# The Anatomy of an Efficient Blackwell GEMM

Antonio Moral Villarín

## Table of contents

# Abstract

# Acknowledgements

# List of Figures and Tables

*To be auto-generated by Quarto.*

# 1 Chapter 1 – Introduction

## 1.1 Motivation and Context: The Need for Hardware–Software Co-Design

## 1.2 Challenges in Efficient Compute for AI and Edge Applications

## 1.3 Objectives and Scope of the Thesis

## 1.4 Methodology Overview

## 1.5 Structure of the Thesis

# 2 Chapter 2 – Background and Related Work

## 2.1 Evolution of GPU Architectures: From Volta to Blackwell

## 2.2 Hardware–Software Co-Design: Principles and Applications

## 2.3 General Matrix-Matrix Multiplication (GEMM) in AI Workloads

## 2.4 Domain-Specific Languages (DSLs) for GPU Programming

## 2.5 Relevant Publications and Tools (NVIDIA Research, Citadel, JAX Scaling Book, etc.)

# 3 Chapter 3 – Architecture Comparison: Hopper vs Blackwell

## 3.1 Overview of Hopper Architecture

## 3.2 Overview of Blackwell Architecture

## 3.3 Key Innovations in Blackwell

### 3.3.1 Ultra Tensor Cores and New Precision Formats (FP8, FP4)

### 3.3.2 Transformer Engine and FP4 Micro Scaling

### 3.3.3 Multi-Die Chip Design and Interconnect (NVLink, NVSwitch)

### 3.3.4 Memory System: HBM3e, L2 Cache, and Shared Memory

## 3.4 Performance/Watt and Area Efficiency Considerations

## 3.5 Summary of Architectural Differences

# 4 Chapter 4 – Metrics for GPU Efficiency

## 4.1 Performance per Watt

## 4.2 Compute Throughput by Data Type