AUC of 1.0 is a perfect classifier, while 0.5 is random chance.



SVM ROC Curves (One-vs-Rest)



CNN ROC Curves (One-vs-Rest)



The Definitive Proof: The significant gap in AUC for the 'Cat' class (0.776 vs. 0.870) proves the CNN has a much stronger, more reliable ability to separate cats from all other classes.

The Verdict: Why the Learner Prevailed



The Craftsman (SVM)

The Craftsman's Limitation

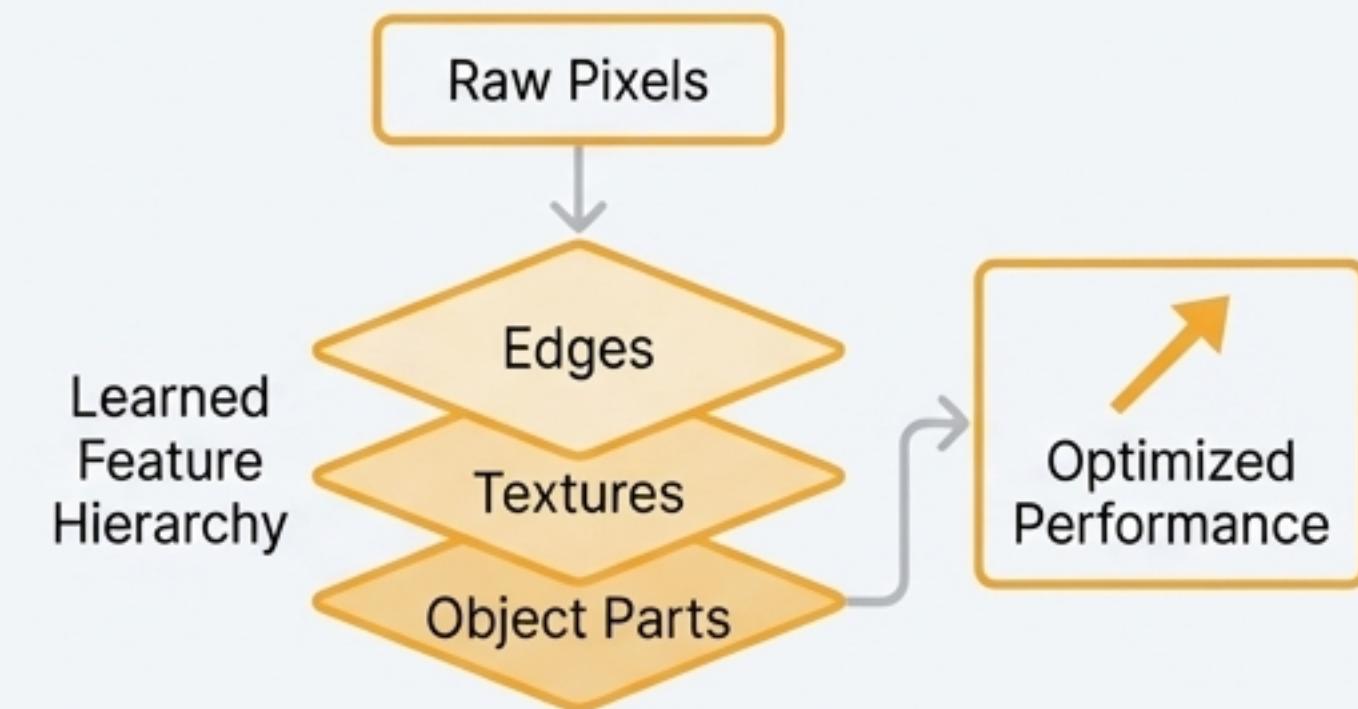
The SVM pipeline's performance was capped by its handcrafted features. HOG is powerful but general-purpose, failing to capture subtle cues in low-resolution images. The model can only be as good as the features it is given.



The Learner (CNN)

The Learner's Advantage

The CNN's victory stems from its ability to perform task-specific feature learning. The network learned a deep hierarchy of features specifically optimized to solve the CIFAR-10 problem.



The Future Belongs to Models That Learn

For complex perception tasks like image recognition, **the era of relying on handcrafted features is over**. The superior performance and efficiency of end-to-end deep learning models have established a new paradigm.

The contest between the Craftsman and the Learner demonstrates a **fundamental shift**: success is no longer defined by the cleverness of human-engineered features, but by the ability to create architectures that can effectively learn their own optimal representations directly from data.

