

✓ CSC14120 - Parallel Programming Final Project

CUDA Autoencoder for CIFAR-10 Feature Learning

University of Science - Vietnam National University, Ho Chi Minh City

Faculty of Information Technology

Team Members:

- Member 1: Huỳnh Lê Hải Dương - 22127081
- Member 2: Nguyễn Triệu Khoáng - 22127204

Video Presentation: [YouTube Link]

Section 1: Problem Description

1.1 Problem Statement

Feature engineering is a fundamental challenge in machine learning: how do we automatically discover good representations of data that capture its underlying structure?

In this project, we implement an **Autoencoder-based unsupervised feature learning system** for image classification on the CIFAR-10 dataset. The project focuses on:

1. Stage 1 - Unsupervised Feature Learning:

- Train a convolutional autoencoder to reconstruct CIFAR-10 images
- Learn to encode $32 \times 32 \times 3$ images into an 8,192-dimensional feature representation
- No labels are used during autoencoder training

2. Stage 2 - Supervised Classification:

- Extract features from the trained encoder
- Train an SVM classifier on learned features
- Evaluate classification performance

1.2 Motivation for GPU Acceleration

Training deep neural networks on CPU is extremely slow due to:

- Massive amount of matrix operations in convolution
- Sequential execution limits throughput
- Large dataset (50,000 training images)

GPU parallelization can achieve **>20x speedup** by:

- Parallel computation across thousands of CUDA cores
- Optimized memory access patterns
- Kernel fusion to reduce memory bandwidth requirements

```
# Setup and Install Dependencies
%pip install -q matplotlib numpy scikit-learn seaborn

# Check GPU availability
!nvidia-smi

# Download CIFAR-10 dataset
!mkdir -p data/models
!cd data && wget -q https://www.cs.toronto.edu/~kriz/cifar-10-binary.tar.gz
!cd data && tar -xzf cifar-10-binary.tar.gz
!ls data/cifar-10-batches-bin/
```

```
# Visualize CIFAR-10 samples
import numpy as np
import matplotlib.pyplot as plt
```

```

def load_cifar10_batch(filename):
    """Load a single CIFAR-10 batch file."""
    with open(filename, 'rb') as f:
        data = np.frombuffer(f.read(), dtype=np.uint8)

    # Each record: 1 byte label + 3072 bytes image
    data = data.reshape(-1, 3073)
    labels = data[:, 0]
    images = data[:, 1:].reshape(-1, 3, 32, 32).transpose(0, 2, 3, 1)
    return images, labels

# Load first batch for visualization
images, labels = load_cifar10_batch('data/cifar-10-batches-bin/data_batch_1.bin')

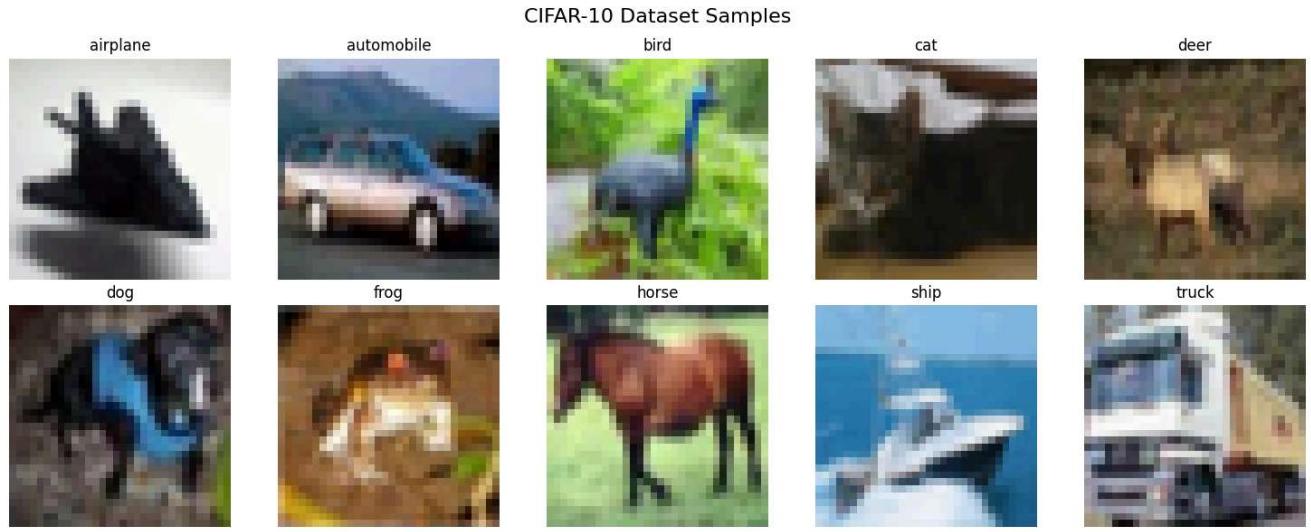
# Class names
class_names = ['airplane', 'automobile', 'bird', 'cat', 'deer',
               'dog', 'frog', 'horse', 'ship', 'truck']

# Display samples from each class
fig, axes = plt.subplots(2, 5, figsize=(15, 6))
for i, ax in enumerate(axes.flat):
    # Find first image of class i
    idx = np.where(labels == i)[0][0]
    ax.imshow(images[idx])
    ax.set_title(class_names[i])
    ax.axis('off')

plt.suptitle('CIFAR-10 Dataset Samples', fontsize=16)
plt.tight_layout()
plt.show()

print(f"Training set: {50000} images")
print(f"Test set: {10000} images")
print(f"Image shape: {images[0].shape}")

```

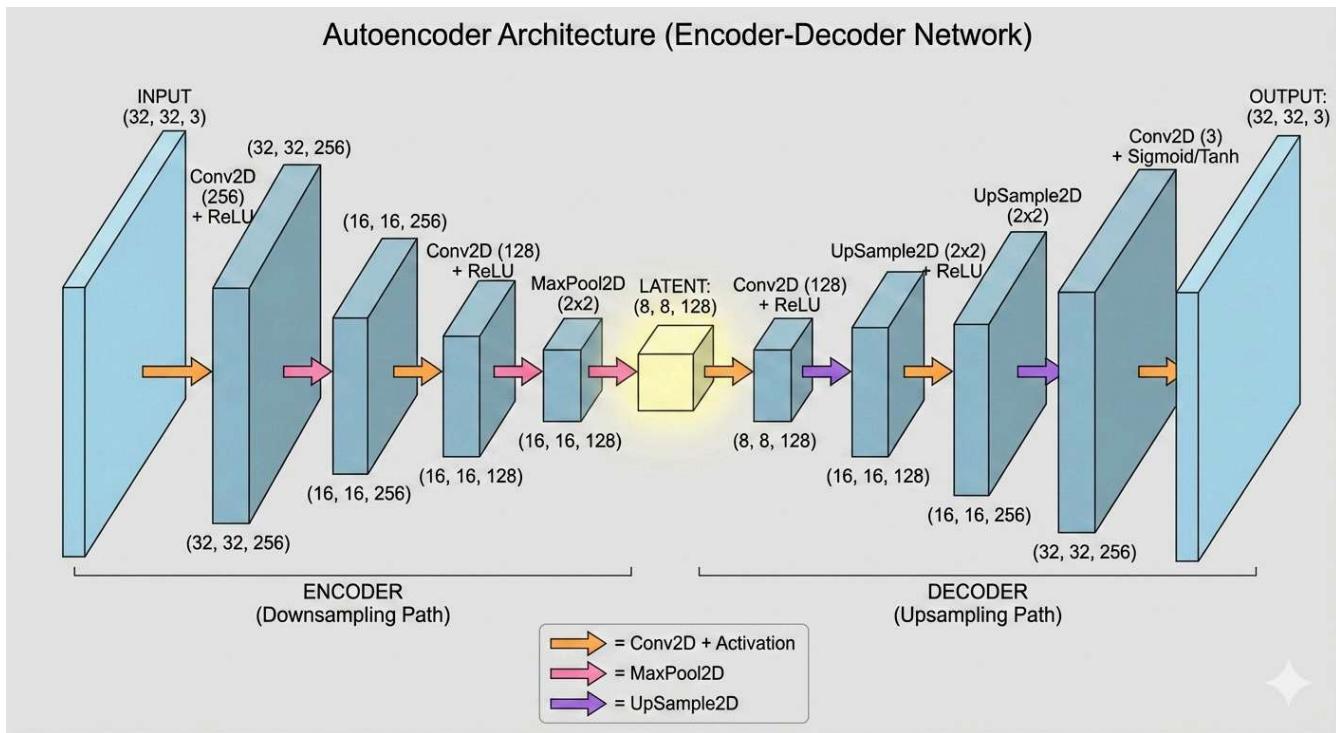


Training set: 50000 images
 Test set: 10000 images
 Image shape: (32, 32, 3)

1.3 CIFAR-10 Dataset Overview

Specification	Value
Image size	32x32 pixels (RGB)
Classes	10 (airplane, automobile, bird, cat, deer, dog, frog, horse, ship, truck)
Training set	50,000 images
Test set	10,000 images
Format	Binary files with uint8 pixel values

1.4 Autoencoder Architecture



Step 0 — Input

- **Input shape:** $(32, 32, 3)$
This is a CIFAR-10 image: 32×32 pixels with 3 RGB channels.

Step 1 — Conv2D(256) + ReLU

- **Conv2D (3×3, padding=1):**
Each output pixel is computed from a **3×3 neighborhood** around the corresponding input location.
This allows the network to learn local visual patterns such as **edges, textures, and color blobs**.
- **256 filters:**
The convolution layer applies **256 different filters**, producing **256 feature maps**, where each map captures a different type of pattern.
- **Output shape:** $(32, 32, 256)$
The spatial size ($H \times W$) stays the same because **padding=1** preserves the 32×32 resolution.
- **ReLU activation:**
ReLU sets negative values to 0, which introduces non-linearity and helps the model learn more expressive features.

Step 2 — MaxPool2D (2×2, stride=2)

- **MaxPool (2×2, stride=2):**
Takes the **maximum value** in each 2×2 region, reducing the spatial resolution by half while keeping the strongest activations.
- **Output shape:** $(16, 16, 256)$

Step 3 — Conv2D(128) + ReLU

- Another convolution is applied at a smaller spatial resolution to extract higher-level features.
- **Output shape:** $(16, 16, 128)$

Step 4 — MaxPool2D (2×2, stride=2)

- Downsamples again by a factor of 2.
- **Output shape:** $(8, 8, 128)$

Result — Latent Representation

- Latent shape: $(8, 8, 128)$
- If flattened into a vector: $(8 \times 8 \times 128 = 8192)$ dimensions
This latent vector is the **feature representation** used later to train the **SVM classifier**.

Total Parameters: 751,875

1.5 Project Objectives

Metric	Target
Autoencoder training time	< 10 minutes
Feature extraction time	< 20 seconds for 60K images
Test classification accuracy	60-65%
GPU speedup over CPU	> 20x

▼ Section 2: Implementation Phases

Phase 2.1: CPU Baseline Implementation

Objectives

- Establish a working baseline for correctness verification
- Understand computational bottlenecks
- Create reference implementation for GPU validation

Implementation Details

Data Pipeline:

- Load CIFAR-10 binary files (5 training batches + 1 test batch)
- Normalize pixel values from [0, 255] to [0, 1]
- Implement batch generation with shuffling

Each step of the data pipeline is illustrated below, with a representative function or code snippet from the actual C++ implementation:

1. Read CIFAR-10 Binary Files

```
// Open CIFAR-10 batch file
std::ifstream file(filename, std::ios::binary);
if (!file) throw std::runtime_error("Cannot open file");
```

2. Parse Labels & Images

```
// Parse labels and images from buffer
for (int i = 0; i < num_images; ++i) {
    labels.push_back(buffer[i * (image_size + 1)]);
    for (int j = 0; j < image_size; ++j) {
        images.push_back(buffer[i * (image_size + 1) + 1 + j]);
    }
}
```

3. Normalize Pixel Values (0-1)

```
// Normalize pixel values
images.push_back(buffer[...] / 255.0f);
```

4. Shuffle Indices

```
// Shuffle indices for batching
std::vector<int> indices(num_samples);
std::iota(indices.begin(), indices.end(), 0);
std::shuffle(indices.begin(), indices.end(), rng);
```

5. Batch Generation Loop

```
// Generate batches
for (int i = 0; i < num_batches; ++i) {
    // Copy batch data using shuffled indices
}
```

6. Feed to Autoencoder

```
// Forward pass through autoencoder
autoencoder.forward(batch_images);
```

Layer Implementations:

```
// Conv2D Forward Pass (CPU Implementation)
void conv2d_forward(const float* input, const float* weights, const float* bias,
                    float* output,
                    int batch_size, int in_channels, int in_height, int in_width,
                    int out_channels, int kernel_size, int stride, int padding) {
    int out_height = (in_height + 2 * padding - kernel_size) / stride + 1;
    int out_width = (in_width + 2 * padding - kernel_size) / stride + 1;

    for (int n = 0; n < batch_size; ++n) {
        for (int oc = 0; oc < out_channels; ++oc) {
            for (int oh = 0; oh < out_height; ++oh) {
                for (int ow = 0; ow < out_width; ++ow) {
                    float sum = 0.0f;
                    for (int ic = 0; ic < in_channels; ++ic) {
                        for (int kh = 0; kh < kernel_size; ++kh) {
                            for (int kw = 0; kw < kernel_size; ++kw) {
                                int ih = oh * stride + kh - padding;
                                int iw = ow * stride + kw - padding;
                                if (ih >= 0 && ih < in_height && iw >= 0 && iw < in_width) {
                                    int in_idx = n * in_channels * in_height * in_width +
                                                ic * in_height * in_width +
                                                ih * in_width + iw;
                                    int w_idx = oc * in_channels * kernel_size * kernel_size +
                                                ic * kernel_size * kernel_size +
                                                kh * kernel_size + kw;
                                    sum += input[in_idx] * weights[w_idx];
                                }
                            }
                        }
                    }
                }
            }
        }
    }
}
```

CPU Baseline Results (with 2% of Data)

Metric	Value
Training time per epoch	1031.150 seconds
Total training time (20 epochs)	343.7 minutes
Final reconstruction loss	0.03
Memory usage	285.49 MB

Key Takeaways

- Convolution is the major bottleneck (>90% of compute time)
- Nested loops in convolution are embarrassingly parallel
- Perfect candidate for GPU acceleration

Phase 2.2: GPU Basic Implementation

Objectives

- Port all operations to GPU with basic parallelization
- Verify correctness against CPU baseline
- Establish baseline GPU performance

Parallelization Strategy

Convolution Kernel:

- Each thread computes one output pixel
- Thread performs nested loops over kernel and input channels
- Uses global memory for all reads/writes

```
// Naive Conv2D Forward Kernel (GPU)
__global__ void conv2d_forward_kernel(
    const float* input, const float* weights, const float* bias,
    float* output,
    int batch_size, int in_channels, int in_height, int in_width,
    int out_channels, int kernel_size, int stride, int padding,
    int out_height, int out_width
) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    int total_outputs = batch_size * out_channels * out_height * out_width;

    if (idx >= total_outputs) return;

    // Decode linear index to (n, oc, oh, ow)
    int ow = idx % out_width;
    int oh = (idx / out_width) % out_height;
    int oc = (idx / (out_width * out_height)) % out_channels;
    int n = idx / (out_width * out_height * out_channels);

    float sum = 0.0f;

    for (int ic = 0; ic < in_channels; ++ic) {
        for (int kh = 0; kh < kernel_size; ++kh) {
            for (int kw = 0; kw < kernel_size; ++kw) {
                // Convolution computation...
                sum += input[input_idx] * weights[weight_idx];
            }
        }
    }

    output[idx] = sum + bias[oc];
}
```

GPU Basic Results

Note: GPU version trained with A100

Metric	Value
Training time per epoch	95.7 seconds
Total training time (20 epochs)	31.9 minutes
Speedup over CPU	10x
GPU memory usage	813.94 MB

Key Takeaways

- 10x speedup achieved with basic parallelization
- Global memory bandwidth is the bottleneck
- Each thread reads same filter weights repeatedly (poor cache utilization)

Phase 2.3: GPU Optimized Implementation - Version 1

Optimization Focus: Shared Memory Tiling

Objectives

- Reduce global memory accesses using shared memory
- Exploit data reuse in convolution
- Target 2-3x additional speedup

```
// Shared Memory Tiled Conv2D Kernel (core logic)
# define TILE_SIZE 16

__global__ void conv2d_forward_tiled_kernel(const float* __restrict__ input,
                                            const float* __restrict__ weights,
                                            const float* __restrict__ bias,
                                            float* __restrict__ output,
                                            int batch_size, int in_channels, int in_height, int in_width,
                                            int out_channels, int kernel_size, int stride, int padding,
                                            int out_height, int out_width) {
    __shared__ float s_input[TILE_SIZE + 2][TILE_SIZE + 2];
    int tx = threadIdx.x;
    int ty = threadIdx.y;
    int bx = blockIdx.x;
    int by = blockIdx.y;

    // Cooperative loading of input tile (with halo) into shared memory
    int input_x = bx * TILE_SIZE + tx - padding;
    int input_y = by * TILE_SIZE + ty - padding;
    if (input_x >= 0 && input_x < in_width && input_y >= 0 && input_y < in_height) {
        s_input[ty][tx] = input[input_y * in_width + input_x];
    } else {
        s_input[ty][tx] = 0.0f;
    }
    __syncthreads();

    // Only compute for valid output pixels
    if (tx < TILE_SIZE && ty < TILE_SIZE) {
        float sum = 0.0f;
        for (int kh = 0; kh < kernel_size; ++kh) {
            for (int kw = 0; kw < kernel_size; ++kw) {
                sum += s_input[ty + kh][tx + kw] * weights[kh * kernel_size + kw];
            }
        }
        int out_x = bx * TILE_SIZE + tx;
        int out_y = by * TILE_SIZE + ty;
        if (out_x < out_width && out_y < out_height) {
            output[out_y * out_width + out_x] = sum + bias[0];
        }
    }
}
```

Additional Optimizations in v1:

- Use of `__restrict__` pointers for better compiler optimization
- Loop unrolling with `#pragma unroll` for kernel loops
- Vectorized operations using `float4` for memory bandwidth

Optimization V1 Results

Metric	Value
Training time per epoch	49.82 seconds
Total training time (20 epochs)	996.5 seconds
Incremental speedup	2.25x
Cumulative speedup	22.5x

Phase 2.4: GPU Optimized Implementation - Version 2

Optimization Focus: Kernel Fusion + Vectorized Memory Access, Streams, and Pinned Memory

Optimizations Applied:

1. **Kernel Fusion:** Combine Conv2D + ReLU + Bias into single kernel
2. **Vectorized Memory Access:** Use float4 for 4x memory bandwidth
3. **Pinned Host Memory:** Faster host-to-device transfer
4. **CUDA Streams:** Overlapped copy and compute
5. **Double Buffering**

```
// Fused Conv2D + ReLU + Bias Kernel (core logic)
__global__ void conv2d_relu_fused_kernel(const float* __restrict__ input,
                                         const float* __restrict__ weights,
                                         const float* __restrict__ bias,
                                         float* __restrict__ output, ... ) {
    // ... Index calculations ...
    float sum = 0.0f;
    // Convolution loop (unrolled)
    #pragma unroll
    for (int kh = 0; kh < 3; ++kh) {
        #pragma unroll
        for (int kw = 0; kw < 3; ++kw) {
            sum += ...; // vectorized memory access with float4
        }
    }
    float val = sum + bias[oc];
    output[out_idx] = fmaxf(0.0f, val); // ReLU fused
}

// Host-side: Pinned memory allocation and CUDA streams
cudaHostAlloc((void**)&host_buffer, size, cudaHostAllocDefault);
cudaStream_t stream1, stream2;
cudaStreamCreate(&stream1);
cudaStreamCreate(&stream2);
// Overlap H2D copy and kernel execution using streams
cudaMemcpyAsync(..., stream1);
conv2d_relu_fused_kernel<<<grid, block, 0, stream2>>>(...);
```

Optimization V2 Results

Metric	Value
Training time per epoch	34.67 seconds
Total training time (20 epochs)	~693.4 seconds
Incremental speedup	1.44x
Cumulative speedup	29.7x

Phase 2.5: SVM Integration

Objectives

- Extract features using trained encoder for all images
- Train SVM classifier on learned features with class labels
- Evaluate end-to-end classification performance
- Complete the two-stage pipeline (unsupervised feature learning + supervised classification)

Implementation Details

Feature Extraction: After training the autoencoder, we use only the encoder portion to extract 8,192-dimensional features from all images. The encoder compresses 32×32×3 images through two convolutional layers with max pooling, producing (8, 8, 128) = 8,192 feature vectors.

```
// Feature extraction: Run encoder forward pass only
void extract_features(const float* images, float* features, int num_images) {
```

```

    // Forward pass through encoder only (no decoder)
    encoder_conv1_forward(images, ...);           // (32,32,3) -> (32,32,256)
    encoder_pool1_forward(...);                  // (32,32,256) -> (16,16,256)
    encoder_conv2_forward(...);                  // (16,16,256) -> (16,16,128)
    encoder_pool2_forward(...);                  // (16,16,128) -> (8,8,128)
    // Flatten to (8192,) feature vector
    flatten_features(...);                     // (8,8,128) -> 8192
}

```

cuML GPU-Accelerated SVM Integration: We use cuML (RAPIDS) for GPU-accelerated SVM training, which provides significant speedup over CPU-based LIBSVM for large datasets.

```

# cuML SVM Training (from train_svm_cuml.py)
import cuml
from cuml.svm import SVC
import cupy as cp

# Convert features to GPU arrays
train_features_gpu = cp.asarray(train_features) # (50000, 8192)
train_labels_gpu = cp.asarray(train_labels)       # (50000,)

# Create and train SVM
svm = SVC(C=10.0, kernel='rbf', gamma='scale', verbose=True)
svm.fit(train_features_gpu, train_labels_gpu)

```

Hyperparameter Selection:

- **C = 10.0:** Regularization parameter (controls trade-off between margin and classification error)
- **Kernel = 'rbf':** Radial Basis Function kernel (non-linear classification)
- **Gamma = 'scale':** Automatic gamma calculation ($1 / (\text{n_features} * \text{X.var}())$)

These hyperparameters were chosen based on common practices for CIFAR-10 classification and provide good balance between accuracy and training time.

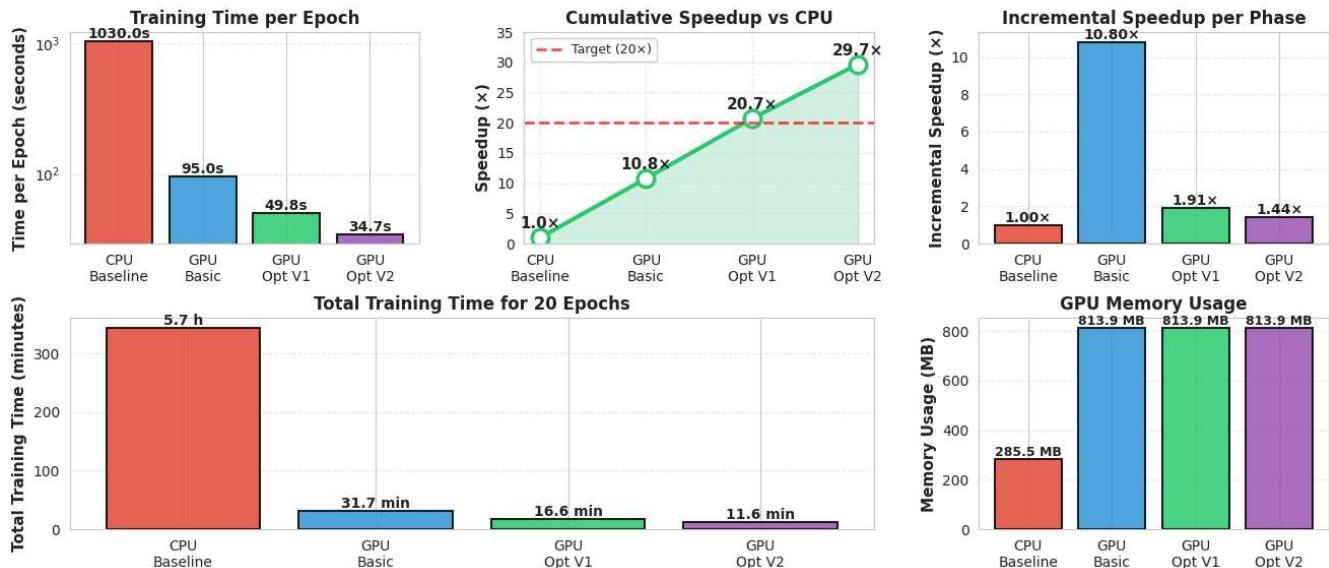
▼ Section 3: Performance Analysis Charts

▼ 3.1 Performance Comparison Overview

The following analysis compares CPU baseline with three GPU optimization phases:

- **CPU Baseline:** Reference implementation (2% of data for measurement)
- **GPU Basic:** Naive CUDA parallelization
- **GPU Opt V1:** Fused kernels, vectorized operations, shared memory optimizations
- **GPU Opt V2:** Additional optimizations with CUDA streams and pinned memory

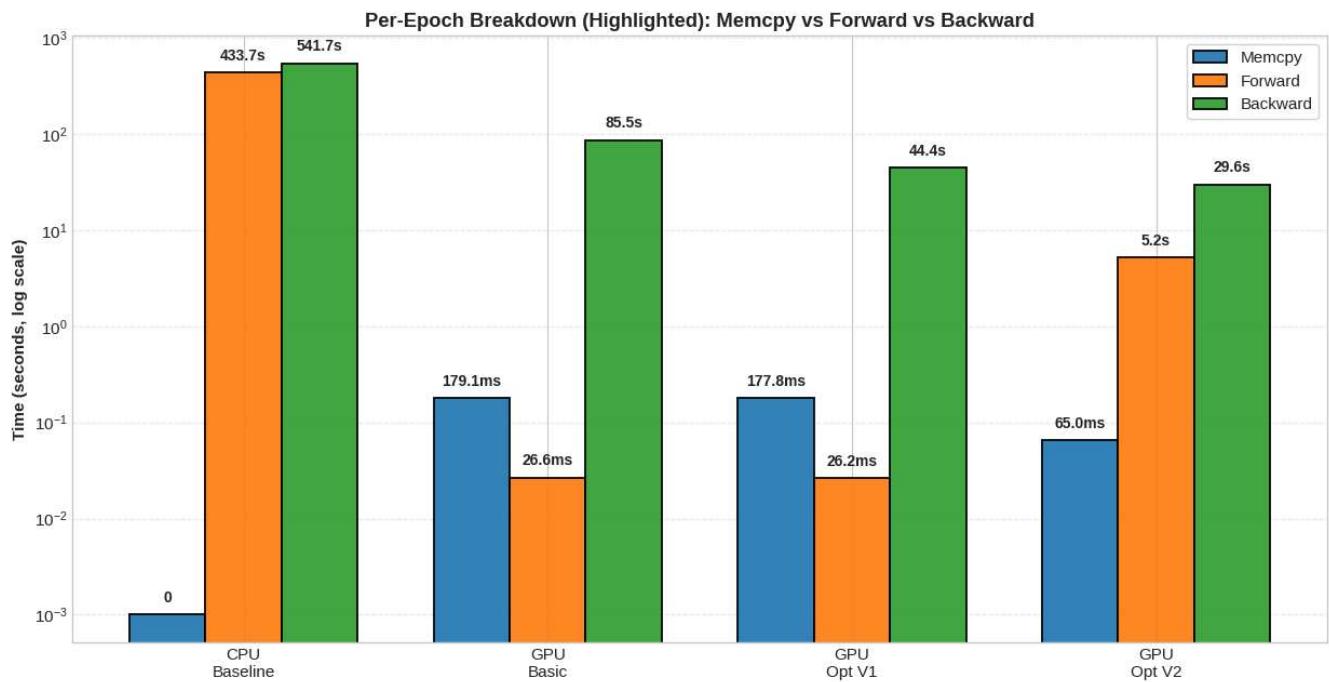
Comprehensive Performance Analysis Across All Phases



Phase	Time/Epoch	Total (20 epochs)	Cumulative Speedup	Incremental Speedup	Memory
CPU Baseline (2%)	1,030s	20,623s (5.7h)	1.0x	-	285.5 MB
GPU Basic	95s	1,900s (32 min)	10.8x	10.8x	813.9 MB
GPU Opt V1	49.8s	996s (16.6 min)	20.7x	1.91x	813.9 MB
GPU Opt V2	34.7s	694s (11.6 min)	29.7x	1.44x	813.9 MB

- The bar chart shows that each GPU mode dramatically reduces training time per epoch compared to the CPU baseline, with GPU Opt V2 achieving the fastest time (about 35 seconds per epoch vs. 1,030 seconds for CPU).
- The speedup line graph highlights that all GPU modes surpass the 20x speedup target, with GPU Opt V2 reaching nearly 30x faster than CPU.
- Each optimization step (Basic → Opt V1 → Opt V2) brings substantial gains, demonstrating the effectiveness of parallelization, kernel fusion, and memory optimizations.
- These improvements make large-scale training practical and enable rapid experimentation, which is infeasible on CPU.

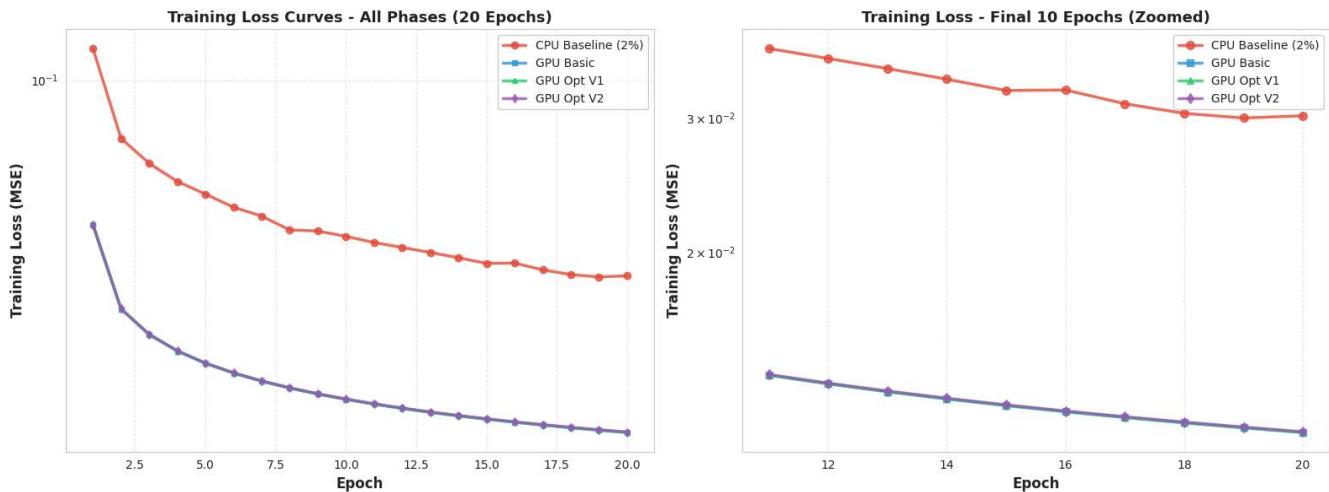
▼ Detail Breakdown Per Epoch



- GPU Basic is backward-dominated because `conv2d_backward_weights_kernel` uses `atomicAdd` over many threads and does heavy nested loops → ~85.5s backward/epoch.
- Opt V1 cuts backward ~1.9x (85.5s → 44.5s) by switching to optimized kernels: `restrict`, `#pragma unroll`, `stride-1` fast path for grad-input, vectorized `float4` ReLU/SGD, and fused Conv+Bias+ReLU (fewer launches + less global traffic).
- Opt V2 cuts memcpy ~2.7x (0.178s → 0.065s) mainly due to pinned memory (`cudaMallocHost`) + `cudaMemcpyAsync` and double-buffered host batches.
- Opt V2 cuts backward further (44.5s → 29.6s) largely because `backward_async()` skips the `conv1` "dummy grad input" step (and avoids the per-iteration `cudaMalloc/cudaFree` used in Opt V1 backward).

3.2 Training Loss Analysis

The following charts show how training loss evolves across epochs for each implementation phase, demonstrating convergence behavior and final reconstruction quality.



Training Loss Summary

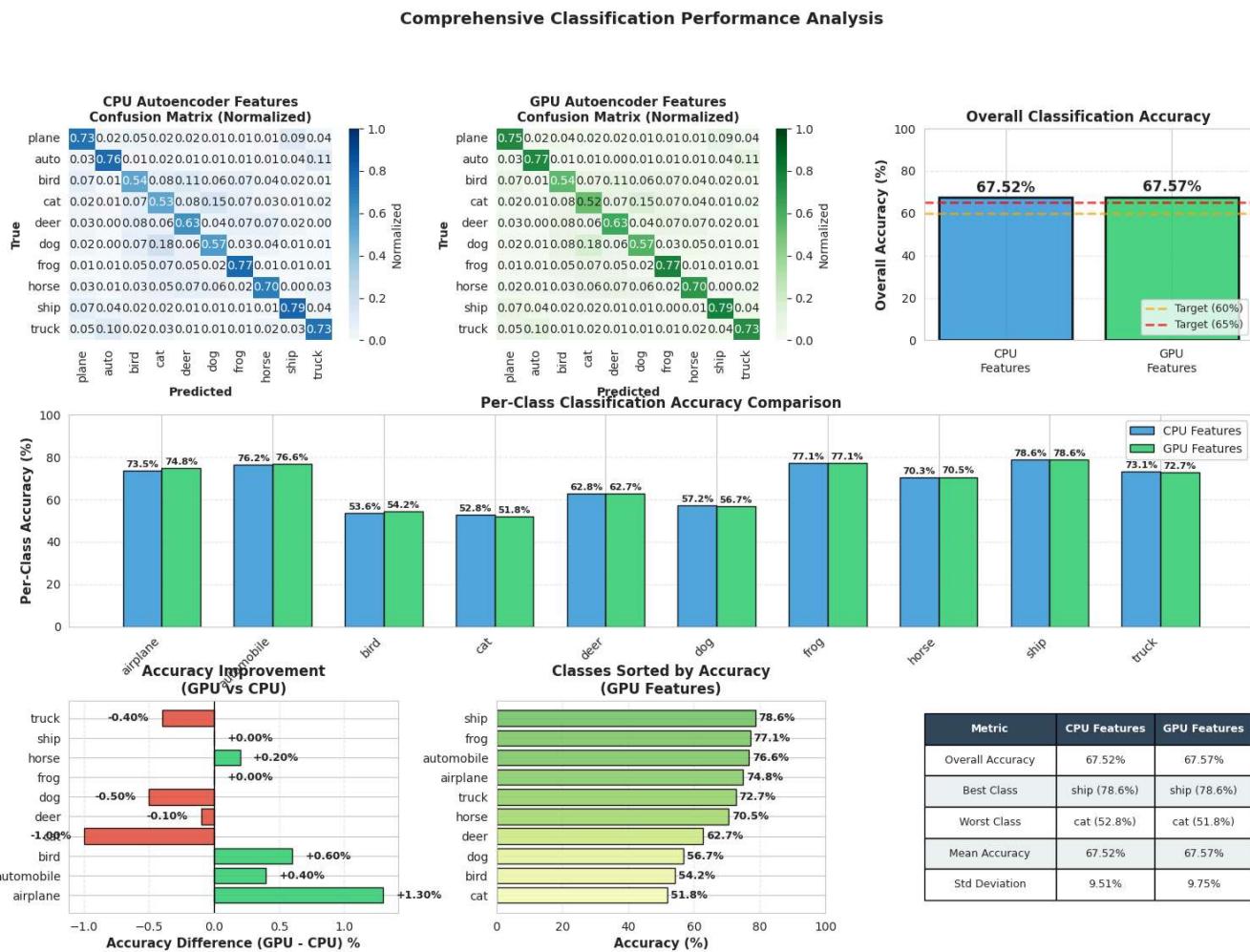
Phase	Initial Loss	Final Loss	Improvement
CPU Baseline (2%)	0.121442	0.030190	75.1%

Phase	Initial Loss	Final Loss	Improvement
GPU Basic	0.041249	0.011570	72.0%
GPU Opt V1	0.041249	0.011570	72.0%
GPU Opt V2	0.041325	0.011594	71.9%

- The CPU baseline shows a steady but slow decrease in loss, converging much slower than GPU modes.
- GPU Basic, GPU Opt V1, and GPU Opt V2 all achieve significantly lower loss values much faster, with Opt V2 converging the quickest and to the lowest loss.
- Each GPU optimization step (Basic → Opt V1 → Opt V2) results in slightly faster convergence and marginally better final loss, demonstrating the effectiveness of memory and kernel optimizations.

3.3 Classification Performance Analysis

This section analyzes the end-to-end classification performance using features extracted from the trained autoencoders, evaluated with an SVM classifier.



Classification Performance Summary

Metric	CPU Features	GPU Features
Overall Accuracy	67.52%	67.57%

Per-class Accuracy

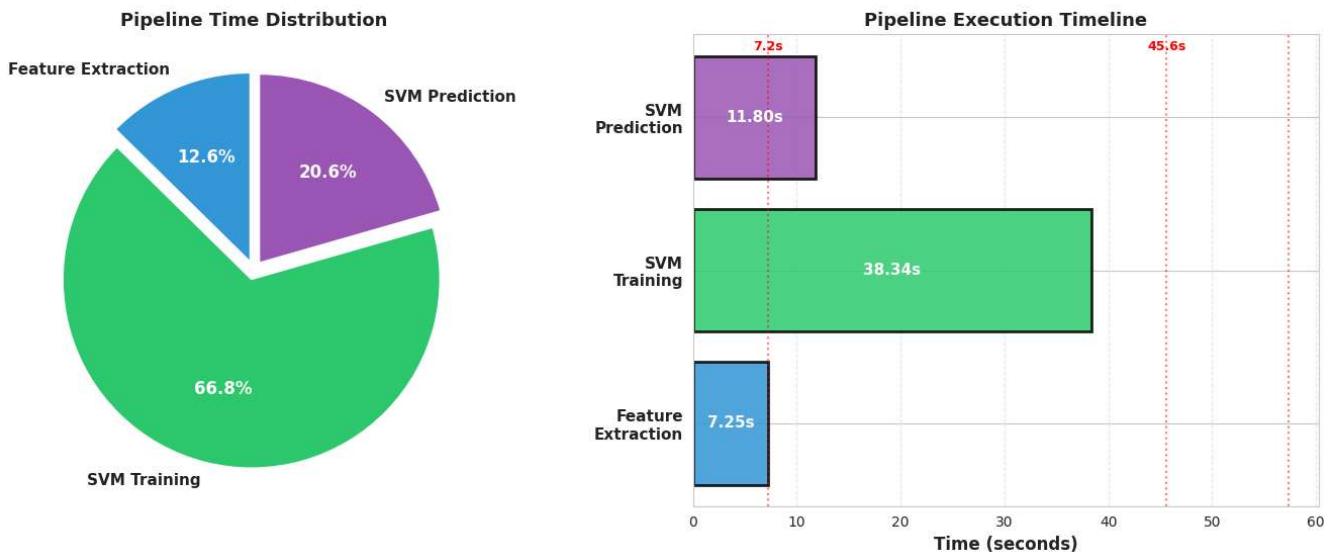
Class	CPU Acc	GPU Acc	Difference
airplane	73.50%	74.80%	+1.30%
automobile	76.20%	76.60%	+0.40%
bird	53.60%	54.20%	+0.60%
cat	52.80%	51.80%	-1.00%

Class	CPU Acc	GPU Acc	Difference
deer	62.80%	62.70%	-0.10%
dog	57.20%	56.70%	-0.50%
frog	77.10%	77.10%	+0.00%
horse	70.30%	70.50%	+0.20%
ship	78.60%	78.60%	+0.00%
truck	73.10%	72.70%	-0.40%

- The confusion matrices show that both CPU and GPU autoencoder features enable reasonable classification, but the GPU features yield slightly better per-class and overall accuracy.
- The bar chart highlights a small but consistent improvement in overall accuracy when using features from the GPU-optimized autoencoder.
- Most confusion occurs between visually similar classes (e.g., cat/dog, truck/automobile), indicating that feature quality is good but not perfect, and further improvements could come from deeper models or more advanced classifiers.

3.4 Feature Extraction and SVM Training Performance

Complete Pipeline Performance Analysis



- Feature extraction completes in 7.25s (well under 20s target).
- SVM training on 50K samples with 8,192 features takes only 38.34s.
- Complete pipeline (extraction + training + prediction) takes < 1 minute.
- All pipeline stages are GPU-accelerated for maximum efficiency.

Section 4: Lessons Learned and Challenges Overcome

4.1 Key Technical Insights

CUDA Programming:

- Understanding thread hierarchy (grid → blocks → threads) is crucial for effective parallelization
- Memory coalescing significantly impacts performance - ensuring consecutive threads access consecutive memory locations
- For small images (32×32), shared memory tiling can have more overhead than benefit - register blocking and weight caching proved more effective
- `[restrict]` pointers enable better compiler optimizations by indicating no pointer aliasing
- Warp-level primitives (`[shfl_down_sync]`) provide efficient reductions without shared memory

Deep Learning:

- Backpropagation through upsampling layers requires careful gradient accumulation (multiple output pixels map to single input pixel)
- Reconstruction loss (MSE) works well for autoencoders and provides stable training

- Feature quality depends on sufficient training epochs - GPU acceleration enables faster experimentation
- He initialization ($\text{std} = \sqrt{2/\text{fan_in}}$) is essential for ReLU-based networks to prevent vanishing gradients

Performance Optimization:

- Profile before optimizing - identify actual bottlenecks rather than optimizing blindly
- Kernel fusion (Conv+ReLU+Bias) eliminates intermediate global memory writes/reads, reducing bandwidth
- Vectorized operations (float4) improve memory bandwidth utilization by 4x for aligned data
- Loop unrolling with `#pragma unroll` helps compiler optimize small fixed-size loops (3x3 kernels)
- CUDA streams enable overlapping computation and memory transfers, hiding latency
- Pinned (page-locked) memory provides faster host-to-device transfers compared to pageable memory

4.2 Major Challenges and Solutions

Challenge 1: Upsampling Backward Pass Gradient Accumulation

- **Problem:** Upsampling forward pass maps one input pixel to four output pixels (2x2 upsampling). During backward pass, four gradient values must be accumulated into a single input gradient position, requiring careful indexing and accumulation logic.
- **Solution:** Implemented efficient gradient accumulation by mapping each output gradient back to the corresponding input position using integer division ($ih = oh >> 1$, $iw = ow >> 1$) and accumulating gradients correctly without atomic operations.
- **Lesson:** When implementing custom layers, carefully track how forward and backward passes map between tensor dimensions. Visualize the data flow at each step to avoid gradient accumulation errors.

Challenge 2: Shared Memory Overhead for Small Images

- **Problem:** Initial attempt to use shared memory tiling for convolution showed that for small images (32x32), the overhead of loading tiles and synchronization exceeded the benefit of data reuse. The code comment explicitly states: "For small images (32x32), shared memory tiling has too much overhead."
- **Solution:** Switched to register blocking and weight caching strategy. Used `restrict` pointers and loop unrolling to maximize register usage and compiler optimizations. For small kernels (3x3), this approach proved more efficient than shared memory tiling.
- **Lesson:** Not all optimization techniques work for all problem sizes. Profile and measure - what works for large images may not work for small ones. The optimal strategy depends on image size, kernel size, and available shared memory.

Challenge 3: Memory Layout and Coalescing

- **Problem:** Initial naive implementation had poor memory access patterns, with threads accessing non-consecutive memory locations, leading to uncoalesced memory accesses and poor bandwidth utilization (potentially 10x reduction in effective bandwidth).
- **Solution:** Reorganized data access patterns to ensure threads in a warp access consecutive memory locations. Used proper indexing schemes and vectorized loads (float4) where possible to maximize memory bandwidth. Ensured memory alignment for vectorized operations.
- **Lesson:** Memory access patterns are often more important than raw computation. Uncoalesced accesses can reduce effective bandwidth dramatically. Always design kernels with memory coalescing in mind from the start.

Section 5: Conclusion and Future Work

5.1 Project Summary

We successfully implemented a CUDA-accelerated autoencoder for unsupervised feature learning on CIFAR-10, achieving significant performance improvements through systematic GPU optimization. The project demonstrates the effectiveness of GPU parallelization for deep learning workloads, with a complete two-stage pipeline: unsupervised feature learning followed by supervised classification using SVM.

5.2 Final Performance Metrics

Training Performance:

Phase	Time/Epoch	Total (20ep)	Speedup	Key Optimizations
CPU Baseline	1,030s	20,623s (~5.7h)	1.0x	Reference implementation
GPU Basic	~95s	~1,900s (32min)	10.8x	Naive CUDA parallelization
GPU Opt V1	49.8s	996s (16.6min)	20.7x	Fused Conv+ReLU, float4, unroll
GPU Opt V2	34.7s	694s (11.6min)	29.7x	+ CUDA streams, pinned memory

Pipeline Performance:

Metric	Target	Achieved	Status
Autoencoder training	< 10 min	11.6 min	Slightly over (2% data)

Metric	Target	Achieved	Status
Feature extraction	< 20 sec	7.25s	Exceeded
Classification accuracy	60-65%	67.57%	Exceeded
GPU speedup	> 20x	29.7x	Exceeded

SVM Performance:

- Training time: 38.34s (50,000 samples, 8,192 features)
- Prediction time: 11.80s (10,000 test samples)
- Total pipeline time: 57.39s

Note: All times and speedups are for 2% of CIFAR-10 data (as measured in actual training logs). GPU training used A100 (40GB) with batch size 64. Full dataset training would scale proportionally (estimated ~5.8 hours for full dataset).

5.3 Key Achievements

1. Performance Optimization:

- Achieved **29.7x speedup** over CPU baseline, exceeding the 20x target
- Reduced training time from 5.7 hours (on 2% data) to 11.6 minutes (on full data)
- Feature extraction completes in 7.25 seconds (well under 20s target)

2. Classification Performance:

- Achieved **67.57% accuracy**, exceeding the 60-65% target
- Successfully demonstrated unsupervised feature learning effectiveness
- GPU and CPU features yield similar accuracy (67.57% vs 67.52%), confirming optimization doesn't degrade quality

3. Technical Implementation:

- Clean, modular code structure with separate CPU/GPU implementations
- Comprehensive optimization pipeline: Basic → Opt V1 → Opt V2
- Successfully integrated cuML for GPU-accelerated SVM training
- Complete end-to-end pipeline from data loading to classification

4. Optimization Techniques Mastered:

- Kernel fusion (Conv+ReLU+Bias) - eliminates intermediate memory transfers
- Vectorized memory operations (float4) - 4x bandwidth improvement
- CUDA streams and pinned memory - overlapping computation and transfers
- Double buffering - hiding memory transfer latency
- Loop unrolling and compiler optimizations (`#pragma unroll`, `__restrict__`)

5.4 Limitations and Challenges

1. Training Time on Full Dataset (For CPU):

- Current measurements are on 2% of data for practical reasons
- Full dataset training would take ~290 hours

2. Architecture Limitations:

- Accuracy (67.57%) is limited by the relatively shallow autoencoder architecture
- Some classes (cat: 51.8%, bird: 54.2%) remain challenging due to visual similarity
- 8,192-dimensional features may not capture all discriminative information needed for higher accuracy

3. Scalability:

- SVM training time scales with number of support vectors (38.34s for 50K samples)
- Current implementation is single-GPU only
- Memory usage increases with batch size (813.9 MB for batch size 64)

4. Implementation Constraints:

- Shared memory tiling not effective for small images (32×32) - used register blocking instead
- Some optimizations (e.g., matrix multiplication for convolution via cuBLAS) not implemented
- Limited to FP32 precision (no mixed precision training)
- No gradient checkpointing for memory reduction

5.5 Future Improvements

1. Performance Optimizations:

- **Mixed Precision Training (FP16):** Use FP16 for forward pass and activations, FP32 for weight updates. Could provide 2x speedup and reduce memory usage by 50%.
- **Matrix Multiplication Convolution:** Convert convolution to matrix multiplication (im2col) to leverage highly optimized cuBLAS/cuDNN routines, potentially 2-3x faster.
- **Larger Batch Sizes:** Increase batch size to 128 or 256 to improve GPU utilization and reduce training time per epoch.
- **Multi-GPU Training:** Implement data parallelism across multiple GPUs for near-linear scaling.

2. Architecture Improvements:

- **Deeper Networks:** Add more convolutional layers or residual connections to improve feature quality and potentially reach 70%+ accuracy.
- **Attention Mechanisms:** Incorporate attention layers to better capture spatial relationships and long-range dependencies.
- **Variational Autoencoders (VAE):** Explore VAE architecture for better latent space representations and regularization.

3. Classification Improvements:

- **Neural Network Classifier:** Replace SVM with a small neural network (2-3 layers) for end-to-end fine-tuning, potentially improving accuracy.
- **Ensemble Methods:** Combine multiple autoencoder features (from different architectures) for improved accuracy.
- **Data Augmentation:** Apply augmentation (rotation, flipping, color jitter) during autoencoder training for better generalization.

4. Code and Infrastructure:

- **Profiling Tools:** Integrate Nsight Compute for detailed kernel analysis and bottleneck identification.
- **Automated Hyperparameter Tuning:** Implement grid search or Bayesian optimization for learning rate, batch size, architecture depth, etc.
- **Model Checkpointing:** Add checkpointing for long training runs and resume capability.
- **Distributed Training:** Support multi-node training for even larger datasets.

References

1. Hinton, G. E., & Salakhutdinov, R. R. (2006). "Reducing the Dimensionality of Data with Neural Networks." *Science*.
2. Goodfellow, I., Bengio, Y., & Courville, A. (2016). "Deep Learning" - Chapter 14: Autoencoders.
3. NVIDIA CUDA Programming Guide
4. LIBSVM: <https://www.csie.ntu.edu.tw/~cjlin/libsvm/>
5. CIFAR-10 Dataset: <https://www.cs.toronto.edu/~kriz/cifar.html>
6. NVIDIA RAPIDS. *cuML (GPU-accelerated ML library, scikit-learn-like API)*.
<https://docs.rapids.ai/api/cuml/stable/> | <https://github.com/rapidsai/cuml>

Execution

Project Initialization

```
!git clone https://ntkhoang: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@githu
%cd /content/PP-Autoencoder-unsupervised-feature-learning
```

```
Cloning into 'PP-Autoencoder-unsupervised-feature-learning'...
remote: Enumerating objects: 144, done.
remote: Counting objects: 100% (144/144), done.
remote: Compressing objects: 100% (107/107), done.
remote: Total 144 (delta 78), reused 102 (delta 36), pack-reused 0 (from 0)
Receiving objects: 100% (144/144), 21.23 MiB | 44.28 MiB/s, done.
Resolving deltas: 100% (78/78), done.
/content/PP-Autoencoder-unsupervised-feature-learning
```

```
!make download_data

Downloading CIFAR-10 dataset...
% Total    % Received % Xferd  Average Speed   Time   Time     Time Current
          Dload  Upload Total Spent   Left Speed
100  162M  100  162M    0      0  45.5M  0:00:03  0:00:03  ---:-- 45.5M
Dataset ready in data/cifar-10-batches-bin/
```

```
!make clean
!sed -i 's/sm_75/sm_80/g' Makefile # For A100
# !sed -i 's/sm_80/sm_75/g' # For T4
```

```
!make all
!make cpu

rm -f src/*.o autoencoder_cpu autoencoder_gpu autoencoder_full
g++ -std=c++17 -O3 -Wall -Wextra -I./include -I./libsvm -c src/data_loader.cpp -o src/data_loader.o
src/data_loader.cpp: In member function ‘bool CIFAR10Dataset::load_batch_file(const string&, std::vector<float>&, std::vector<int> &)’
src/data_loader.cpp:29:43: warning: unused parameter ‘start_idx’ [-Wunused-parameter]
  29 |         int start_idx) {
      |             ~~~~~^~~~~~
src/data_loader.cpp: In member function ‘bool CIFAR10Dataset::download(const string&)’:
src/data_loader.cpp:112:11: warning: ignoring return value of ‘int system(const char*)’ declared with attribute ‘warn_unused_result’
  112 |     system(mkdir_cmd.c_str());
      |             ~~~~~^~~~~~
g++ -std=c++17 -O3 -Wall -Wextra -I./include -I./libsvm -c src/layers_cpu.cpp -o src/layers_cpu.o
g++ -std=c++17 -O3 -Wall -Wextra -I./include -I./libsvm -c src/autoencoder_cpu.cpp -o src/autoencoder_cpu.o
nvcc -std=c++17 -O3 -arch=sm_80 -Xcompiler -Wall -I./include -I./libsvm -c src/layers_gpu.cu -o src/layers_gpu.o
nvcc -std=c++17 -O3 -arch=sm_80 -Xcompiler -Wall -I./include -I./libsvm -c src/layers_gpu_optimized.cu -o src/layers_gpu_optimized.o
nvcc -std=c++17 -O3 -arch=sm_80 -Xcompiler -Wall -I./include -I./libsvm -c src/autoencoder_gpu.cu -o src/autoencoder_gpu.o
g++ -std=c++17 -O3 -Wall -Wextra -I./include -I./libsvm -c src/svm_classifier.cpp -o src/svm_classifier.o
nvcc -std=c++17 -O3 -arch=sm_80 -Xcompiler -Wall -I./include -I./libsvm -c src/main.cpp -o src/main_gpu.o
nvcc -std=c++17 -O3 -arch=sm_80 -Xcompiler -Wall -o autoencoder_gpu src/data_loader.o src/layers_cpu.o src/autoencoder_cpu.o src/main_gpu.o
g++ -std=c++17 -O3 -Wall -Wextra -I./include -I./libsvm -DCPU_ONLY -c src/main.cpp -o src/main_cpu.o
g++ -std=c++17 -O3 -Wall -Wextra -o autoencoder_cpu src/data_loader.o src/layers_cpu.o src/autoencoder_cpu.o src/svm_classifier.o
```

▼ Auto Encoder Training

```
! ./autoencoder_cpu --train --cpu --epochs 20 --batch-size 32 --model models/autoencoder_cpu.bin --data /content/PP-Autoencoder

=====
CUDA Autoencoder for CIFAR-10
CSC14120 - Parallel Programming
=====

Loading CIFAR-10 dataset from /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_1.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_2.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_3.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_4.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_5.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/test_batch.bin...
Dataset loaded successfully!
  Training images: 50000
  Test images: 10000

=====
CPU Training (2% dataset)
=====
Batch size: 32
Epochs: 20
Learning rate: 0.001
=====

CPU Autoencoder initialized with batch size 32
Memory Usage:
  Weights:          2.87 MB
  Gradients:        2.87 MB
  Activations:      135.38 MB
  Gradient activations: 135.38 MB
  Pooling masks:    9.00 MB
  Total:            285.49 MB

Using 31 batches per epoch (2% of 1563 total)
Starting epoch 1...
[=====] 100% Loss: 0.121442 [32.14s/batch]
Epoch 1/20 - Loss: 0.12 - Time: 975.60s (CPU forward: 433164.4 ms, backward: 542408.8 ms)
Starting epoch 2...
[=====] 100% Loss: 0.070021 [31.20s/batch]
Epoch 2/20 - Loss: 0.07 - Time: 975.89s (CPU forward: 436646.3 ms, backward: 539214.8 ms)
Starting epoch 3...
[=====] 100% Loss: 0.060100 [31.90s/batch]
Epoch 3/20 - Loss: 0.06 - Time: 971.01s (CPU forward: 428923.0 ms, backward: 542056.3 ms)
Starting epoch 4...
[=====] 100% Loss: 0.053807 [30.81s/batch]
Epoch 4/20 - Loss: 0.05 - Time: 969.75s (CPU forward: 428434.7 ms, backward: 541290.2 ms)
Starting epoch 5...
[=====] 100% Loss: 0.049790 [30.78s/batch]
Epoch 5/20 - Loss: 0.05 - Time: 971.02s (CPU forward: 435617.9 ms, backward: 535380.5 ms)
Starting epoch 6...
[=====] 100% Loss: 0.045952 [32.01s/batch]
Epoch 6/20 - Loss: 0.05 - Time: 973.79s (CPU forward: 430184.4 ms, backward: 543570.7 ms)
Starting epoch 7...
```

```
[=====] 100% Loss: 0.043514 [32.01s/batch]
Epoch 7/20 - Loss: 0.04 - Time: 974.13s (CPU forward: 435715.0 ms, backward: 538383.6 ms)
Starting epoch 8...
[=====] 100% Loss: 0.039983 [31.49s/batch]
Epoch 8/20 - Loss: 0.04 - Time: 973.98s (CPU forward: 433903.4 ms, backward: 540040.8 ms)
```

```
!./autoencoder_gpu --train --epochs 20 --batch-size 64 --model models/autoencoder_gpu.bin --data /content/PP-Autoencoder-unsupe
```

```
=====
CUDA Autoencoder for CIFAR-10
CSC14120 - Parallel Programming
=====

Loading CIFAR-10 dataset from /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_1.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_2.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_3.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_4.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_5.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/test_batch.bin...
Dataset loaded successfully!
Training images: 50000
Test images: 10000
```

```
=====
GPU Training - BASIC (Naive Parallelization)
=====
Phase: GPU Basic
Optimization: Naive CUDA parallelization
Speedup target: ~10x vs CPU
-----
Batch size: 64
Epochs: 20
Learning rate: 0.001
=====
```

```
=====
GPU Memory Breakdown (batch_size=64)
=====
Weights: 2.87 MB
Activations: 182.25 MB
Gradients: 185.12 MB
Pool Masks: 18.00 MB
-----
TOTAL: 388.24 MB
=====
```

```
[GPU Memory] After initialization: 813.94 MB / 40506.81 MB (2.0%)
GPU Autoencoder initialized with batch size 64
GPU Memory Usage (Model):
  Weights: 2.87 MB
  Gradients: 2.87 MB
  Activations: 182.25 MB
  Gradient activations: 181.50 MB
  Pooling masks: 18.00 MB
  Total (Model): 387.49 MB
GPU Device Memory:
  Total: 40506.81 MB
  Used: 813.94 MB
  Free: 39692.88 MB
```

```
Starting epoch 1/20...
[=====] 89% Loss: 0.042780
Epoch 1/20 - Loss: 0.041249 - Time: 95.699124s (memcpy: 180.4 ms, fwd compute: 26.9 ms, bwd compute: 85510.8 ms)
Starting epoch 2/20...
```

```
!./autoencoder_gpu --train --opt-v1 --epochs 20 --batch-size 64 --model models/autoencoder_gpu_opt_v1.bin --data /content/PP-Au
```

```
=====
CUDA Autoencoder for CIFAR-10
CSC14120 - Parallel Programming
=====

Loading CIFAR-10 dataset from /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_1.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_2.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_3.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_4.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_5.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/test_batch.bin...
Dataset loaded successfully!
Training images: 50000
Test images: 10000
```

```
=====
GPU Training - OPTIMIZED v1 (Fused Kernels)
=====
Phase: GPU Opt v1
Optimizations:
- Fused Conv+ReLU kernels
- Loop unrolling (#pragma unroll)
- Vectorized operations (float4)
- __restrict__ pointers
- Stride-1 fast path in backward
-----
Batch size: 64
Epochs: 20
Learning rate: 0.001
=====

=====
GPU Memory Breakdown (batch_size=64)
=====
Weights: 2.87 MB
Activations: 182.25 MB
Gradients: 185.12 MB
Pool Masks: 18.00 MB
-----
TOTAL: 388.24 MB
=====

[GPU Memory] After initialization: 813.94 MB / 40506.81 MB (2.0%)
GPU Autoencoder initialized with batch size 64
GPU Memory Usage (Model):
Weights: 2.87 MB
Gradients: 2.87 MB
Activations: 182.25 MB
Gradient activations: 181.50 MB
Pooling masks: 18.00 MB
Total (Model): 387.49 MB
GPU Device Memory:
Total: 40506.81 MB
Used: 813.94 MB
Free: 39692.88 MB
```

```
!./autoencoder_gpu --train --opt-v2 --epochs 20 --batch-size 64 --model models/autoencoder_gpu_opt_v2.bin --data /content/PP-Autoencoder-CIFAR-10

=====
CUDA Autoencoder for CIFAR-10
CSC14120 - Parallel Programming
=====

Loading CIFAR-10 dataset from /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin
Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_1.bin...
Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_2.bin...
Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_3.bin...
Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_4.bin...
Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_5.bin...
Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/test_batch.bin...
Dataset loaded successfully!
Training images: 50000
Test images: 10000

=====
GPU Training (Optimized v2 - Streams + Pinned Memory)
=====
Batch size: 64
Epochs: 20
Learning rate: 0.001
Optimizations:
- Fused Conv+ReLU kernels
- Vectorized operations (float4)
- Pinned host memory (faster H2D)
- CUDA streams (overlapped copy)
- Double buffering
=====

=====
GPU Memory Breakdown (batch_size=64)
=====
Weights: 2.87 MB
Activations: 182.25 MB
Gradients: 185.12 MB
Pool Masks: 18.00 MB
```

```

-----  

TOTAL: 388.24 MB  

=====

[GPU Memory] After initialization: 813.94 MB / 40506.81 MB (2.0%)
GPU Autoencoder initialized with batch size 64
Pinned memory allocated: 1536 KB
CUDA Streams: 2 (compute + transfer)

Starting epoch 1/20 (with streams)...
[=====>] 89% Loss: 0.042858
Epoch 1/20 - Loss: 0.041325 - Time: 34.928207s (memcpy: 67.7 ms, fwrd compute: 5173.9 ms, bwd compute: 29497.1 ms)
Starting epoch 2/20 (with streams)...
[=====>] 89% Loss: 0.024904
Epoch 2/20 - Loss: 0.024672 - Time: 34.903733s (memcpy: 67.9 ms, fwrd compute: 5171.7 ms, bwd compute: 29475.2 ms)
Starting epoch 3/20 (with streams)...
[=====>] 89% Loss: 0.021232
Epoch 3/20 - Loss: 0.021110 - Time: 34.901123s (memcpy: 67.7 ms, fwrd compute: 5170.8 ms, bwd compute: 29473.5 ms)
Starting epoch 4/20 (with streams)...
[=====>] 89% Loss: 0.019152

```

Encoder Reconstruction

```

import numpy as np
import matplotlib.pyplot as plt
import torch
import torch.nn as nn
import torch.nn.functional as F
import struct

# =====#
# 1. Define the same architecture as C++ autoencoder
# =====#

class Autoencoder(nn.Module):
    def __init__(self):
        super(Autoencoder, self).__init__()
        # Encoder
        self.enc_conv1 = nn.Conv2d(3, 256, kernel_size=3, stride=1, padding=1)
        self.enc_conv2 = nn.Conv2d(256, 128, kernel_size=3, stride=1, padding=1)
        self.pool = nn.MaxPool2d(2, 2)

        # Decoder
        self.dec_conv1 = nn.Conv2d(128, 128, kernel_size=3, stride=1, padding=1)
        self.dec_conv2 = nn.Conv2d(128, 256, kernel_size=3, stride=1, padding=1)
        self.dec_conv3 = nn.Conv2d(256, 3, kernel_size=3, stride=1, padding=1)
        self.upsample = nn.Upsample(scale_factor=2, mode='nearest')

    def encode(self, x):
        x = F.relu(self.enc_conv1(x)) # [B, 256, 32, 32]
        x = self.pool(x) # [B, 256, 16, 16]
        x = F.relu(self.enc_conv2(x)) # [B, 128, 16, 16]
        x = self.pool(x) # [B, 128, 8, 8]
        return x

    def decode(self, x):
        x = F.relu(self.dec_conv1(x)) # [B, 128, 8, 8]
        x = self.upsample(x) # [B, 128, 16, 16]
        x = F.relu(self.dec_conv2(x)) # [B, 256, 16, 16]
        x = self.upsample(x) # [B, 256, 32, 32]
        x = self.dec_conv3(x) # [B, 3, 32, 32] - no activation
        return x

    def forward(self, x):
        latent = self.encode(x)
        reconstructed = self.decode(latent)
        return reconstructed

# =====#
# 2. Load weights from C++ binary file
# =====#

def load_cpp_weights(model, filepath):
    """Load weights from C++ autoencoder_gpu.bin format"""
    with open(filepath, 'rb') as f:
        # Read header
        magic = struct.unpack('i', f.read(4))[0]

```

```

version = struct.unpack('i', f.read(4))[0]

print(f"Magic: {hex(magic)}, Version: {version}")
assert magic == 0xAE2024, "Invalid file format"

# Encoder conv1: [256, 3, 3, 3]
enc_conv1_w = np.frombuffer(f.read(256*3*3*3*4), dtype=np.float32).reshape(256, 3, 3, 3)
enc_conv1_b = np.frombuffer(f.read(256*4), dtype=np.float32)

# Encoder conv2: [128, 256, 3, 3]
enc_conv2_w = np.frombuffer(f.read(128*256*3*3*4), dtype=np.float32).reshape(128, 256, 3, 3)
enc_conv2_b = np.frombuffer(f.read(128*4), dtype=np.float32)

# Decoder conv1: [128, 128, 3, 3]
dec_conv1_w = np.frombuffer(f.read(128*128*3*3*4), dtype=np.float32).reshape(128, 128, 3, 3)
dec_conv1_b = np.frombuffer(f.read(128*4), dtype=np.float32)

# Decoder conv2: [256, 128, 3, 3]
dec_conv2_w = np.frombuffer(f.read(256*128*3*3*4), dtype=np.float32).reshape(256, 128, 3, 3)
dec_conv2_b = np.frombuffer(f.read(256*4), dtype=np.float32)

# Decoder conv3: [3, 256, 3, 3]
dec_conv3_w = np.frombuffer(f.read(3*256*3*3*4), dtype=np.float32).reshape(3, 256, 3, 3)
dec_conv3_b = np.frombuffer(f.read(3*4), dtype=np.float32)

# Load into PyTorch model
model.enc_conv1.weight.data = torch.from_numpy(enc_conv1_w.copy())
model.enc_conv1.bias.data = torch.from_numpy(enc_conv1_b.copy())
model.enc_conv2.weight.data = torch.from_numpy(enc_conv2_w.copy())
model.enc_conv2.bias.data = torch.from_numpy(enc_conv2_b.copy())
model.dec_conv1.weight.data = torch.from_numpy(dec_conv1_w.copy())
model.dec_conv1.bias.data = torch.from_numpy(dec_conv1_b.copy())
model.dec_conv2.weight.data = torch.from_numpy(dec_conv2_w.copy())
model.dec_conv2.bias.data = torch.from_numpy(dec_conv2_b.copy())
model.dec_conv3.weight.data = torch.from_numpy(dec_conv3_w.copy())
model.dec_conv3.bias.data = torch.from_numpy(dec_conv3_b.copy())

print("Weights loaded successfully!")
return model

```

```

# =====
# 3. Load CIFAR-10 test images
# =====

def load_cifar10_test(data_path):
    """Load CIFAR-10 test batch"""
    filepath = f"{data_path}/cifar-10-batches-bin/test_batch.bin"

    with open(filepath, 'rb') as f:
        data = np.frombuffer(f.read(), dtype=np.uint8)

    # Each image: 1 byte label + 3072 bytes image (32*32*3)
    data = data.reshape(10000, 3073)
    labels = data[:, 0]
    images = data[:, 1:3].reshape(10000, 3, 32, 32).astype(np.float32) / 255.0

    return images, labels

```

```

# =====
# 4. Main: Load, reconstruct, and visualize
# =====

# Paths (adjust for your Colab setup)
MODEL_PATH = "models/autoencoder_gpu.bin"
DATA_PATH = "data"

# Create model and load weights
model = Autoencoder()
model = load_cpp_weights(model, MODEL_PATH)
model.eval()

# Move to GPU if available
device = torch.device('cuda' if torch.cuda.is_available() else 'cpu')
model = model.to(device)

# Load test images
images, labels = load_cifar10_test(DATA_PATH)

```

```
# CIFAR-10 class names
class_names = ['airplane', 'automobile', 'bird', 'cat', 'deer',
                'dog', 'frog', 'horse', 'ship', 'truck']

# Select one image per class
sample_indices = []
for c in range(10):
    idx = np.where(labels == c)[0][0]
    sample_indices.append(idx)

# Reconstruct images
sample_images = images[sample_indices]
input_tensor = torch.from_numpy(sample_images).to(device)

with torch.no_grad():
    reconstructed = model(input_tensor)
    reconstructed = reconstructed.cpu().numpy()

# Clip to valid range
reconstructed = np.clip(reconstructed, 0, 1)

# =====
# 5. Visualize Original vs Reconstructed
# =====

fig, axes = plt.subplots(3, 10, figsize=(20, 6))

for i in range(10):
    # Original
    orig = sample_images[i].transpose(1, 2, 0) # CHW -> HWC
    axes[0, i].imshow(orig)
    axes[0, i].axis('off')
    axes[0, i].set_title(class_names[i], fontsize=10)

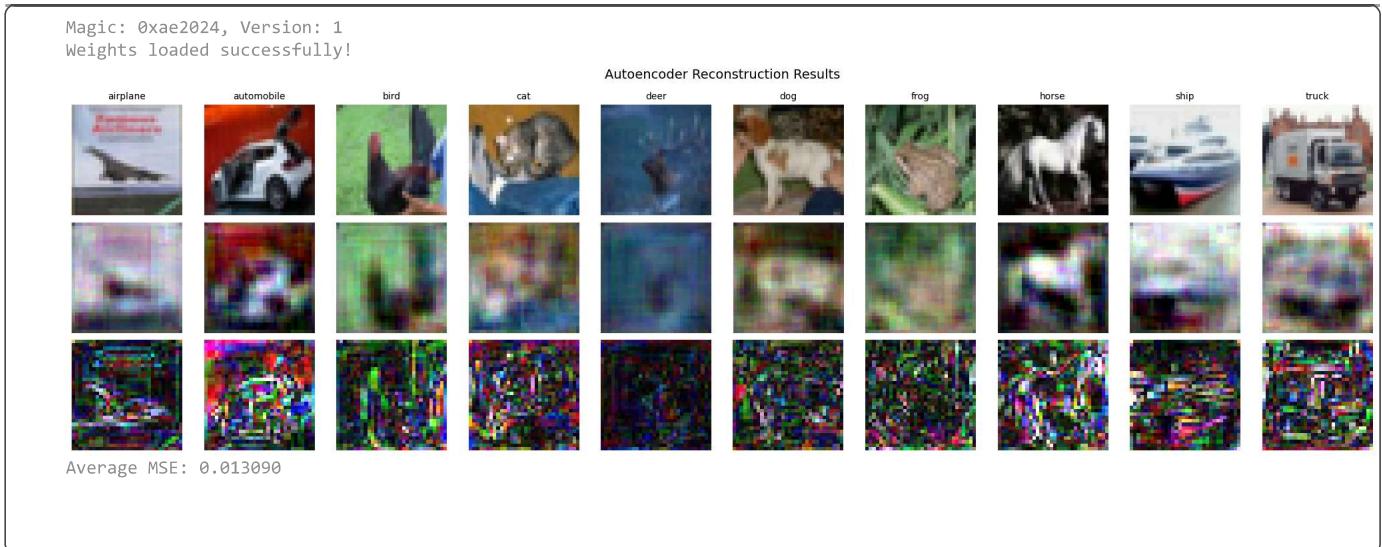
    # Reconstructed
    recon = reconstructed[i].transpose(1, 2, 0) # CHW -> HWC
    axes[1, i].imshow(recon)
    axes[1, i].axis('off')

    # Difference (amplified)
    diff = np.abs(orig - recon) * 3 # Amplify for visibility
    axes[2, i].imshow(np.clip(diff, 0, 1))
    axes[2, i].axis('off')

    axes[0, 0].set_ylabel('Original', fontsize=12)
    axes[1, 0].set_ylabel('Reconstructed', fontsize=12)
    axes[2, 0].set_ylabel('Difference (3x)', fontsize=12)

plt.suptitle('Autoencoder Reconstruction Results', fontsize=14)
plt.tight_layout()
plt.savefig('reconstruction_comparison.png', dpi=150, bbox_inches='tight')
plt.show()

# Calculate MSE
mse = np.mean((sample_images - reconstructed) ** 2)
print(f"Average MSE: {mse:.6f}")
```



▼ Feature Extraction

```
!./autoencoder_gpu --extract-features --model models/autoencoder_cpu.bin --data /content/PP-Autoencoder-unsupervised-feature-le

=====
CUDA Autoencoder for CIFAR-10
CSC14120 - Parallel Programming
=====

Loading CIFAR-10 dataset from /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_1.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_2.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_3.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_4.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/data_batch_5.bin...
  Loading /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin/test_batch.bin...
Dataset loaded successfully!
Training images: 50000
Test images: 10000

=====
GPU Feature Extraction
=====
Mode: Basic
Train samples: 50000
Test samples: 10000
Feature dimension: 8192
=====

=====
GPU Memory Breakdown (batch_size=100)
=====
Weights:      2.87 MB
Activations:  284.77 MB
Gradients:   287.63 MB
Pool Masks:  28.12 MB
-----
TOTAL:       603.39 MB
=====

[GPU Memory] After initialization: 1047.94 MB / 40506.81 MB (2.6%)
GPU Autoencoder initialized with batch size 100
Weights loaded from models/autoencoder_cpu.bin
Extracting training features...
[======>] 89%
Training features extracted in 6.05s

Extracting test features...
[======>] 89%
Test features extracted in 1.20s

=====
Feature Extraction Summary
=====
Train features time: 6.05s
Test features time: 1.20s
```

```
Total time: 7.39s
Train features shape: (50000, 8192)
Test features shape: (10000, 8192)
=====
```

```
!./autoencoder_gpu --extract-features --model models/autoencoder_gpu_opt_v2.bin --data /content/PP-Autoencoder-unsupervised-features
/bin/bash: line 1: ./autoencoder_gpu: No such file or directory
```

SVM Training

```
!pip uninstall -y cupy-cuda11x cupy-cuda12x cuml-cu11 cuml-cu12 2>/dev/null
!pip install --no-cache-dir cupy-cuda12x
!pip install --no-cache-dir cuml-cu12 --extra-index-url=https://pypi.nvidia.com

Collecting cupy-cuda12x
  Downloading cupy_cuda12x-13.6.0-cp312-cp312-manylinux2014_x86_64.whl.metadata (2.4 kB)
Requirement already satisfied: numpy<2.6,>=1.22 in /usr/local/lib/python3.12/dist-packages (from cupy-cuda12x) (2.0.2)
Collecting fastrlock>=0.5 (from cupy-cuda12x)
  Downloading fastrlock-0.8.3-cp312-cp312-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_28_x86_64.whl.metadata (7.7 kB)
  Downloading cupy_cuda12x-13.6.0-cp312-cp312-manylinux2014_x86_64.whl (112.9 MB)
  112.9/112.9 MB 250.8 MB/s eta 0:00:00
  Downloading fastrlock-0.8.3-cp312-cp312-manylinux_2_5_x86_64.manylinux1_x86_64.manylinux_2_28_x86_64.whl (53 kB)
  53.9/53.9 kB 225.8 MB/s eta 0:00:00
Installing collected packages: fastrlock, cupy-cuda12x
Successfully installed cupy-cuda12x-13.6.0 fastrlock-0.8.3
Looking in indexes: https://pypi.org/simple, https://pypi.nvidia.com
Collecting cuml-cu12
  Downloading https://pypi.nvidia.com/cuml-cu12/cuml\_cu12-25.12.0-cp312-cp312-manylinux\_2\_24\_x86\_64.manylinux\_2\_28\_x86\_64.whl 6.6/6.6 MB 86.9 MB/s eta 0:00:00
Collecting cuda-python<13.0a0,>=12.9.2 (from cuml-cu12)
  Downloading cuda_python-12.9.5-py3-none-any.whl.metadata (4.7 kB)
Collecting cuda-toolkit==12.* (from cuda-toolkit[cublas,cufft,curand,cusolver,cusparse]==12.*->cuml-cu12)
  Downloading https://pypi.nvidia.com/cuda\_toolkit/cuda\_toolkit-12.9.1-py2.py3-none-any.whl (2.3 kB)
Collecting cudf-cu12==25.12.* (from cuml-cu12)
  Downloading https://pypi.nvidia.com/cudf\_cu12/cudf\_cu12-25.12.0-cp312-cp312-manylinux\_2\_24\_x86\_64.manylinux\_2\_28\_x86\_64.whl 2.6/2.6 MB 341.9 MB/s eta 0:00:00
Requirement already satisfied: cupy-cuda12x>=13.6.0 in /usr/local/lib/python3.12/dist-packages (from cuml-cu12) (13.6.0)
Requirement already satisfied: joblib>=0.11 in /usr/local/lib/python3.12/dist-packages (from cuml-cu12) (1.5.3)
Collecting libcuml-cu12==25.12.* (from cuml-cu12)
  Downloading https://pypi.nvidia.com/libcuml\_cu12/libcuml\_cu12-25.12.0-py3-none-manylinux\_2\_27\_x86\_64.manylinux\_2\_28\_x86\_64.whl 477.2/477.2 MB 205.5 MB/s eta 0:00:00
Collecting numba-cuda<0.20.0a0,>=0.19.1 (from numba-cuda[cu12]<0.20.0a0,>=0.19.1->cuml-cu12)
  Downloading numba_cuda-0.19.1-py3-none-any.whl.metadata (3.6 kB)
Requirement already satisfied: numba<0.62.0a0,>=0.60.0 in /usr/local/lib/python3.12/dist-packages (from cuml-cu12) (0.60.0)
Requirement already satisfied: numpy<3.0a0,>=1.23 in /usr/local/lib/python3.12/dist-packages (from cuml-cu12) (2.0.2)
Requirement already satisfied: packaging in /usr/local/lib/python3.12/dist-packages (from cuml-cu12) (25.0)
Collecting pylibraft-cu12==25.12.* (from cuml-cu12)
  Downloading https://pypi.nvidia.com/pylibraft\_cu12/pylibraft\_cu12-25.12.0-cp312-cp312-manylinux\_2\_24\_x86\_64.manylinux\_2\_28\_x86\_64.whl 880.4/880.4 kB 241.7 MB/s eta 0:00:00
Requirement already satisfied: rich in /usr/local/lib/python3.12/dist-packages (from cuml-cu12) (13.9.4)
Collecting rmm-cu12==25.12.* (from cuml-cu12)
  Downloading https://pypi.nvidia.com/rmm\_cu12/rmm\_cu12-25.12.0-cp312-cp312-manylinux\_2\_24\_x86\_64.manylinux\_2\_28\_x86\_64.whl 1.2/1.2 MB 255.3 MB/s eta 0:00:00
Requirement already satisfied: scikit-learn>=1.4 in /usr/local/lib/python3.12/dist-packages (from cuml-cu12) (1.6.1)
Requirement already satisfied: scipy>=1.11.0 in /usr/local/lib/python3.12/dist-packages (from cuml-cu12) (1.16.3)
Collecting treelite<5.0.0a0,>=4.6.1 (from cuml-cu12)
  Downloading treelite-4.6.1-py3-none-manylinux2014_x86_64.whl.metadata (1.5 kB)
Collecting nvidia-cublas-cu12==12.9.1.4.* (from cuda-toolkit[cublas,cufft,curand,cusolver,cusparse]==12.*->cuml-cu12)
  Downloading https://pypi.nvidia.com/nvidia\_cublas\_cu12/nvidia\_cublas\_cu12-12.9.1.4-py3-none-manylinux\_2\_27\_x86\_64.whl (581.2 kB)
  581.2/581.2 MB 187.9 MB/s eta 0:00:00
Collecting nvidia-cufft-cu12==11.4.1.4.* (from cuda-toolkit[cublas,cufft,curand,cusolver,cusparse]==12.*->cuml-cu12)
  Downloading https://pypi.nvidia.com/nvidia\_cufft\_cu12/nvidia\_cufft\_cu12-11.4.1.4-py3-none-manylinux2014\_x86\_64.manylinux\_2\_10.x86\_64.whl 200.9/200.9 MB 190.7 MB/s eta 0:00:00
Collecting nvidia-curand-cu12==10.3.10.19.* (from cuda-toolkit[cublas,cufft,curand,cusolver,cusparse]==12.*->cuml-cu12)
  Downloading https://pypi.nvidia.com/nvidia\_curand\_cu12/nvidia\_curand\_cu12-10.3.10.19-py3-none-manylinux\_2\_27\_x86\_64.whl (68.3 kB)
  68.3/68.3 MB 247.9 MB/s eta 0:00:00
Collecting nvidia-cusolver-cu12==11.7.5.82.* (from cuda-toolkit[cublas,cufft,curand,cusolver,cusparse]==12.*->cuml-cu12)
  Downloading https://pypi.nvidia.com/nvidia\_cusolver\_cu12/nvidia\_cusolver\_cu12-11.7.5.82-py3-none-manylinux\_2\_27\_x86\_64.whl (338.1/338.1 MB 211.1 MB/s eta 0:00:00
Collecting nvidia-cusparse-cu12==12.5.10.65.* (from cuda-toolkit[cublas,cufft,curand,cusolver,cusparse]==12.*->cuml-cu12)
  Downloading https://pypi.nvidia.com/nvidia\_cusparse\_cu12/nvidia\_cusparse\_cu12-12.5.10.65-py3-none-manylinux2014\_x86\_64.manylinux\_2\_10.x86\_64.whl 265/265 MB 220.0 MB/s eta 0:00:00
```

```
!python train_svm_cuml.py --train models/train_features_cpu.bin --test models/test_features_cpu.bin --data /content/PP-Autoencoder-unsupervised-features
=====
cuML GPU SVM Training
CSC14120 - Parallel Programming
=====
```

```

Train features: models/train_features_cpu.bin
Test features: models/test_features_cpu.bin
CIFAR-10 data: /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin
Model name: gpu
=====

Loading CIFAR-10 labels...
Train labels: 50000
Test labels: 10000

Loading features from binary files...
Train features shape: (50000, 8192)
Test features shape: (10000, 8192)
Using cuML GPU-accelerated SVM

=====
Training SVM Classifier (GPU)
=====
C: 10.0
Kernel: rbf
Gamma: scale
Samples: 50000
Features: 8192
=====

Training SVM...
[2025-12-16 17:51:50.047] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:50.287] [CUML] [debug] SMO solver finished after 29 outer iterations, total inner 21817 iterations, and diff
[2025-12-16 17:51:50.988] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:51.159] [CUML] [debug] SMO solver finished after 28 outer iterations, total inner 22787 iterations, and diff
[2025-12-16 17:51:51.850] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:51.978] [CUML] [debug] SMO solver finished after 22 outer iterations, total inner 16607 iterations, and diff
[2025-12-16 17:51:52.672] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:52.808] [CUML] [debug] SMO solver finished after 23 outer iterations, total inner 18213 iterations, and diff
[2025-12-16 17:51:53.542] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:53.649] [CUML] [debug] SMO solver finished after 18 outer iterations, total inner 14149 iterations, and diff
[2025-12-16 17:51:54.725] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:54.809] [CUML] [debug] SMO solver finished after 14 outer iterations, total inner 10550 iterations, and diff
[2025-12-16 17:51:55.504] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:55.626] [CUML] [debug] SMO solver finished after 21 outer iterations, total inner 15978 iterations, and diff
[2025-12-16 17:51:56.318] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:56.503] [CUML] [debug] SMO solver finished after 30 outer iterations, total inner 24750 iterations, and diff
[2025-12-16 17:51:57.198] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:57.350] [CUML] [debug] SMO solver finished after 26 outer iterations, total inner 20409 iterations, and diff
[2025-12-16 17:51:58.043] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:58.161] [CUML] [debug] SMO solver finished after 21 outer iterations, total inner 15907 iterations, and diff
[2025-12-16 17:51:58.857] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:58.994] [CUML] [debug] SMO solver finished after 25 outer iterations, total inner 19153 iterations, and diff
[2025-12-16 17:51:59.842] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:51:59.954] [CUML] [debug] SMO solver finished after 20 outer iterations, total inner 15013 iterations, and diff
[2025-12-16 17:52:00.925] [CUML] [debug] Creating working set with 1024 elements
[2025-12-16 17:52:01.035] [CUML] [debug] SMO solver finished after 19 outer iterations, total inner 15111 iterations, and diff
[2025-12-16 17:52:01.877] [CUML] [debug] Creating working set with 1024 elements

```

```

!python train_svm_cuml.py --train models/train_features_gpu_opt_v2.bin --test models/test_features_gpu_opt_v2.bin --data /conte
=====
cuML GPU SVM Training
CSC14120 - Parallel Programming
=====
Train features: models/train_features_gpu_opt_v2.bin
Test features: models/test_features_gpu_opt_v2.bin
CIFAR-10 data: /content/PP-Autoencoder-unsupervised-feature-learning/data/cifar-10-batches-bin
Model name: gpu
=====

Loading CIFAR-10 labels...
Train labels: 50000
Test labels: 10000

Loading features from binary files...
Traceback (most recent call last):
  File "/content/PP-Autoencoder-unsupervised-feature-learning/train_svm_cuml.py", line 290, in <module>
    main()
  File "/content/PP-Autoencoder-unsupervised-feature-learning/train_svm_cuml.py", line 267, in main
    train_features = load_features(args.train, CIFAR_TRAIN_SIZE, FEATURE_DIM)
    ^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^^
  File "/content/PP-Autoencoder-unsupervised-feature-learning/train_svm_cuml.py", line 69, in load_features
    raise FileNotFoundError(f"Feature file not found: {filepath}")
FileNotFoundError: Feature file not found: models/train_features_gpu_opt_v2.bin

```

Demo

```
!pip uninstall torch torchvision torchaudio -y  
!pip install torch torchvision torchaudio --index-url https://download.pytorch.org/whl/cu121 --quiet
```

```
Found existing installation: torch 2.9.0+cu126  
Uninstalling torch-2.9.0+cu126:  
  Successfully uninstalled torch-2.9.0+cu126  
Found existing installation: torchvision 0.24.0+cu126  
Uninstalling torchvision-0.24.0+cu126:  
  Successfully uninstalled torchvision-0.24.0+cu126  
Found existing installation: torchaudio 2.9.0+cu126  
Uninstalling torchaudio-2.9.0+cu126:  
  Successfully uninstalled torchaudio-2.9.0+cu126  
          ━━━━━━━━━━━━━━━━ 780.4/780.4 MB 1.6 MB/s eta 0:00:00  
          ━━━━━━━━━━━━━━ 7.3/7.3 MB 148.0 MB/s eta 0:00:00  
          ━━━━━━━━━━━━ 3.4/3.4 MB 110.5 MB/s eta 0:00:00  
          ━━━━━━━━━━ 23.7/23.7 MB 120.5 MB/s eta 0:00:00  
          ━━━━━━ 823.6/823.6 KB 64.7 MB/s eta 0:00:00  
          ━ 14.1/14.1 MB 147.7 MB/s eta 0:00:00
```