

# Error-Aware Baselines for Generative Visual Recognition: A CIFAR-10 Study from Handcrafted Features to Transfer Learning

Anonymous CVPR submission

Paper ID \*\*\*\*

## Abstract

001 *Generative image models are increasingly used as data*  
002 *sources for visual recognition, but their value is difficult*  
003 *to judge without strong real-data baselines and careful er-*  
004 *ror analysis. This paper presents an empirical CIFAR-10*  
005 *study designed as a baseline and diagnostic framework for*  
006 *generative visual recognition. Starting from an existing clas-*  
007 *sification report, we reorganize the study into a double-blind*  
008 *CVPR workshop style paper and evaluate a broad spectrum*  
009 *of methods: HOG features with KNN, SVM, and random*  
010 *forests; convolutional neural networks including VGG19-BN,*  
011 *ResNet, Inception-v4, and DenseNet; Pyramid Vision Trans-*  
012 *former; deep-feature extraction followed by classical clas-*  
013 *sifiers; and ImageNet-based transfer learning. The results*  
014 *show a consistent transition from handcrafted descriptors to*  
015 *learned representations: HOG-based methods reach at most*  
016 *55.13% test accuracy, optimized ResNet reaches 93.60%,*  
017 *and DenseNet transfer learning reaches 93.67%. However,*  
018 *high average accuracy does not remove systematic failure*  
019 *modes. Confusion matrices and high-confidence errors re-*  
020 *veal persistent ambiguities between cat/dog, deer/horse, au-*  
021 *tomobile/truck, and airplane/ship, often caused by low reso-*  
022 *lution, background correlation, and semantic overlap. We*  
023 *use these findings to define an error-aware protocol for fu-*  
024 *ture generative augmentation: synthetic images should be*  
025 *evaluated not only by overall accuracy gains, but also by*  
026 *their ability to reduce class-wise confusion without introduc-*  
027 *ing label or background bias.*

## 028 1. Introduction

029 Image classification remains one of the clearest ways to  
030 study visual representation learning. Even though modern  
031 networks can obtain strong performance on small bench-  
032 marks, datasets such as CIFAR-10 [12] are still useful be-  
033 cause they expose the interaction between feature design,  
034 model capacity, optimization, transfer learning, and class  
035 ambiguity. CIFAR-10 contains only  $32 \times 32$  images, but its

ten categories include several visually overlapping pairs. The  
low resolution makes the benchmark a compact test bed for  
asking whether a representation captures object semantics or  
relies on coarse texture and background cues.

At the same time, generative image models have rapidly  
changed the way visual data can be produced and reused.  
Generative adversarial networks [6], variational autoen-  
coders [11], and diffusion models [9, 13] have demon-  
strated increasingly strong image synthesis capabilities. This  
progress has motivated a natural question: can synthetic  
images generated by modern generative models improve  
downstream visual recognition? Prior work has explored  
synthetic data for representation learning, robustness, and  
data augmentation [1, 8, 14], but its benefit is often task-  
dependent. Synthetic images may improve diversity and  
reduce overfitting, but they may also introduce distribution  
shift, label noise, background bias, or unrealistic class se-  
mantics.

In this work, we revisit CIFAR-10 as a controlled small-  
scale benchmark for studying this question from the recog-  
nition side. Rather than beginning with a generative model  
and assuming that additional images should help, we first  
build a set of classical and deep classification baselines to  
understand where errors occur and which types of visual am-  
biguity remain difficult. Our study begins with handcrafted  
HOG features [4] combined with classical classifiers, includ-  
ing K-nearest neighbors, support vector machines [3], and  
random forests [2]. We then compare several neural archite-  
ctures, including VGG [15], ResNet [7], Inception [16, 17],  
DenseNet [10], and Pyramid Vision Transformer [18]. We  
also examine deep features with non-neural classifiers and  
transfer learning from ImageNet [5]. Finally, we analyze  
error cases and formulate a generative augmentation proto-  
col for evaluating the effect of class-conditioned synthetic  
images.

Our motivation is that small-scale classification provides  
a simple but informative setting for generative augmentation.  
If synthetic images are useful, they should not only improve  
overall accuracy but also reduce meaningful class confusions.  
For example, a useful augmentation method should help

076	distinguish cats from dogs, deer from horses, automobiles	126
077	from trucks, and airplanes from ships, rather than merely	127
078	increasing the number of training images. Conversely, if	128
079	synthetic images emphasize spurious background cues or	129
080	produce semantically inconsistent samples, they may worsen	130
081	exactly these ambiguous cases. This distinction is central for	131
082	generative visual recognition: the value of generated images	132
083	should be judged by their effect on the decision boundary,	133
084	not only by visual realism or sample count.	134
085	This paper makes three main contributions:	135
086	• We provide a unified empirical comparison of handcrafted	
087	features, classical classifiers, deep networks, transformer-	
088	based models, deep-feature classifiers, and transfer learn-	
089	ing for CIFAR-10 classification.	
090	• We turn the original project results into a compact error	
091	analysis, using confusion matrices and high-confidence	
092	failure cases to identify persistent ambiguities relevant to	
093	synthetic data evaluation.	
094	• We propose an error-aware generative augmentation proto-	
095	col for studying when class-conditioned synthetic images	
096	help small-scale visual recognition and when they may	
097	introduce semantic or background bias.	
098	The goal of this workshop paper is not to claim that syn-	
099	thetic data always improves CIFAR-10 classification. The	
100	available experiments are real-data recognition baselines and	
101	diagnostic analyses. We therefore position the generative	
102	component as a protocol and set of evaluation criteria for	
103	the next experimental stage. This is a useful contribution	
104	for a non-archival workshop setting because it connects an	
105	existing recognition study to the central question of how ge-	
106	nerative models can benefit computer vision, while keeping	
107	the claims faithful to the evidence.	
108	<b>2. Related Work</b>	
109	<b>2.1. Classical Image Features and Classifiers</b>	
110	Before the success of deep learning, visual recognition sys-	
111	tems often relied on handcrafted descriptors and classical	
112	machine learning models. HOG [4] represents local gradient	
113	orientation statistics and has been widely used for object	
114	detection and recognition. Classical classifiers such as K-	
115	nearest neighbors, support vector machines [3], and random	
116	forests [2] provide simple but useful baselines for evaluating	
117	the discriminative power of image features. Although these	
118	methods are generally outperformed by modern deep net-	
119	works, they remain valuable for understanding the transition	
120	from handcrafted to learned representations.	
121	<b>2.2. Deep Networks for Image Classification</b>	
122	Convolutional neural networks have substantially improved	
123	image classification performance. VGG [15] demonstrated	
124	the effectiveness of using deep networks with small convo-	
125	lutional filters. ResNet [7] introduced residual connections,	
	enabling much deeper models to be optimized effectively.	126
	Inception networks [16, 17] explored multi-branch convolu-	127
	tional modules and factorized convolutions to improve com-	128
	putational efficiency. DenseNet [10] further encouraged fea-	129
	ture reuse by connecting each layer to all subsequent layers.	130
	More recently, transformer-based models have been adapted	131
	to vision tasks. Pyramid Vision Transformer [18] introduces	132
	a hierarchical transformer backbone and spatial-reduction	133
	attention, making transformer representations more suitable	134
	for dense prediction and visual recognition.	135
	<b>2.3. Transfer Learning</b>	136
	Transfer learning has become a standard technique for im-	137
	proving performance when the target dataset is limited. Mod-	138
	els pretrained on large-scale datasets such as ImageNet [5]	139
	often learn transferable visual representations that can be	140
	adapted to downstream tasks. A common strategy is to ini-	141
	tialize a network with pretrained weights and fine-tune either	142
	the classifier head or the full network. This approach has	143
	been widely used in image classification and visual recog-	144
	nition [19, 20]. In our study, transfer learning serves as a	145
	strong baseline for evaluating whether additional synthetic	146
	images are useful beyond pretrained representations.	147
	<b>2.4. Generative Models and Synthetic Data</b>	148
	Generative models have achieved significant progress in im-	149
	age synthesis. GANs [6] introduced adversarial learning for	150
	generating realistic images. VAEs [11] provided a proba-	151
	bilistic latent-variable framework for generative modeling.	152
	Diffusion models [9] and latent diffusion models [13] have	153
	further improved image quality and controllability. These	154
	advances have motivated the use of synthetic images for	155
	training visual recognition models. Recent work has investi-	156
	gated whether generated images can replace or complement	157
	real images for downstream classification and representation	158
	learning [1, 8, 14]. However, the effectiveness of synthetic	159
	data depends on image fidelity, diversity, label consistency,	160
	and distribution alignment with the target task. Our work	161
	focuses on this issue in a controlled CIFAR-10 setting.	162
	<b>3. Method</b>	163
	<b>3.1. Problem Setup</b>	164
	We study image classification on CIFAR-10 [12], which	165
	contains ten object categories: airplane, automobile, bird, cat,	166
	deer, dog, frog, horse, ship, and truck. Given a training set	167
	$\mathcal{D}_r = \{(x_i, y_i)\}_{i=1}^N$ of real images and labels, the goal is to	168
	learn a classifier $f_\theta$ that predicts the class label $y$ for an input	169
	image $x$ . We use the standard 50,000/10,000 train/test split	170
	and report top-1 accuracy. All images are treated as RGB	171
	inputs at the native CIFAR-10 resolution unless a network	172
	requires resizing. The study has two goals. First, we compare	173
	a sequence of increasingly expressive recognition pipelines.	174

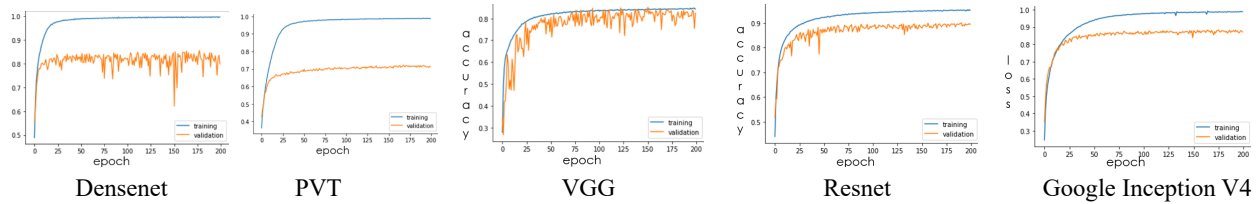


Figure 1. **Training curves of representative deep classification models on CIFAR-10.** We compare the training and validation behavior of VGG19-BN, ResNet, Inception-v4, and DenseNet with transfer learning. The curves show that learned deep representations substantially outperform handcrafted feature baselines, while optimized training schedules and transfer learning lead to improved convergence and better test accuracy.

175 Second, we use their errors to specify how future generated  
176 images should be evaluated.

### 177 3.2. Classical Feature-Based Baselines

178 We first evaluate handcrafted feature representations. For  
179 each image, we extract HOG descriptors [4], which summa-  
180 rize local gradient orientation distributions. Following the  
181 project setting, each  $8 \times 8$  pixel cell contributes an orienta-  
182 tion histogram, and each  $2 \times 2$  group of cells is normalized as  
183 a block before concatenation into a descriptor. The extracted  
184 feature vector is then used as input to classical classifiers:

$$185 \hat{y} = g(\phi_{\text{HOG}}(x)), \quad (1)$$

186 where  $\phi_{\text{HOG}}$  denotes the HOG feature extractor and  $g$  is a  
187 classifier such as K-nearest neighbors, support vector ma-  
188 chine, or random forest. For KNN, we vary the number of  
189 neighbors and compare uniform and distance weighting. For  
190 SVM, we use a linear classifier and tune the regularization  
191 parameter  $C$ . For random forests, we vary the number of  
192 estimators. These baselines provide a reference for under-  
193 standing the limitations of handcrafted features on small  
194 natural images.

### 195 3.3. Deep Classification Baselines

196 We then train deep neural networks directly on CIFAR-10.  
197 For a neural classifier  $f_{\theta}$ , we optimize the standard cross-  
198 entropy loss:

$$199 \mathcal{L}_{\text{ce}}(\theta) = -\frac{1}{N} \sum_{i=1}^N \sum_{c=1}^C \mathbf{1}[y_i = c] \log p_{\theta}(c|x_i), \quad (2)$$

200 where  $C = 10$  for CIFAR-10 and  $p_{\theta}(c|x_i)$  is the predicted  
201 probability for class  $c$ . We compare representative convolu-  
202 tional and transformer-based architectures, including VGG,  
203 ResNet, Inception, DenseNet, and Pyramid Vision Trans-  
204 former. The first group uses a common training recipe from  
205 the original project: Adam optimization, categorical cross  
206 entropy, batch size 128, dropout before the final fully con-  
207 nected layer, and  $L_2$  regularization. We then compare this  
208 recipe with a stronger SGD-based schedule using step decay,  
209 motivated by residual-network training practice [7].

### 210 3.4. Deep Feature Extraction with Classical Classi- 211 fiers

212 To separate feature learning from classifier design, we also  
213 use trained neural networks as feature extractors. Given a  
214 trained model, we take the activation before the final classi-  
215 fication layer as a deep feature:

$$216 z_i = \phi_{\theta}(x_i). \quad (3)$$

217 We then train classical classifiers on  $z_i$ , including SVM,  
218 random forest, and KNN. This setting evaluates whether  
219 the final softmax classifier is optimal or whether classical  
220 classifiers can better exploit the learned representation. It  
221 also separates two sources of improvement: representation  
222 learning in the backbone and decision-boundary learning in  
223 the classifier.

### 224 3.5. Transfer Learning

225 We further consider transfer learning using an ImageNet-  
226 pretrained DenseNet model. The classifier head is first  
227 trained while the backbone is fixed, and then the full network  
228 is fine-tuned. This strategy allows the model to benefit from  
229 large-scale visual pretraining while adapting to CIFAR-10.  
230 Transfer learning serves as a strong baseline because it al-  
231 ready incorporates external visual knowledge. In the project  
232 report, the top layers are trained first and the full model  
233 is then unfrozen for fine-tuning. This two-stage procedure  
234 obtains the strongest result in our study.

### 235 3.6. Generative Augmentation Protocol

236 The workshop theme motivates a natural extension: using  
237 generated images to improve recognition. Because  
238 the available project report does not contain completed  
239 generative experiments, we define the protocol explicitly  
240 rather than reporting unverified synthetic-data results. Let  
241  $\mathcal{D}_s = \{(\tilde{x}_j, \tilde{y}_j)\}_{j=1}^M$  be a synthetic training set generated by  
242 a class-conditioned model. For each class  $c$ , a generative  
243 model produces synthetic samples:

$$244 \tilde{x} \sim p_{\psi}(x|y = c), \quad (4)$$

245 where  $p_\psi$  denotes a generative image model such as a GAN  
246 or diffusion model. The augmented training set is

$$247 \quad \mathcal{D}_{aug} = \mathcal{D}_r \cup \lambda \mathcal{D}_s, \quad (5)$$

248 where  $\lambda$  controls the relative contribution of synthetic images  
249 during training. We then train classifiers under four settings:

- 250 • **Real only**: training only on CIFAR-10 real images.
- 251 • **Synthetic only**: training only on generated images.
- 252 • **Real + synthetic**: training on the union of real and syn-  
253 thetic images.
- 254 • **Real + filtered synthetic**: training on real images and syn-  
255 thetic images filtered by confidence or similarity criteria.

256 The purpose of this protocol is to determine whether  
257 synthetic images improve downstream classification, espe-  
258 cially in ambiguous classes. Rather than only reporting  
259 average accuracy, we evaluate class-wise confusion and high-  
260 confidence failure cases. For a test image  $x_i$ , we define a  
261 failure margin

$$262 \quad m_i = \max_{c \neq y_i} p_\theta(c|x_i) - p_\theta(y_i|x_i). \quad (6)$$

263 Large positive margins identify confident mistakes. Future  
264 generative augmentation is considered useful only if it im-  
265 proves accuracy while reducing these margins for ambiguous  
266 class pairs.

## 267 4. Experiments

### 268 4.1. Dataset and Evaluation Metrics

269 We conduct experiments on CIFAR-10 [12]. The dataset  
270 contains 50,000 training images and 10,000 test images from  
271 ten object classes. Each image is an RGB image of size  
272  $32 \times 32$ . We report top-1 classification accuracy as the pri-  
273 mary metric. To analyze class-level behavior, we also report  
274 confusion matrices and inspect high-confidence failure cases.  
275 Unless otherwise stated, accuracies are reported in percent.  
276 The numbers in this section are transcribed from the project  
277 report and reorganized into a common evaluation format.

### 278 4.2. Classical Baselines

279 We first evaluate HOG features combined with KNN, SVM,  
280 and random forest classifiers. These methods provide a non-  
281 deep-learning baseline for CIFAR-10 classification. The  
282 best KNN setting uses distance weighting with  $k = 6$  and  
283 reaches 47.99% test accuracy. For the linear SVM, increas-  
284 ing  $C$  improves performance up to the tested range, with the  
285 best reported test accuracy of 50.64%. The random forest  
286 baseline reaches 55.13% but obtains 100% training accuracy,  
287 indicating severe overfitting. Overall, HOG-based classi-  
288 fiers achieve substantially lower accuracy than deep neural  
289 networks. This result suggests that handcrafted gradient  
290 statistics are insufficient for capturing semantic variation  
291 in CIFAR-10, especially for animal categories and visually  
292 ambiguous object classes.

Table 1. Classical baselines using HOG features on CIFAR-10.

Method	Feature	Train Acc.	Test Acc.
KNN	HOG	100.00	47.99
SVM	HOG	50.40	50.64
Random Forest	HOG	100.00	55.13

Table 2. Deep neural network baselines on CIFAR-10.

Model	Training Acc.	Test Acc.
VGG19-BN	85.78	82.64
ResNet	95.10	89.36
Inception-v4	98.87	87.13
PVT-Tiny	98.58	71.43
DenseNet	99.36	78.40
DenseNet Transfer	99.82	93.67

Table 3. Effect of hyperparameter optimization.

Model	Original Test Acc.	Optimized Test Acc.
VGG19-BN	82.64	92.29
ResNet	89.36	93.60
Inception-v4	87.13	89.70

### 293 4.3. Deep Network Baselines

294 We next compare representative deep neural architectures.  
295 Compared with HOG-based methods, deep networks achieve  
296 significantly stronger performance, indicating the impor-  
297 tance of learned hierarchical representations. Among the  
298 evaluated models, ResNet and DenseNet provide strong re-  
299 sults, while transfer learning with DenseNet achieves the  
300 best performance in our preliminary study.

### 301 4.4. Effect of Hyperparameter Optimization

302 We observe that hyperparameter choices have a large effect  
303 on deep classification performance. In particular, replacing  
304 a fixed learning rate with a step-decay schedule and using  
305 stochastic gradient descent improves the performance of sev-  
306 eral architectures. This observation is consistent with the  
307 training strategy used in residual networks [7]. The results  
308 highlight that comparisons between architectures should be  
309 interpreted carefully unless the optimization settings are  
310 sufficiently tuned. The improvement is most visible for  
311 VGG19-BN and ResNet. VGG19-BN rises from 82.64% to  
312 92.29% test accuracy in the original report’s optimized run,  
313 while ResNet reaches 93.60%. Inception-v4 also improves,  
314 although less strongly than ResNet under the reported sched-  
315 ular.

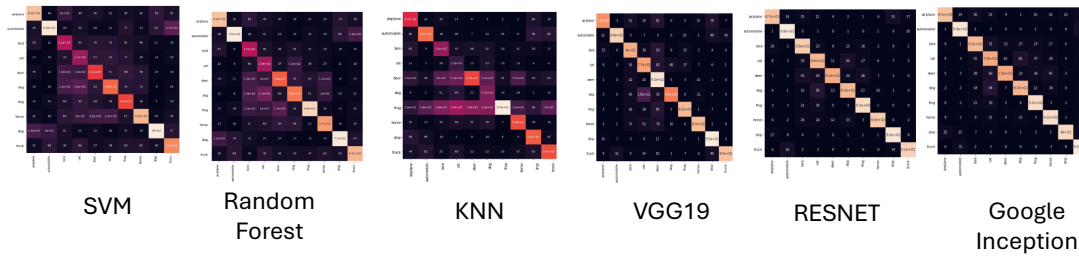


Figure 2. Confusion matrix of the strongest CIFAR-10 classifier.

Table 4. Deep feature extraction followed by classical classifiers.

Feature Backbone	SVM	Random Forest	KNN
VGG19-BN	85.63	85.79	84.92
ResNet	90.59	90.91	90.20
Inception-v4	88.17	88.12	88.10

Table 5. Final ranking of the strongest reported settings.

Setting	Test Acc.
DenseNet transfer learning	93.67
Optimized ResNet	93.60
Optimized VGG19-BN	92.29
ResNet features + random forest	90.91
Optimized Inception-v4	89.70
HOG + random forest	55.13

#### 316 4.5. Deep Features with Classical Classifiers

317 We also evaluate whether features extracted from trained  
 318 neural networks can be better classified by classical classi-  
 319 fiers. For each trained model, we extract features before the  
 320 final softmax layer and train SVM, random forest, and KNN  
 321 classifiers. The results show that deep features combined  
 322 with classical classifiers can slightly improve over the origi-  
 323 nal softmax classifier in some cases. This suggests that the  
 324 learned representation contains discriminative information  
 325 that may not be fully exploited by the final classifier layer.  
 326 The strongest feature-classifier combination in this setting  
 327 is ResNet features with a random forest, reaching 90.91%  
 328 test accuracy. Although the gains over softmax are mod-  
 329 est, the result is useful diagnostically because it shows that  
 330 representation quality and classifier choice can be studied  
 331 separately.

332 Fig. 1 shows the training curves of the evaluated deep  
 333 models. Compared with the original optimization setting, im-  
 334 proved learning-rate scheduling leads to more stable conver-  
 335 gence and higher test accuracy. Among the evaluated models,  
 336 DenseNet with transfer learning achieves the strongest final  
 337 performance. To further understand the failure modes, we  
 338 visualize the confusion matrix of the strongest model in  
 339 Fig. 4.3. Although the model achieves high overall accuracy,  
 340 several class pairs remain difficult to separate. In particu-  
 341 lar, cat/dog, deer/horse, automobile/truck, and airplane/ship  
 342 show recurring confusion, indicating that background corre-  
 343 lation and semantic similarity still affect the classifier.

#### 344 4.6. Transfer Learning and Model Ranking

345 The best result in the report is obtained by DenseNet transfer  
 346 learning. The model is initialized from ImageNet, trained

first through newly added top layers, and then fine-tuned 347  
 end to end. After the first stage, the model reaches 79.69% 348  
 test accuracy; after full fine-tuning, it reaches 93.67%. This 349  
 result is close to the optimized ResNet baseline but slightly 350  
 higher, making DenseNet transfer learning the strongest 351  
 reported model. The comparison is important for generative 352  
 augmentation: synthetic images should be compared against 353  
 this strong transfer baseline, not only against weaker from- 354  
 scratch models. 355

#### 356 4.7. What These Results Mean for Synthetic Data

The baseline results constrain what a useful generative aug- 357  
 mentation study must show. First, generated data should be 358  
 evaluated against a strong transfer-learning baseline because 359  
 transfer already imports external visual knowledge. Second, 360  
 improvements should be class-wise, not only average. For ex- 361  
 ample, if synthetic samples increase the number of airplane 362  
 images but preserve the same blue-sky background bias, 363  
 they may not reduce airplane/bird/ship confusion. Third, 364  
 synthetic data should be filtered when the conditioning label 365  
 and image semantics disagree. This is especially important 366  
 for classes with fuzzy boundaries, such as automobile/truck 367  
 and cat/dog. 368

### 369 5. Error Analysis and Discussion

#### 370 5.1. Class-Wise Confusion

Overall accuracy provides a useful summary but does not 371  
 fully explain how models fail. Therefore, we analyze con- 372

373 fusion matrices and high-confidence failure cases. The  
374 strongest baseline still makes systematic mistakes between  
375 visually similar categories. For example, cats and dogs are  
376 frequently confused because of similar texture, pose, and  
377 facial structure. Deer and horses are also difficult to distin-  
378 guish when horns are not visible or when the object occupies  
379 only a small region of the image. Automobiles and trucks  
380 form another ambiguous pair, especially for vans or vehicles  
381 captured from partial viewpoints. Airplanes, birds, and ships  
382 can be confused when the background is dominated by sky  
383 or water.

384 These errors suggest that CIFAR-10 classification is not  
385 only a problem of model capacity. Some errors are caused by  
386 low resolution, ambiguous labels, background correlation,  
387 and semantic overlap between classes. This observation is  
388 important for evaluating generative augmentation. If syn-  
389 thetic images merely increase the number of background-  
390 correlated samples, they may not reduce these errors. A  
391 useful synthetic augmentation method should instead im-  
392 prove class-specific visual diversity while preserving label  
393 consistency.

## 394 5.2. High-Confidence Failure Cases

395 We further inspect high-confidence wrong predictions. A  
396 high-confidence failure occurs when the model assigns a  
397 high probability to an incorrect class while assigning a low  
398 probability to the true class. These cases are particularly  
399 informative because they reveal systematic biases rather than  
400 random uncertainty. For example, an airplane with a blue  
401 background may be predicted as a bird or ship, while a small  
402 truck may be predicted as an automobile. Some examples  
403 are also visually ambiguous to humans, which indicates that  
404 the ground-truth label may not always capture the perceptual  
405 ambiguity of the image. The original report also compares  
406 the strongest DenseNet transfer model with the second-best  
407 optimized ResNet model. This comparison is useful because  
408 common errors across two different architectures are less  
409 likely to be incidental. For airplanes, many common errors  
410 involve small objects against sky or water backgrounds. For  
411 automobiles and trucks, the common failures often involve  
412 vans or partial vehicle views. For animal classes, the recur-  
413 ring errors are concentrated around cat/dog and deer/horse  
414 boundaries, where texture, pose, and missing discriminative  
415 parts make the class label uncertain.

416 This analysis motivates the use of class-wise and failure-  
417 aware metrics for synthetic data evaluation. If generative  
418 augmentation improves only easy classes while worsening  
419 ambiguous classes, its overall benefit may be limited. Con-  
420 versely, if synthetic images reduce high-confidence errors in  
421 ambiguous categories, they may provide meaningful comple-  
422 mentary information.

## 5.3. Architecture-Level Observations 423

424 The comparison across architectures gives several practical  
425 lessons. First, handcrafted HOG features are too local and  
426 edge-oriented to describe the semantic structure of CIFAR-  
427 10 objects. They work better for classes with distinctive  
428 shape and background patterns, such as airplane, ship, and  
429 truck, but struggle on animals. Second, convolutional archi-  
430 tectures remain strong on this benchmark. ResNet benefits  
431 from residual optimization, while DenseNet benefits from  
432 feature reuse and transfer learning. Third, the PVT-Tiny re-  
433 sult is weaker than expected from large-scale recognition lit-  
434 erature. This does not imply that transformers are unsuitable  
435 for CIFAR-10; rather, it suggests that small low-resolution  
436 datasets and limited training recipes can make transformer  
437 optimization fragile. Fourth, replacing the final softmax clas-  
438 sifier with a classical classifier over deep features produces  
439 small gains for several backbones. This indicates that the  
440 penultimate representation is discriminative, but the final  
441 decision layer may not always be the best possible classifier  
442 under the chosen training setting.

## 5.4. Implications for Generative Augmentation 443

444 The proposed generative augmentation protocol should be  
445 evaluated with three criteria.

**Fidelity.** Synthetic images should be visually plausible  
446 and consistent with the target class. Low-quality images may  
447 introduce noise and reduce classifier performance. 448

**Diversity.** Synthetic images should cover variations not  
449 sufficiently represented in the real training set. If the gen-  
450 erated samples are too similar to existing images, they may  
451 provide little benefit. 452

**Label consistency.** The generated image should match  
453 the intended class label. For ambiguous categories such as  
454 cat/dog and automobile/truck, label inconsistency can be  
455 especially harmful. 456

457 These criteria suggest that synthetic data should not be  
458 used blindly. Instead, it should be filtered or weighted ac-  
459 cording to confidence, similarity, or semantic consistency.  
460 For example, a pretrained classifier or vision-language model  
461 can be used to remove generated samples whose predicted  
462 class disagrees with the conditioning label. Feature-space  
463 similarity can also be used to reject samples that are too  
464 far from the real data distribution. For CIFAR-10, a practi-  
465 cal protocol would generate a fixed number of images per  
466 class, train the same classifier under real-only and real-plus-  
467 synthetic settings, and then report both overall accuracy and  
468 per-class confusion changes. The most important compar-  
469 isons should focus on the ambiguous pairs identified above.  
470 If synthetic images improve easy classes but do not improve  
471 these pairs, the augmentation is unlikely to address the main  
472 recognition failure modes. If the generated images reduce  
473 common DenseNet/ResNet errors, the evidence for useful  
474 generative augmentation is much stronger.

## 475 5.5. Submission Relevance

476 This study is relevant to a workshop on generative models  
477 for computer vision because it frames synthetic data as a  
478 recognition problem rather than only a synthesis problem.  
479 The central question is not whether a generative model can  
480 produce plausible CIFAR-like images, but whether those  
481 images contain the variations needed to improve recognition.  
482 The baseline and error analysis in this paper provide the  
483 measurement scaffold for that question. In a non-archival  
484 workshop setting, this makes the work suitable as an ongoing  
485 study: the current paper establishes the recognition baselines,  
486 while the proposed protocol defines the next set of generative  
487 experiments.

## 488 5.6. Limitations

489 This study has several limitations. First, CIFAR-10 is a  
490 small and low-resolution dataset, so conclusions may not  
491 directly transfer to large-scale or high-resolution recogni-  
492 tion tasks. Second, the current baseline study focuses on  
493 classification accuracy and error analysis; the full generative  
494 augmentation experiments require additional synthetic data  
495 generation and filtering. Third, different generative models  
496 may produce synthetic images with different levels of fidelity  
497 and diversity. Therefore, future work should compare GAN-  
498 based, diffusion-based, and retrieval-augmented synthetic  
499 data sources under the same evaluation protocol.

## 500 6. Conclusion

501 We presented a preliminary CIFAR-10 study for analyz-  
502 ing small-scale visual classification with classical features,  
503 deep networks, transfer learning, and generative augmenta-  
504 tion. Our baseline experiments show that deep learned rep-  
505 resentations substantially outperform handcrafted HOG fea-  
506 tures, and that DenseNet-based transfer learning provides the  
507 strongest performance among the evaluated models. More  
508 importantly, our error analysis reveals persistent class am-  
509 biguities that are not fully explained by average accuracy  
510 alone. These findings motivate a generative augmentation  
511 protocol that evaluates synthetic images not only by their  
512 effect on overall classification accuracy, but also by their  
513 impact on class-wise confusion and high-confidence failure  
514 cases. We hope this study provides a simple and controlled  
515 framework for understanding when synthetic images help  
516 visual recognition and when they may introduce additional  
517 semantic or background bias.

## 518 References

519 [1] Shekoofeh Azizi, Simon Kornblith, Chitwan Saharia, Mo-  
520 hammad Norouzi, and David J. Fleet. Synthetic data from  
521 diffusion models improves imagenet classification. In *Trans-*  
522 *actions on Machine Learning Research*, 2023. 1, 2

- [2] Leo Breiman. Random forests. *Machine Learning*, 45(1): 533  
534 5–32, 2001. 1, 2
- [3] Corinna Cortes and Vladimir Vapnik. Support-vector net- 525  
works. *Machine Learning*, 20(3):273–297, 1995. 1, 2 526
- [4] Navneet Dalal and Bill Triggs. Histograms of oriented gradi- 527  
ents for human detection. In *Proceedings of the IEEE Con-*  
528 *ference on Computer Vision and Pattern Recognition*, pages  
529 886–893, 2005. 1, 2, 3 530
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li 531  
Fei-Fei. Imagenet: A large-scale hierarchical image database.  
532 In *Proceedings of the IEEE Conference on Computer Vision*  
533 *and Pattern Recognition*, pages 248–255, 2009. 1, 2 534
- [6] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing 535  
Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and  
536 Yoshua Bengio. Generative adversarial nets. In *Advances in*  
537 *Neural Information Processing Systems*, 2014. 1, 2 538
- [7] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 539  
Deep residual learning for image recognition. In *Proceedings*  
540 *of the IEEE Conference on Computer Vision and Pattern*  
541 *Recognition*, pages 770–778, 2016. 1, 2, 3, 4 542
- [8] Ruifei He, Shuyang Sun, Xin Yu, Chuhui Xue, Wenqing 543  
Zhang, Philip H. S. Torr, Song Bai, and Xiaojuan Qi. Is  
544 synthetic data from generative models ready for image recog-  
545 nition? *arXiv preprint arXiv:2210.07574*, 2022. 1, 2 546
- [9] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffu- 547  
sion probabilistic models. In *Advances in Neural Information*  
548 *Processing Systems*, 2020. 1, 2 549
- [10] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kil- 550  
ian Q. Weinberger. Densely connected convolutional net-  
551 works. In *Proceedings of the IEEE Conference on Computer*  
552 *Vision and Pattern Recognition*, pages 4700–4708, 2017. 1, 2 553
- [11] Diederik P. Kingma and Max Welling. Auto-encoding vari- 554  
ational bayes. *arXiv preprint arXiv:1312.6114*, 2013. 1, 555  
2 556
- [12] Alex Krizhevsky. Learning multiple layers of features from 557  
tiny images. Technical report, University of Toronto, 2009. 1, 558  
2, 4 559
- [13] Robin Rombach, Andreas Blattmann, Dominik Lorenz, 560  
Patrick Esser, and Bjorn Ommer. High-resolution image  
561 synthesis with latent diffusion models. In *Proceedings of*  
562 *the IEEE/CVF Conference on Computer Vision and Pattern*  
563 *Recognition*, pages 10684–10695, 2022. 1, 2 564
- [14] Mert Bulent Sariyildiz, Karteek Alahari, Diane Larlus, and 565  
Yannis Kalantidis. Fake it till you make it: Learning trans-  
566 ferable representations from synthetic imagenet clones. In  
567 *Proceedings of the IEEE/CVF Conference on Computer Vi-*  
568 *sion and Pattern Recognition*, pages 8011–8021, 2023. 1, 569  
2 570
- [15] Karen Simonyan and Andrew Zisserman. Very deep convo- 571  
lutional networks for large-scale image recognition. *arXiv*  
572 *preprint arXiv:1409.1556*, 2014. 1, 2 573
- [16] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, 574  
Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent  
575 Vanhoucke, and Andrew Rabinovich. Going deeper with  
576 convolutions. In *Proceedings of the IEEE Conference on*  
577 *Computer Vision and Pattern Recognition*, pages 1–9, 2015. 578  
1, 2 579

- 580 [17] Christian Szegedy, Sergey Ioffe, Vincent Vanhoucke, and  
581 Alexander A. Alemi. Inception-v4, inception-resnet and the  
582 impact of residual connections on learning. In *Proceedings*  
583 *of the AAAI Conference on Artificial Intelligence*, 2017. 1, 2
- 584 [18] Wenhai Wang, Enze Xie, Xiang Li, Deng-Ping Fan, Kaitao  
585 Song, Ding Liang, Tong Lu, Ping Luo, and Ling Shao. Pyra-  
586 mid vision transformer: A versatile backbone for dense predic-  
587 tion without convolutions. In *Proceedings of the IEEE/CVF*  
588 *International Conference on Computer Vision*, pages 568–  
589 578, 2021. 1, 2
- 590 [19] Jason Yosinski, Jeff Clune, Yoshua Bengio, and Hod Lipson.  
591 How transferable are features in deep neural networks? In  
592 *Advances in Neural Information Processing Systems*, 2014. 2
- 593 [20] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi,  
594 Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. A  
595 comprehensive survey on transfer learning. *Proceedings of*  
596 *the IEEE*, 109(1):43–76, 2020. 2