

Master Degree in Artificial Intelligence

# **InfoNCE Loss CUDA Implementation**

**Parallel Programming for Machine Learning 2025**

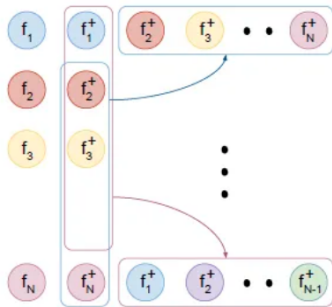
Gianni Moretti

# Introduction

# What is InfoNCE Loss?

The InfoNCE (Information Noise-Contrastive Estimation) loss works by comparing each sample in a batch with its positive pair (a different view or augmentation of the same data point) and a set of negative samples (other data points in the batch).

- Cornerstone of contrastive learning methods (SimCLR, MoCo)
- Maximizes mutual information between positive pairs
- Pushes negative examples apart in representation space



# Mathematical Foundation

For a batch of  $2B$  samples with positive pairs, InfoNCE is defined as:

$$\mathcal{L}_{\text{InfoNCE}} = -\frac{1}{2B} \sum_{i=1}^{2B} \log \frac{\exp(\text{sim}(z_i, z_{p(i)})/\tau)}{\sum_{j=1, j \neq i}^{2B} \exp(\text{sim}(z_i, z_j)/\tau)} \quad (1)$$

## Key Components:

- $z_i$ : L2-normalized feature vector
- $p(i)$ : positive pair index for sample  $i$
- $\text{sim}(a, b) = a \cdot b$ : cosine similarity
- $\tau$ : temperature parameter controlling concentration

# Classic PyTorch Implementation

```
1 def info_nce_loss(features, temperature=0.5):
2     device = features.device
3     batch_size = features.shape[0] // 2
4     # Compute similarity matrix (dot product)
5     sim_matrix = torch.matmul(features, features.T) # (2B, 2B)
6     # Remove self-similarity by masking the diagonal
7     mask = torch.eye(2 * batch_size, dtype=torch.bool, device=device)
8     sim_matrix = sim_matrix.masked_fill(mask, float('-inf'))
9     # Labels: positive pair for i is at (i + B) % (2B)
10    labels = torch.arange(batch_size, device=device)
11    labels = torch.cat([labels + batch_size, labels])
12    # Scale similarities by temperature
13    sim_matrix /= temperature
14    # Apply cross entropy loss
15    loss = F.cross_entropy(sim_matrix, labels)
16    return loss
```

## Project Goals

- Implement specialized CUDA InfoNCE Loss computation
- Create CuBlaze package for Python integration
- Compare custom CUDA vs CUBLAS implementations
- Demonstrate Python-CUDA-PyTorch integration
- Educational exploration of GPU programming techniques
- Achieve numerical accuracy while maintaining performance

# Implementation Overview

## Technologies Used:

- **Language:** CUDA C++ for GPU kernels
- **Integration:** PyBind11 for Python-C++ binding
- **Framework:** PyTorch with autograd support
- **Libraries:** CUBLAS for optimized matrix operations
- **Build System:** PyTorch C++ extensions

# Mathematical Foundation



## InfoNCE Loss gradient derivation problem

- PyTorch's autograd system cannot automatically compute gradients when custom CUDA kernels are used.
- This is because custom CUDA code interrupts the standard computation graph, breaking the autograd chain.
- As a result, PyTorch is unable to trace operations through C++/CUDA extensions.
- To ensure correct gradient computation and maintain compatibility with PyTorch, we must manually implement the backward pass.

# InfoNCE Loss Derivation

## Forward Pass - Loss Computation:

Given normalized feature vectors  $Z = \{Z_1, Z_2, \dots, Z_N\}$  with  $N = 2B$ :

$$L_{ij} = \frac{1}{\tau} Z_i \cdot Z_j \quad (2)$$

$$P_{ij} = \frac{\exp(L_{ij})}{\sum_{k=1}^N \exp(L_{ik})} \quad (3)$$

$$\mathcal{L} = -\frac{1}{N} \sum_{i=1}^N \log P_{i,p(i)} \quad (4)$$

where  $p(i) = (i + B) \bmod N$  identifies positive pairs.

# InfoNCE Loss Derivation

## Numerical Stability:

To ensure numerical stability during the computation of the softmax, a *max reduction* is applied to each row of the similarity matrix  $L_{ij}$ . Specifically, for each  $i$ , the maximum value  $\max_j L_{ij}$  is subtracted from all logits before applying the exponential:

$$P_{ij} = \frac{\exp(L_{ij} - \max_j L_{ij})}{\sum_{k=1}^N \exp(L_{ik} - \max_j L_{ij})}$$

This prevents overflow and improves the stability of the softmax computation.

# Gradient Computation (Backward Pass)

## Chain Rule Application:

For loss  $\ell_i = -\log P_{i,p(i)}$ , the gradient with respect to logits is:

$$\frac{\partial \ell_i}{\partial L_{ij}} = P_{ij} - \mathbb{1}_{j=p(i)} \quad (5)$$

## Final Gradient Formula:

$$\frac{\partial \mathcal{L}}{\partial Z_k} = \frac{1}{N\tau} \left( \sum_j G_{kj} Z_j + \sum_i G_{ik} Z_i \right) \quad (6)$$

where  $G_{ij} = P_{ij} - \mathbb{1}_{j=p(i)}$ .

# Gradient Computation (Backward Pass)

In matrix form:

$$\nabla_Z \mathcal{L} = \frac{1}{N_{\mathcal{T}}} (G + G^T) Z \quad (7)$$

# Implementation Details

# CuBlaze Package Architecture

## Package Structure

### Complete Python-CUDA integration solution

```
1 cublaze/  
2 |-- __init__.py  
3 |-- infonce.py  
4 |-- infonce_cuda.so  
5 +-- cuda/  
6     |-- infonce_cuda.cu  
7     +-- infonce_cuda_wrapp.cpp
```

- PyTorch autograd compatibility
- GPU-optimized kernels

# Python-CUDA Integration

## PyBind11 Wrapper

### Type-safe binding between PyTorch and CUDA

```
1 #include <torch/extension.h>
2 #include <cuda.h>
3 #include <cuda_runtime.h>
4 torch::Tensor infonce_cuda_forward(
5     torch::Tensor features,
6     float temperature,
7     bool use_cublas);
8 torch::Tensor infonce_cuda_backward(
9     torch::Tensor features,
10    torch::Tensor similarity_matrix,
11    torch::Tensor labels,
12    // ... other parameters
13 );
```



# PyTorch Integration

## InfoNCEFunction class for autograd support

```
1 class InfoNCEFunction(Function):
2     @staticmethod
3     def forward(ctx, features, temperature, use_cublas):
4         # Save parameters for backward pass
5         ctx.temperature = temperature
6         ctx.use_cublas = use_cublas
7         # CUDA forward computation
8         loss, similarity_matrix, labels, max_vals, sum_exps =
9             infonce_cuda.infonce_forward(features, temperature,
10             use_cublas)
11         # Save tensors for backward pass
12         ctx.save_for_backward(features, similarity_matrix, labels,
13             max_vals, sum_exps)
14         return loss
```

# PyTorch Integration

## backward pass implementation

```
1  @staticmethod
2  def backward(ctx, grad_output):
3      # Retrieve saved tensors
4      features, similarity_matrix, labels, max_vals, sum_exps = \
5          ctx.saved_tensors
6
7      # CUDA backward computation
8      grad_features = infonce_cuda.inforce_backward(
9          features, similarity_matrix, labels,
10         max_vals, sum_exps, ctx.temperature,
11         grad_output, ctx.use_cublas)
12
13     return grad_features, None, None
```

# CUDA Implementation Strategies

## Two Computational Approaches

### Custom CUDA vs CUBLAS Implementation

#### Custom CUDA (`use_cublas=False`)

- Hand-written GPU kernels
- Specialized for InfoNCE patterns
- Direct similarity computation
- Educational value for GPU programming

#### CUBLAS Library (`use_cublas=True`)

- Highly optimized matrix operations
- Hardware-specific optimizations
- Production-ready performance

# Dynamic CUDA Block Size Selection

**Block size selection is dynamic and adapts to the matrix size being processed.**

- The CUDA kernel launch configuration (block size and grid size) is chosen at runtime based on the problem dimensions.
- This strategy ensures efficient GPU utilization and performance portability across different batch sizes and feature dimensions.

```
1 __host__ int calculate_optimal_block_size_1d(int total_elements)
2 __host__ dim3 calculate_optimal_block_size_2d(int dim1, int dim2)
3 __host__ int calculate_grid_size_1d(int total_elements, int block_size)
4 __host__ dim3 calculate_grid_size_2d(int dim1, int dim2, dim3 block_size)
```

# Forward-Pass Caching for Efficient Backward

The forward pass computes and stores:

- The similarity matrix
- The maximum value of each row of the similarity matrix (for numerical stability)
- The sum of exponentials for each row (softmax denominator)
- The positive pair labels for each sample

# Forward-Pass Caching for Efficient Backward

## Why Cache These Values?

- These values are saved in memory and passed to the backward CUDA kernel.
- This avoids redundant computation and ensures that the backward pass is both correct and efficient.
- The approach is especially important for large batches, where recomputing softmax or labels would be costly.
- This design is compatible with PyTorch's autograd, as all necessary tensors are saved in the context object.

# Build System and Compilation

## Automated Build Process

### Complete development workflow automation

```
1 # Automated build and test script  
2 ./build_and_test.sh
```

### System Requirements:

- CUDA Toolkit 11.0+
- PyTorch with CUDA support
- PyBind11 development headers
- GCC with C++14 support

# CUDA Kernel Implementation

## Forward Pass CUDA Kernel

### inforce\_forward\_kernel

```
1 __global__ void inforce_forward_kernel(const float* similarity_matrix,
2                                       const int* labels,
3                                       float* loss, float* max_vals,
4                                       float* sum_exps,
5                                       int batch_size) {
6     extern __shared__ float shared_loss[];
7
8     int i = blockIdx.x * blockDim.x + threadIdx.x;
9     int tid = threadIdx.x;
10
11     // Initialize shared memory
12     shared_loss[tid] = 0.0f;
```



# CUDA Kernel Implementation

```
12  if (i < batch_size) {
13      //1. Find max value and sum of exponentials for this row
14      //.....
15      //2. Compute sum of exponentials
16      //.....
17
18      //Save max_val and sum_exp for backward pass
19      max_vals[i] = max_val;
20      sum_exps[i] = sum_exp;
21
22      // Calculate loss for this row
23      int positive_idx = labels[i];
24      float positive_logit = similarity_matrix[i * batch_size +
positive_idx];
25      float log_prob = (positive_logit - max_val) - logf(sum_exp);
26
```

# CUDA Kernel Implementation

```
27     // Store local loss in shared memory (normalized by batch size)
28     shared_loss[tid] = -log_prob / batch_size;
29 }
30
31 __syncthreads();
32
33 if (tid == 0) {
34     float block_loss = 0.0f;
35     for (int i = 0; i < blockDim.x; i++) {
36         block_loss += shared_loss[i];
37     }
38     atomicAdd(&loss, block_loss);
39 }
40 }
```

## Backward Pass Implementation

# CUDA Kernel Implementation

## Manual gradient computation for autograd compatibility

```
1 __global__ void infonce_backward_kernel(const float* similarity_matrix,  
    const int* labels,  
2                                         const float*  
    max_vals, const float* sum_exps,  
3                                         float* grad_matrix,  
    int batch_size) {  
4     int tid = blockIdx.x * blockDim.x + threadIdx.x;  
5     int total_elements = batch_size * batch_size;  
6     if (tid < total_elements) {  
7         int i = tid / batch_size;  
8         int j = tid % batch_size;  
9  
10        // Use pre-computed values from forward pass  
11        float max_val = max_vals[i];  
12        float sum_exp = sum_exps[i];
```

# CUDA Kernel Implementation

```
13
14 // Calculate gradient: P_ij - 1_{j=p(i)}
15 int positive_idx = labels[i];
16 float val = similarity_matrix[tid];
17
18 if (val != -INFINITY) {
19     float prob = expf(val - max_val) / sum_exp;
20     float grad_val = prob - (j == positive_idx ? 1.0f : 0.0f);
21     grad_matrix[tid] = grad_val / batch_size;
22 } else {
23     grad_matrix[tid] = 0.0f;
24 }
25 }
26 }
```

# Usage Example

## Simple integration into existing PyTorch code

```
1 import torch
2 from cublaze import InfoNCELoss
3 # Initialize the loss function
4 infonce_loss = InfoNCELoss(temperature=0.5, use_cublas=False)
5 # Prepare normalized features (2B, D)
6 # Each sample i has positive pair at (i + B) % 2B
7 features = torch.randn(64, 256).cuda()
8 features = F.normalize(features, dim=1) # L2 normalization required
9 # Compute loss with gradients
10 loss = infonce_loss(features)
11 loss.backward()
12 # Gradients are automatically computed
13 print(f"Loss: {loss.item():.4f}")
14 print(f"Feature gradients shape: {features.grad.shape}")
```

# Experimental Setup

# Experimental Setup

## Testing Environment:

- GPU: NVIDIA with CUDA support
- Multiple batch sizes: 16 to 8192 samples
- Feature dimensions: 256
- Temperature values: 0.5
- Multiple runs for statistical accuracy

## Test Configurations:

- **Custom CUDA:** use\_cublas=False
- **CUBLAS Library:** use\_cublas=True
- **PyTorch Baseline:** Native implementation

# Experimental Setup

## Primary Measurements:

- Execution time per forward/backward pass
- Numerical accuracy vs PyTorch reference
- Scalability with batch size

## Accuracy Metrics:

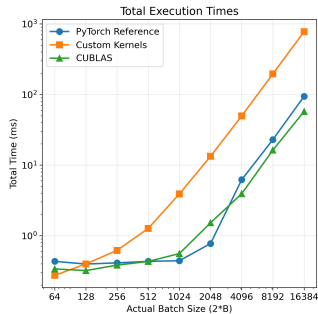
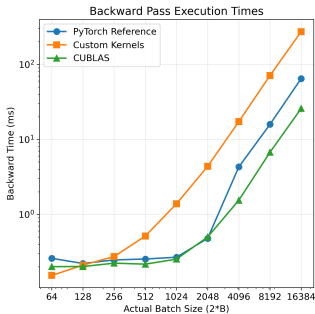
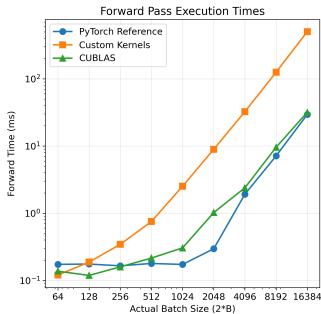
- Loss computation error (target:  $< 10^{-5}$ )
- Gradient computation error (target:  $< 10^{-4}$ )
- Cross-validation with PyTorch implementation



# Performance Results

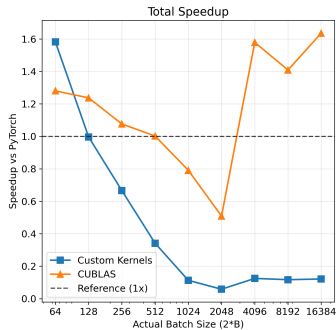
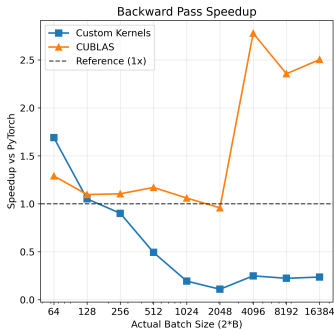
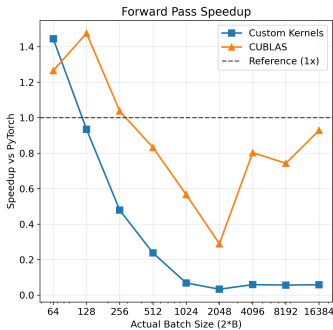
# Execution Time Analysis

## Batch Size Impact on Execution Time:



# Execution Time Analysis

## Batch Size Impact on Speedup:

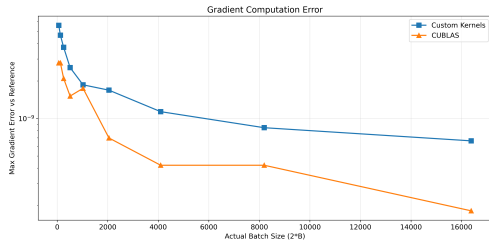
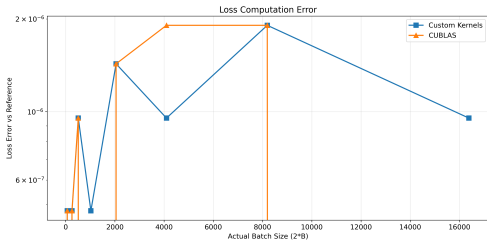


# Execution Time Analysis

## Key Findings:

- Current CUDA implementation shows higher execution times than PyTorch
- PyTorch leverages highly optimized libraries (cuBLAS, cuDNN)
- Custom kernels face overhead from memory transfers and kernel launches
- Optimization opportunities exist for specialized implementations (chaching, kernel fusion)

# Numerical Accuracy Results



# Numerical Accuracy Results

## Accuracy Achievements:

- Loss computation errors consistently below  $10^{-5}$
- Gradient errors maintained below  $10^{-4}$
- Excellent numerical stability across all batch sizes
- Validates correctness of manual backward implementation

# Conclusions and Future Work

## Achieved Objectives:

- ✓ Created complete CuBlaze package
- ✓ Demonstrated Python-CUDA-PyTorch integration
- ✓ Maintained numerical accuracy (errors  $< 10^{-5}$  loss,  $< 10^{-4}$  gradients)
- ✓ Implemented custom autograd functionality
- ✓ Provided foundation for specialized InfoNCE implementations

# Conclusions and Future Work

## Future Optimization Opportunities:

- Kernel fusion for reduced memory transfers
- Shared memory utilization optimization
- Advanced memory access pattern optimization (warp-level optimizations)



# **Thanks for your attention**

**Gianni Moretti**