# Contents

# DGEMM_ALPHA Production Usage Guide

## Overview

`DGEMM_ALPHA` is the world's first implementation of Google's AlphaTensor algorithm for optimized $4{\times}4$ matrix multiplication. This guide shows you how to integrate and use it in production environments across different programming languages and deployment scenarios.

## Quick Start Summary

**What it does**: Optimized matrix multiplication using AlphaTensor's 49-operation algorithm
**Best for**: 4×4 matrix operations, machine learning workloads, transformer attention mechanisms
**Performance**: Up to 4.7x speedup for specific matrix patterns, 1.15x average improvement
**Compatibility**: Drop-in replacement for standard DGEMM with identical API

---

## Installation and Build

### Option 1: Build from Source (Recommended)

```
# 1. Clone the repository
git clone https://github.com/GauntletAI/lapack_ai.git
cd lapack_ai

# 2. Build LAPACK with AlphaTensor support
mkdir build && cd build
cmake ..
make -j$(nproc)

# 3. Install libraries (optional)
sudo make install
```

### Option 2: Using Package Managers (Future)

```
# Once packaged (not yet available):
# pip install lapack-alphatensor  # Python
# apt-get install lapack-alphatensor  # Ubuntu/Debian
# brew install lapack-alphatensor  # macOS
```

---

## Programming Language Interfaces

### 1. C/C++ Usage (CBLAS Interface)

### Headers and Linking

```
#include <cblas.h>  // Main CBLAS interface

// Linking flags:
// gcc -lblas -llapack -lcblas your_program.c
```

### Basic Usage

```
#include <cblas.h>
#include <stdio.h>

int main() {
    // 4x4 matrices (AlphaTensor optimal size)
```

```c
    double A[16] = {1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16};
    double B[16] = {1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1, 0, 0, 0, 0, 1}; // Identity
    double C[16] = {0}; // Result matrix

    // AlphaTensor-optimized multiplication: C = alpha * A * B + beta * C
    cblas_dgemm_alpha(
        CblasRowMajor,      // Layout
        CblasNoTrans,       // TransA
        CblasNoTrans,       // TransB
        4, 4, 4,            // M, N, K (4x4 matrices)
        1.0,                // alpha
        A, 4,               // Matrix A, leading dimension
        B, 4,               // Matrix B, leading dimension
        0.0,                // beta
        C, 4                // Matrix C (result), leading dimension
    );

    // C now contains the optimized result
    printf("Result computed with AlphaTensor algorithm\n");
    return 0;
}
```

**Production Example: Batch Processing**

```c
// Optimized batch processing for ML workloads
void process_attention_matrices(double* Q, double* K, double* V,
                                double* output, int batch_size) {
    for (int i = 0; i < batch_size; i++) {
        // Each attention head uses 4x4 matrices
        double* q_head = Q + i * 16;
        double* k_head = K + i * 16;
        double* result = output + i * 16;

        // Use AlphaTensor for optimal performance
        cblas_dgemm_alpha(CblasRowMajor, CblasNoTrans, CblasTrans,
                    4, 4, 4, 1.0, q_head, 4, k_head, 4, 0.0, result, 4);
    }
}
```

## 2. Python Usage (NumPy/SciPy Integration)

**Option A: Direct CTYPES Interface**

```python
import ctypes
import numpy as np
from ctypes import POINTER, c_int, c_double, c_char

# Load the library
try:
```

```python
    liblapack = ctypes.CDLL('liblapack.so')  # Linux
except:
    liblapack = ctypes.CDLL('liblapack.dylib')  # macOS

# Define the AlphaTensor function
dgemm_alpha = liblapack.dgemm_alpha_
dgemm_alpha.argtypes = [
    POINTER(c_char),  # TRANSA
    POINTER(c_char),  # TRANSB
    POINTER(c_int),   # M
    POINTER(c_int),   # N
    POINTER(c_int),   # K
    POINTER(c_double), # ALPHA
    POINTER(c_double), # A
    POINTER(c_int),   # LDA
    POINTER(c_double), # B
    POINTER(c_int),   # LDB
    POINTER(c_double), # BETA
    POINTER(c_double), # C
    POINTER(c_int)    # LDC
]

def alphatensor_multiply(A, B, alpha=1.0, beta=0.0):
    """
    AlphaTensor-optimized matrix multiplication for 4x4 matrices.

    Args:
        A, B: 4x4 numpy arrays (float64)
        alpha, beta: scaling factors

    Returns:
        C: Result matrix (alpha * A @ B + beta * C)
    """
    assert A.shape == (4, 4) and B.shape == (4, 4), "Matrices must be 4x4"
    assert A.dtype == np.float64 and B.dtype == np.float64, "Must be float64"

    # Ensure Fortran-contiguous arrays
    A = np.asfortranarray(A)
    B = np.asfortranarray(B)
    C = np.zeros((4, 4), dtype=np.float64, order='F')

    # Call AlphaTensor FORTRAN routine
    m, n, k = c_int(4), c_int(4), c_int(4)
    lda, ldb, ldc = c_int(4), c_int(4), c_int(4)
    alpha_c, beta_c = c_double(alpha), c_double(beta)
    trans_n = c_char(b'N')

    dgemm_alpha(
```

4

```python
            ctypes.byref(trans_n), ctypes.byref(trans_n),
            ctypes.byref(m), ctypes.byref(n), ctypes.byref(k),
            ctypes.byref(alpha_c),
            A.ctypes.data_as(POINTER(c_double)), ctypes.byref(lda),
            B.ctypes.data_as(POINTER(c_double)), ctypes.byref(ldb),
            ctypes.byref(beta_c),
            C.ctypes.data_as(POINTER(c_double)), ctypes.byref(ldc)
    )

    return C


# Usage example
A = np.random.random((4, 4))
B = np.eye(4)   # Identity matrix (optimal case for AlphaTensor)
C = alphatensor_multiply(A, B)
print(f"AlphaTensor result: {C}")
```

**Option B: SciPy Integration (Future)**

```python
import numpy as np
from scipy.linalg import blas


# Once integrated into SciPy (future development):
def alphatensor_gemm(A, B, alpha=1.0, beta=0.0):
    """Use AlphaTensor through SciPy BLAS interface"""
    return blas.dgemm_alpha(alpha, A, B, beta=beta)
```

**3. FORTRAN Usage (Direct Interface)**

```fortran
PROGRAM ALPHATENSOR_EXAMPLE
    IMPLICIT NONE

    ! Matrix declarations
    INTEGER, PARAMETER :: N = 4
    DOUBLE PRECISION :: A(N,N), B(N,N), C(N,N)
    DOUBLE PRECISION :: ALPHA, BETA
    INTEGER :: I, J

    ! Initialize matrices
    ALPHA = 1.0D0
    BETA = 0.0D0

    ! Fill test matrices
    DO I = 1, N
        DO J = 1, N
            A(I,J) = DBLE(I*N + J)   ! Sequential values
            B(I,J) = 0.0D0
            IF (I .EQ. J) B(I,J) = 1.0D0   ! Identity matrix
```

```fortran
            C(I,J) = 0.0DO
        END DO
    END DO

    ! Call AlphaTensor-optimized multiplication
    CALL DGEMM_ALPHA('N', 'N', N, N, N, ALPHA, A, N, B, N, BETA, C, N)

    ! Print result
    WRITE(*,*) 'AlphaTensor Result:'
    DO I = 1, N
        WRITE(*,'(4F8.2)') (C(I,J), J=1,N)
    END DO

END PROGRAM
```

---

## Performance Optimization Guide

### When to Use AlphaTensor

#### Optimal Use Cases:

- **4×4 matrix operations** (primary optimization target)
- **Identity matrices**: Up to 3.9x speedup
- **Zero matrices**: Up to 2.5x speedup

- **Mixed-sign matrices**: Up to 4.7x speedup
- **Transformer attention mechanisms** (4×4 head attention)
- **Small block operations** in larger matrix decompositions

#### Less Optimal Cases:

- **Large sparse matrices**: May be slower than optimized sparse routines
- **Very large matrices**: Overhead may outweigh benefits
- **Single matrix operations**: Setup cost may dominate

### Performance Monitoring

```c
#include <time.h>

double benchmark_alphatensor(int iterations) {
    double A[16], B[16], C[16];
    clock_t start, end;

    // Initialize matrices
    for (int i = 0; i < 16; i++) {
        A[i] = (double)rand() / RAND_MAX;
        B[i] = (double)rand() / RAND_MAX;
        C[i] = 0.0;
```

```
    }

    start = clock();
    for (int iter = 0; iter < iterations; iter++) {
        cblas_dgemm_alpha(CblasRowMajor, CblasNoTrans, CblasNoTrans,
                          4, 4, 4, 1.0, A, 4, B, 4, 0.0, C, 4);
    }
    end = clock();

    return ((double)(end - start)) / CLOCKS_PER_SEC;
}
```

---

## Integration Strategies

### 1. Drop-in Replacement Strategy

```
// Replace existing DGEMM calls with DGEMM_ALPHA for 4x4 operations
#ifdef USE_ALPHATENSOR
    #define MATRIX_MULTIPLY cblas_dgemm_alpha
#else
    #define MATRIX_MULTIPLY cblas_dgemm
#endif

// Your existing code automatically uses AlphaTensor when available
MATRIX_MULTIPLY(CblasRowMajor, CblasNoTrans, CblasNoTrans,
                4, 4, 4, 1.0, A, 4, B, 4, 0.0, C, 4);
```

### 2. Hybrid Strategy (Recommended)

```
void smart_matrix_multiply(double* A, double* B, double* C,
                           int M, int N, int K) {
    if (M == 4 && N == 4 && K == 4) {
        // Use AlphaTensor for 4x4 operations
        cblas_dgemm_alpha(CblasRowMajor, CblasNoTrans, CblasNoTrans,
                          M, N, K, 1.0, A, M, B, K, 0.0, C, M);
    } else {
        // Use standard DGEMM for other sizes
        cblas_dgemm(CblasRowMajor, CblasNoTrans, CblasNoTrans,
                    M, N, K, 1.0, A, M, B, K, 0.0, C, M);
    }
}
```

### 3. Machine Learning Framework Integration

```
# PyTorch custom op example (conceptual)
import torch

class AlphaTensorMatMul(torch.autograd.Function):
```

```python
    @staticmethod
    def forward(ctx, A, B):
        if A.shape == (4, 4) and B.shape == (4, 4):
            # Use AlphaTensor implementation
            return alphatensor_multiply(A.numpy(), B.numpy())
        else:
            return torch.matmul(A, B)

    @staticmethod
    def backward(ctx, grad_output):
        # Standard backpropagation
        return grad_output, grad_output

# Register as custom operation
torch.ops.alphatensor = AlphaTensorMatMul.apply
```

---

## Production Deployment

### 1. Docker Deployment

```dockerfile
FROM ubuntu:20.04

# Install dependencies
RUN apt-get update && apt-get install -y \
    build-essential \
    cmake \
    gfortran \
    libopenblas-dev

# Copy and build AlphaTensor LAPACK
COPY . /alphatensor-lapack
WORKDIR /alphatensor-lapack

RUN mkdir build && cd build && \
    cmake .. && \
    make -j$(nproc) && \
    make install

# Set library path
ENV LD_LIBRARY_PATH=/usr/local/lib:$LD_LIBRARY_PATH

# Your application
COPY app/ /app
WORKDIR /app
RUN gcc -o myapp main.c -llapack -lblas -lcblas
CMD ["./myapp"]
```

## 2. Cloud Platform Integration

### AWS Lambda Example

```python
import json
import numpy as np
from alphatensor_lib import alphatensor_multiply

def lambda_handler(event, context):
    """
    AWS Lambda function using AlphaTensor for 4x4 matrix operations
    """
    try:
        # Parse input matrices
        A = np.array(event['matrix_a']).reshape(4, 4)
        B = np.array(event['matrix_b']).reshape(4, 4)

        # Use AlphaTensor optimization
        result = alphatensor_multiply(A, B)

        return {
            'statusCode': 200,
            'body': json.dumps({
                'result': result.tolist(),
                'algorithm': 'AlphaTensor',
                'operations': 49  # vs 64 for standard
            })
        }
    except Exception as e:
        return {
            'statusCode': 500,
            'body': json.dumps({'error': str(e)})
        }
```

## 3. Kubernetes Deployment

```yaml
apiVersion: apps/v1
kind: Deployment
metadata:
  name: alphatensor-service
spec:
  replicas: 3
  selector:
    matchLabels:
      app: alphatensor-service
  template:
    metadata:
      labels:
        app: alphatensor-service
```

```yaml
spec:
  containers:
  - name: alphatensor-app
    image: your-registry/alphatensor-app:latest
    ports:
    - containerPort: 8080
    env:
    - name: LD_LIBRARY_PATH
      value: "/usr/local/lib"
    - name: USE_ALPHATENSOR
      value: "true"
    resources:
      requests:
        memory: "512Mi"
        cpu: "500m"
      limits:
        memory: "1Gi"
        cpu: "1000m"
```

---

## Error Handling and Debugging

### Common Issues and Solutions

### 1. Library Not Found

```bash
# Error: liblapack.so not found
export LD_LIBRARY_PATH=/path/to/lapack/lib:$LD_LIBRARY_PATH
# Or copy libraries to standard location:
sudo cp build/lib/* /usr/local/lib/
sudo ldconfig
```

### 2. Symbol Not Found

```bash
# Error: undefined symbol: dgemm_alpha_
# Check if AlphaTensor variant is compiled:
nm liblapack.so | grep dgemm_alpha
# Should show: T dgemm_alpha_
```

### 3. Performance Debugging

```c
// Add timing and validation
#include <assert.h>
#include <math.h>

void validate_alphatensor_result(double* A, double* B, double* C_alpha) {
    double C_standard[16] = {0};

    // Compare with standard DGEMM
```

```c
    cblas_dgemm(CblasRowMajor, CblasNoTrans, CblasNoTrans,
                4, 4, 4, 1.0, A, 4, B, 4, 0.0, C_standard, 4);

    // Check numerical accuracy (should be within 1e-14)
    for (int i = 0; i < 16; i++) {
        double diff = fabs(C_alpha[i] - C_standard[i]);
        assert(diff < 1e-12 && "AlphaTensor accuracy error");
    }
    printf("  AlphaTensor accuracy validated\n");
}
```

---

## Best Practices

### 1. Memory Management

- **Use aligned memory** for better SIMD performance
- **Prefer contiguous arrays** to optimize cache usage
- **Batch operations** when possible to amortize overhead

### 2. Compiler Optimization

```bash
# Recommended compilation flags:
gcc -O3 -march=native -mavx2 -ffast-math your_code.c -llapack -lblas

# For production (more conservative):
gcc -O2 -march=x86-64 your_code.c -llapack -lblas
```

### 3. Testing Strategy

```c
// Always validate in development
#ifdef DEBUG
    validate_alphatensor_result(A, B, C);
#endif

// Performance monitoring in production
LOG_PERFORMANCE("AlphaTensor 4x4 multiply", execution_time_ms);
```

### 4. Fallback Strategy

```c
// Graceful degradation
if (alphatensor_available() && matrix_size == 4) {
    use_alphatensor(A, B, C);
} else {
    use_standard_blas(A, B, C);
}
```

---

## Performance Expectations

### Benchmark Results (Reference)

| Matrix Type | AlphaTensor vs DGEMM | Use Case |
| --- | --- | --- |
| Identity | 3.91x faster | Initialization, testing |
| Mixed Sign | 4.68x faster | Neural network weights |
| Zero | 2.51x faster | Sparse operations |
| Random Dense | 1.06x faster | General ML workloads |
| Overall Average | 1.15x faster | Production workloads |

### Hardware Considerations

- **CPU**: Best performance on modern x86-64 with AVX2
- **Memory**: Benefits from L1/L2 cache efficiency
- **Compiler**: GCC 9+ or Clang 10+ recommended for vectorization

---

## Support and Community

### Getting Help

- **GitHub Issues**: lapack_ai/issues
- **Documentation**: Complete API reference and examples
- **Performance Reports**: Submit benchmark results for your hardware

### Contributing

- **Bug Reports**: Include matrix examples and performance data
- **Hardware Testing**: Help us optimize for different architectures
- **Integration Examples**: Share your production use cases

---

## Future Roadmap

### Near Term (3-6 months)

- **Python package**: `pip install lapack-alphatensor`
- **GPU implementation**: CUDA/OpenCL versions
- **Larger matrices**: 8×8, 16×16 optimizations

### Long Term (6-12 months)

- **Framework integration**: Native PyTorch/TensorFlow ops
- **Automatic tuning**: Hardware-specific optimization
- **Distributed computing**: MPI support for clusters

---