

3LC x AWS Hackathon - Final Write-Up

Chihuahua vs. Muffin - A Data-Centric Binary Image Classification Pipeline Using 3LC

1. Introduction

This project implements a binary image classification system for distinguishing between Chihuahuas and Muffins using a data-centric AI workflow powered by 3LC. Instead of focusing on model complexity, the goal was to steadily improve classification performance by improving data quality, label accuracy, and dataset structure through iterative refinement.

The workflow leverages 3LC's capabilities - including embeddings, confidence-based filtering, misclassification analysis, clustering, and table versioning - to build a clean, consistent, and highly informative training dataset.

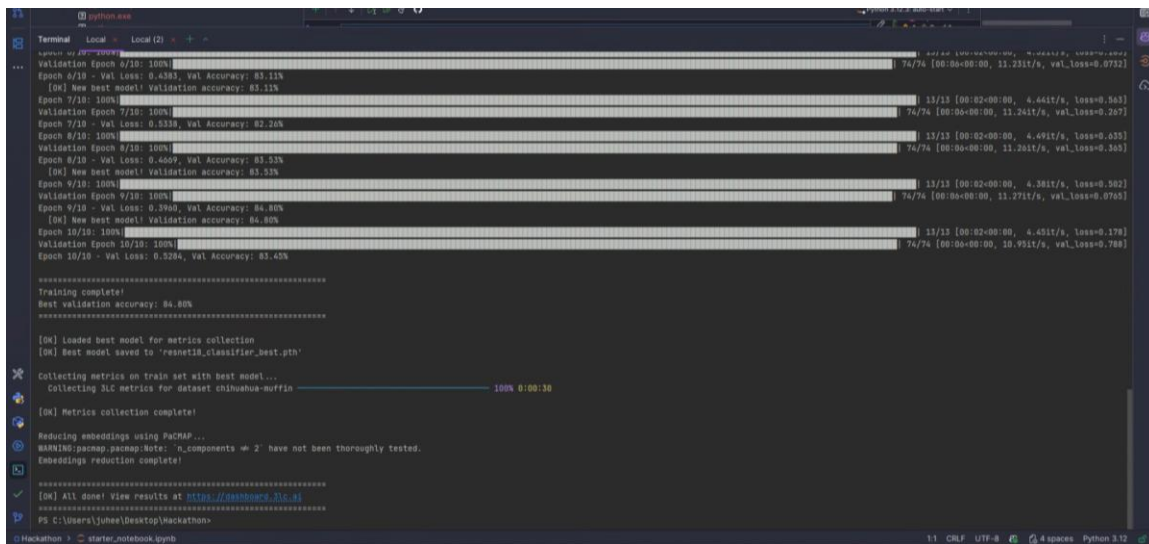
2. Baseline Setup

We began with:

- 200 labeled images (Chihuahua / Muffin)
- 4,533 unlabeled images under the undefined class
- ResNet-18 (no pretrained weights)
- Basic augmentations
- 20 training epochs, learning rate 0.0001

This created the baseline accuracy, serving as the foundation for all subsequent data-centric improvements.

Initial Accuracy: 84.80%



```
validation Epoch 6/10: 100% | 76/76 [00:00<00:00, 11.231t/s, val_loss=0.0732]
Epoch 6/10 - Val Loss: 0.4383, Val Accuracy: 85.11%
[OK] New best model! Validation accuracy: 85.11%
Epoch 7/10: 100% | 76/76 [00:00<00:00, 11.241t/s, val_loss=0.2697]
validation Epoch 7/10: 100% | 76/76 [00:00<00:00, 11.241t/s, val_loss=0.2697]
Epoch 7/10 - Val Loss: 0.5238, Val Accuracy: 82.16%
Epoch 8/10: 100% | 76/76 [00:00<00:00, 11.251t/s, val_loss=0.4357]
validation Epoch 8/10: 100% | 76/76 [00:00<00:00, 11.251t/s, val_loss=0.4357]
Epoch 8/10 - Val Loss: 0.4609, Val Accuracy: 83.53%
[OK] New best model! Validation accuracy: 83.53%
Epoch 9/10: 100% | 76/76 [00:00<00:00, 11.261t/s, val_loss=0.3802]
validation Epoch 9/10: 100% | 76/76 [00:00<00:00, 11.261t/s, val_loss=0.3802]
Epoch 9/10 - Val Loss: 0.3960, Val Accuracy: 84.80%
[OK] New best model! Validation accuracy: 84.80%
Epoch 10/10: 100% | 76/76 [00:00<00:00, 10.931t/s, val_loss=0.7807]
validation Epoch 10/10: 100% | 76/76 [00:00<00:00, 10.931t/s, val_loss=0.7807]
Epoch 10/10 - Val Loss: 0.5284, Val Accuracy: 83.45%

=====
Training complete!
Best validation accuracy: 84.80%
=====

[OK] Loaded best model for metrics collection
[OK] Best model saved to 'resnet18_classifier_best.pth'

Collecting metrics on train set with best model...
Collecting 3LC metrics for dataset chihuahua-muffin ----- 100% 0:00:10

[OK] Metrics collection complete!

Reducing embeddings using PaCMAP...
WARNING:pacmap.pacmap:Note: 'n_components' = 2 have not been thoroughly tested.
Embeddings reduction complete!

=====
[OK] All done! View results at https://analysis.3lc.ai
=====

PS C:\Users\johnd\Desktop\hackathon>
```

3. Data-Centric Iterative Loop

Our process followed the structured data-centric loop:

Train → Analyze → Fix → Retrain

Across 15 dataset versions, we consistently used model feedback, confidence metrics, embedding clusters, and misclassification reasoning to improve label quality and dataset balance.

4. Version-by-Version Dataset Improvements

Version 1 — High Confidence Auto-Labeling (>0.95)

High-confidence unlabeled samples were automatically assigned Chihuahua/Muffin labels. This expanded the dataset while maintaining precision.

Version 2 — High-Loss Sample Removal

Samples with high loss (~21.7) were identified and removed or corrected, reducing noisy and harmful examples.

Versions 3-7 — Iterative Labeling Using Confidence (>0.94 / >0.95)

Across multiple passes:

- Predicted labels were applied to unlabeled samples.
- ID-based filtering was performed.
- Outliers were re-labeled as undefined.

The dataset steadily grew with high-quality additions.

Versions 8-9 — Strict Threshold Labeling (>0.99)

Only extremely confident predictions (>0.99) were accepted. These samples were manually validated. Ambiguous or visually odd samples remained undefined.

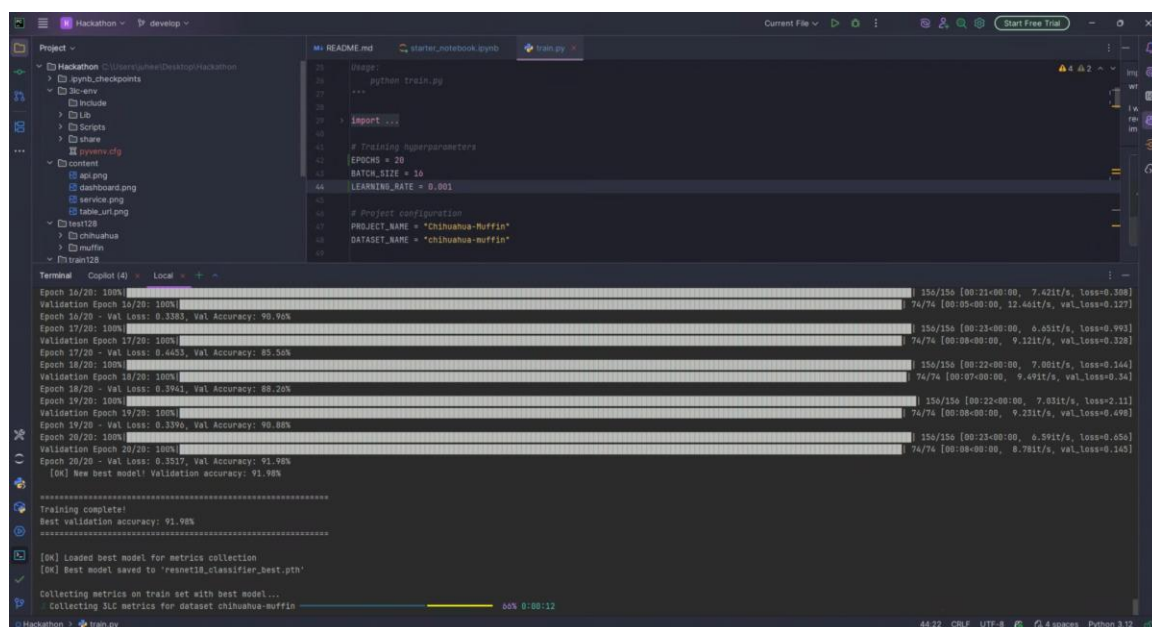
Versions 10-11 — Low Confidence Mining (<0.68 / <0.70)

- Embedding plot colored by confidence
- Low-confidence clusters were corrected
- Misabeled and inconsistent images were removed or fixed

Hyperparameter Improvements

- Epochs increased from 20 → 50
- Learning rate increased to 0.001

Mid-Pipeline Accuracy (after Iteration 10-11): 91.98%



The screenshot shows a Jupyter Notebook interface with a file explorer on the left, a code editor in the center, and a terminal at the bottom. The code editor displays a Python script for training a ResNet18 model. The terminal shows the training progress over 20 epochs, with validation accuracy reaching 91.98%.

```
1 # Imports
2 import torch
3 from torch import nn
4 from torchvision import datasets, transforms
5
6 # Training hyperparameters
7 EPOCHS = 20
8 BATCH_SIZE = 16
9 LEARNING_RATE = 0.001
10
11 # Project configuration
12 PROJECT_NAME = "Chihuahua-Muffin"
13 DATASET_NAME = "chihuahua-muffin"
```

Terminal output:

```
Epoch 16/20: 100% | 156/156 [00:21:00:00, 7.421t/s, loss=0.308]
Validation Epoch 16/20: 100% | 74/74 [00:05:00:00, 12.441t/s, val_loss=0.127]
Epoch 16/20 - Val Loss: 0.1283, Val Accuracy: 90.96%
Epoch 17/20: 100% | 156/156 [00:21:00:00, 6.851t/s, loss=0.993]
Validation Epoch 17/20: 100% | 74/74 [00:06:00:00, 9.121t/s, val_loss=0.328]
Epoch 17/20 - Val Loss: 0.4453, Val Accuracy: 85.50%
Epoch 18/20: 100% | 156/156 [00:22:00:00, 7.001t/s, loss=0.344]
Validation Epoch 18/20: 100% | 74/74 [00:06:00:00, 9.491t/s, val_loss=0.34]
Epoch 18/20 - Val Loss: 0.3941, Val Accuracy: 88.26%
Epoch 19/20: 100% | 156/156 [00:22:00:00, 7.031t/s, loss=2.11]
Validation Epoch 19/20: 100% | 74/74 [00:06:00:00, 9.231t/s, val_loss=0.498]
Epoch 19/20 - Val Loss: 0.3946, Val Accuracy: 90.80%
Epoch 20/20: 100% | 156/156 [00:23:00:00, 6.591t/s, loss=0.654]
Validation Epoch 20/20: 100% | 74/74 [00:06:00:00, 8.781t/s, val_loss=0.145]
Epoch 20/20 - Val Loss: 0.3337, Val Accuracy: 91.98%
[00] New Best model! Validation accuracy: 91.98%
Training complete!
Best validation accuracy: 91.98%
[00] Loaded best model for metrics collection
[00] Best model saved to 'resnet18_classifier_best.pth'
Collecting metrics on train set with best model...
Collecting 3LC metrics for dataset chihuahua-muffin
```

Version 12 — Cluster-Based Label Repairs

Low-confidence clusters were corrected, mislabels were fixed, and similar examples were grouped.

Version 13 — Reassigning Undefined Samples

Undefined samples were inspected using predicted label, confidence, and cluster location. Confident predictions were promoted to true labels.

Augmentation Upgrades

To improve generalization:

- RandomRotation
- ColorJitter
- RandomResizedCrop
- RandomHorizontalFlip

Version 14 — Outlier Removal & Fine Cluster Cleaning

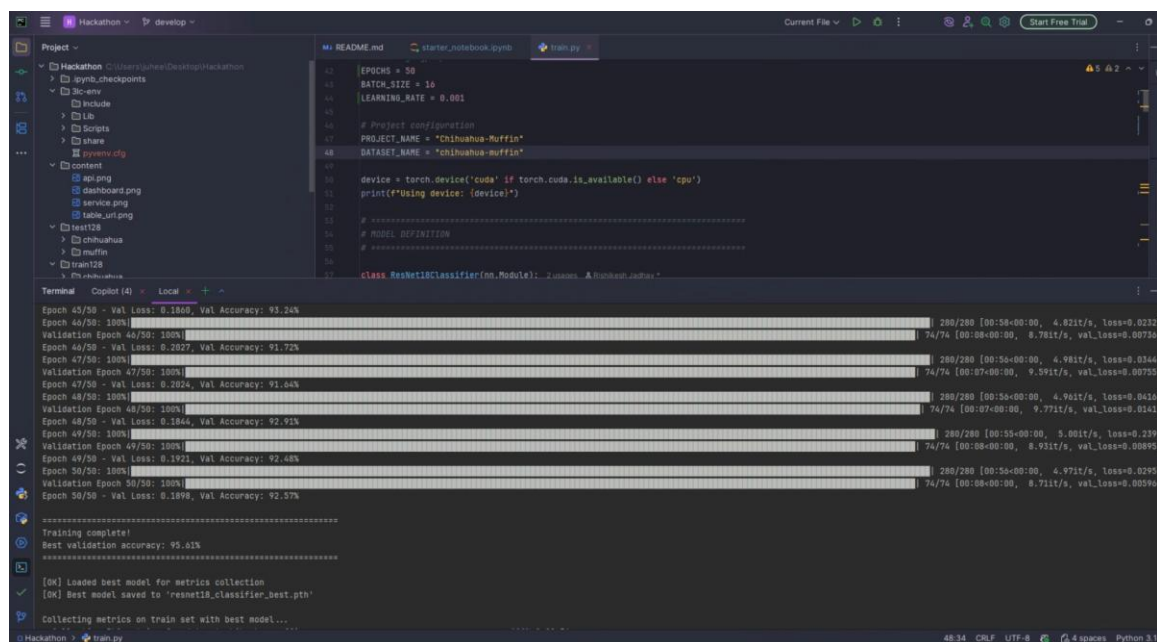
Removed noisy images, incorrect groupings, and visually inconsistent samples.

Version 15 — Final Dataset Cleanup + 50 Epoch Training

A final sweep ensured all labels were correct. The final model was trained with:

- 50 epochs
- Strong augmentations
- Dynamic learning rate scheduling

Final Accuracy: 95.61%



The screenshot shows a code editor with a project named 'Hackathon'. The file explorer on the left shows a directory structure with files like 'README.md', 'starter_notebook.ipynb', and 'train.py'. The main editor window displays the 'train.py' file, which contains the following code:

```
12 EPOCHS = 50
13 BATCH_SIZE = 16
14 LEARNING_RATE = 0.001
15
16 # Project configuration
17 PROJECT_NAME = "Chihuahua-Muffin"
18 DATASET_NAME = "Chihuahua-Muffin"
19
20 device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
21 print(f"Using device: {device}")
22
23 # =====
24 # MODEL DEFINITION
25 # =====
26 class ResNet18Classifier(nn.Module):
27     def __init__(self):
28         super(ResNet18Classifier, self).__init__()
29         self.resnet18 = torchvision.models.resnet18(pretrained=True)
30         self.fc = nn.Linear(512, 1000)
31         self.relu = nn.ReLU()
32         self.softmax = nn.Softmax(1)
33
34     def forward(self, x):
35         x = self.relu(self.resnet18(x))
36         x = self.fc(x)
37         x = self.softmax(x)
38         return x
```

The terminal window at the bottom shows the training progress for 50 epochs. The logs indicate that the training is complete and the best validation accuracy is 95.61%.

```
Epoch 45/50 - Val Loss: 0.1860, Val Accuracy: 93.24%
Epoch 46/50: 100% | 280/280 [00:54:00:00, 4.821t/s, loss=0.0232]
Validation Epoch 46/50: 100% | 74/74 [00:06:00:00, 8.781t/s, val_loss=0.00736]
Epoch 46/50 - Val Loss: 0.2027, Val Accuracy: 91.72%
Epoch 47/50: 100% | 280/280 [00:54:00:00, 4.981t/s, loss=0.0344]
Validation Epoch 47/50: 100% | 74/74 [00:07:00:00, 9.591t/s, val_loss=0.00753]
Epoch 47/50 - Val Loss: 0.2024, Val Accuracy: 91.64%
Epoch 48/50: 100% | 280/280 [00:54:00:00, 4.961t/s, loss=0.0416]
Validation Epoch 48/50: 100% | 74/74 [00:07:00:00, 9.771t/s, val_loss=0.0141]
Epoch 48/50 - Val Loss: 0.1844, Val Accuracy: 92.91%
Epoch 49/50: 100% | 280/280 [00:55:00:00, 5.801t/s, loss=0.239]
Validation Epoch 49/50: 100% | 74/74 [00:08:00:00, 8.931t/s, val_loss=0.00895]
Epoch 49/50 - Val Loss: 0.1923, Val Accuracy: 92.48%
Epoch 50/50: 100% | 280/280 [00:55:00:00, 4.971t/s, loss=0.0295]
Validation Epoch 50/50: 100% | 74/74 [00:08:00:00, 8.711t/s, val_loss=0.00596]
Epoch 50/50 - Val Loss: 0.1898, Val Accuracy: 92.57%
=====
Training complete!
Best validation accuracy: 95.61%
=====
[OK] Loaded best model for metrics collection
[OK] Best model saved to 'resnet18_classifier_best.pth'
Collecting metrics on train set with best model...
```

5. Challenges & Learnings

- The biggest challenge was working with a very small labeled dataset and a large pool of unlabeled images.
- Many unlabeled images looked extremely similar, which made confidence-based labeling tricky at times.
- The model initially overfitted because the dataset was not clean, so improving labels became more important than tuning the model.
- Visualizing embeddings in 3LC helped clearly identify mislabeled samples and confusing clusters.
- A key learning was that data quality matters more than model complexity—small corrections created noticeable accuracy improvements.

6. Future Improvements

- Explore semi-supervised techniques like refined pseudo-labeling to improve confidence in difficult regions.
- Build a small manually curated validation set for more reliable accuracy tracking.
- Experiment with additional augmentation strategies for more robustness.
- Extend this workflow to multi-class problems or real-world datasets.
- Automate parts of the 3LC labeling workflow to speed up future iterations.

7. Conclusion

This project demonstrates a complete data-centric AI workflow powered by 3LC, showing how dataset refinement, not increased model complexity is the key driver of performance gains. Through iterative labeling, noise reduction, confidence analysis, clustering, and misclassification review, a robust dataset was created that enabled strong performance from a simple ResNet-18 classifier.