

# Accelerating Low-Bit LLMs with Hybrid Photonic-Electrical Networking: Architectural Design for High-Efficiency Computation

## Introduction

Recent advancements in low-bit quantized computation, like the BitNet b1.58 model, have reshaped our understanding of how large language models (LLMs) can operate efficiently. By moving away from traditional floating-point weights and adopting a ternary approach—where weights can be only -1, 0, or 1—BitNet b1.58 drastically reduces the computational and memory requirements for LLMs. This paradigm shift provides an opportunity to develop an entirely new type of hardware architecture that optimizes both power and performance beyond what traditional GPU-based systems can achieve. In this document, we will explore the design of a single-chip solution that uses streamlined compute cores and integrates hybrid photonic-electrical networking to accelerate these ternary models at scale. By utilizing specialized optical components in strategic areas and simplified processing cores, we outline how this design can improve LLM economics, scalability, and power consumption compared to conventional GPU approaches.

## Architectural Overview

The proposed architecture represents a fusion of existing RISC-V compatible elements with highly specialized photonic and electrical components, along with streamlined compute cores optimized for ternary operations. This design seeks to address the unique needs of low-bit LLMs through a combination of high-density core integration, low-power compute units, and efficient data movement facilitated by hybrid electrical-photonic interconnects. Instead of complex CISC (Complex Instruction Set Computing) processors with deep pipelines and speculative execution, we propose a shift to lightweight RISC-V cores that prioritize throughput, efficient memory access, and scalable communication. The challenge is to integrate hundreds, if not thousands, of these cores into a single chip, with minimal energy consumption and optimal performance.

## Core Design for Ternary Computation

A fundamental part of the architecture is the Ternary Arithmetic Unit (TAU). In contrast to traditional GPU or CPU designs, the TAU is specifically built to handle ternary matrix multiplications involving weights that take on values from the set  $\{-1, 0, 1\}$ . This significantly reduces the number of operations needed compared to floating-point calculations. Traditional LLM models running on GPUs often involve floating-point matrix multiplications that are computationally and energy-intensive. By contrast, the ternary approach removes the need for expensive floating-point units (FPUs), leading to lower power consumption and a simpler computational pipeline.

Each TAU utilizes lookup tables (LUTs) for ternary multiplication, allowing multiplication results to be directly retrieved rather than computed in real-time. The LUT is precomputed for combinations of weights and activations, effectively reducing the multiply operation to a simple table lookup. This enables the pipeline to operate with a single cycle for most arithmetic operations, which

substantially increases throughput while reducing latency. The simplified operation can be mathematically described as follows:

$$y_{ij} = LUT[w_j + 1][x_i + 128]$$

Where LUT is indexed by adjusting the ternary value to match an appropriate offset in the table. This approach ensures that operations involving zero values are skipped entirely, enhancing computational efficiency by focusing only on non-zero contributions.

Each core is also equipped with localized Static Random-Access Memory (SRAM) clusters, shared among groups of cores, ensuring that frequently accessed data such as model weights, intermediate activations, and temporary results remain physically close to the compute units. This reduces access latency, which is often a bottleneck in GPU-based systems, where data must travel over larger distances between compute units and off-chip memory. By integrating SRAM clusters for multiple cores, memory access is kept efficient, reducing the complexity of distributed memory management.

### **Hybrid Photonic-Electrical Networking Integration**

To efficiently connect thousands of lightweight cores on a single chip, we leverage hybrid photonic-electrical networking for inter-core communication. Optical links are employed for long-distance data transfers due to their minimal heat generation and high efficiency, while electrical interconnects manage short-distance communication within core clusters to reduce design complexity and enhance manufacturing yields.

A key innovation in this architecture is Wavelength-Division Multiplexing (WDM), which allows multiple core clusters to communicate simultaneously over a single optical channel. By assigning different wavelengths of light to carry distinct data streams on the same waveguide, WDM significantly increases bandwidth without additional wiring—essential for thousands of cores requiring frequent data exchanges. For example, during LLM inference, clusters accessing and updating model weights concurrently benefit from WDM's ability to support parallel operations, reducing contention and improving throughput.

The optical interconnect topology is a scalable 3D torus or mesh, where each cluster connects to several neighbors, enabling deterministic, low-latency paths across the chip. Micro-ring resonators serve as add/drop filters, selectively coupling specific wavelengths to efficiently route data packets between clusters without converting optical signals to electrical form.

Photonic RDMA (Remote Direct Memory Access) further enhances the architecture, enabling clusters to access each other's memory directly without CPU involvement, minimizing latency. This feature is crucial for distributed LLM computations, where weights and activations must be rapidly shared for synchronized processing. Dedicated optical RDMA links handle high-bandwidth transfers, ensuring efficient communication during critical operations while maintaining system simplicity.

Each chiplet within a cluster integrates embedded SRAM for low-latency memory and processes ternary states via amplitude or phase modulation. Clusters, comprising 50–100 chiplets optimized for current Dense WDM (DWDM) technology, dynamically filter and process assigned wavelengths using tunable filters or silicon photonic resonators. Intra-cluster communication is facilitated by a photonic RDMA network with a non-blocking Clos topology for efficient signal routing, while inter-cluster connectivity is managed by a global non-blocking RDMA network.

Low-loss silicon photonic waveguides, with 50 GHz channel spacing (or 25 GHz for higher density), ensure efficient intra-cluster transfers, while guard bands and precise filtering minimize crosstalk. Together, these features deliver a robust, scalable communication framework tailored to the demands of large-scale LLM inference and high-performance computing.

### **Production Process and Packaging**

The proposed architecture will be fabricated using advanced 5nm or 3nm process nodes. These nodes are chosen to achieve a high density of transistors, enabling the integration of the massive number of cores and photonic components required for this design. However, integrating both electronic and photonic components on a single chip presents significant manufacturing challenges. Silicon photonics involves intricate waveguides and optical modulators that must be precisely aligned with the electronic circuits to ensure minimal signal loss and accurate data transmission. To simplify integration, we employ a chiplet-based approach, where photonic and electronic components are manufactured separately and later assembled using advanced packaging techniques such as chip-on-wafer or 2.5D integration. This modular approach reduces the risk of misalignment and allows for easier updates and replacements of specific components.

To overcome cooling challenges, we initially utilize traditional convection-based cooling with heat spreaders, along with active liquid cooling. Microfluidic cooling can be introduced in later versions, once the core design has been validated and production yields are stable. This approach helps maintain thermal stability without adding unnecessary complexity in the prototype phase.

### **Accelerating 1.58-Bit Models with High Efficiency**

The architectural choices made in this design are uniquely tailored to maximize the performance of 1.58-bit ternary models. One of the key benefits of using a hybrid photonic-electrical interconnect is its ability to handle the massive data movement that LLMs require. During matrix multiplication, for example, each cluster may need to communicate intermediate results with several neighboring clusters. In traditional architectures, this kind of communication can create significant bottlenecks, especially as the number of cores increases. However, with hybrid photonic networking, data can be transmitted almost instantaneously over long distances, with minimal heat generation and energy expenditure.

Furthermore, the photonic RDMA capability allows for direct memory transfers between clusters. This is particularly useful in the context of LLMs, where weights may need to be updated simultaneously across many clusters during both forward and backward passes. With photonic

RDMA, these updates can occur without the overhead of routing data through centralized memory controllers, thus reducing latency and enabling a higher degree of parallelism. For example, when a cluster computes the output of a layer, it can directly broadcast the results to all other relevant clusters using the optical interconnect, ensuring that subsequent operations are not delayed by communication bottlenecks.

The integration of shared SRAM clusters within each cluster further enhances this architecture's performance for 1.58-bit models. By keeping weights and activations close to the processing units, we reduce the need for repeated memory accesses to off-chip storage, which are typically high-latency and energy-intensive. This local storage strategy ensures that each cluster has rapid access to the data it needs, effectively removing one of the most significant bottlenecks in GPU-based LLM computations.

### **Economic Viability and Power Efficiency**

The economic benefits of this design stem from its energy efficiency and scalable production. GPUs, while powerful, consume considerable energy when performing large-scale matrix operations with floating-point precision. By contrast, this architecture leverages the simplicity of ternary operations, which not only require fewer transistors per arithmetic unit but also consume significantly less energy. The elimination of speculative execution and hyper-threading further reduces the power footprint, enabling the deployment of thousands of cores without prohibitive energy requirements.

For instance, consider an LLM inference workload that involves billions of parameter multiplications and additions. In a typical GPU, each operation involves multiple floating-point units, which consume a large amount of power and generate significant heat. In contrast, our ternary cores perform the same calculations using simple lookup tables and integer arithmetic, which drastically reduces power usage. Moreover, because the data is moved using hybrid photonic-electrical interconnects, the energy cost of communication—which is a significant portion of the overall energy budget in large-scale models—is also minimized.

### **Production Risks and Development Gaps**

While this architecture promises significant advantages, there are notable risks and challenges in the production process. One of the main risks is the complexity of integrating photonic and electronic components using chiplets. Silicon photonics is still a relatively new field, and achieving the precision required for aligning waveguides, resonators, and modulators with electronic components is challenging. By adopting a modular chiplet-based approach, we aim to mitigate these risks, but precise alignment will still be crucial to ensure minimal signal losses.

Another area of concern is thermal management. While photonic components generate less heat compared to electronic interconnects, the high density of processing cores and the integrated nature of the SRAM still result in considerable thermal output. Traditional cooling methods will be

employed initially, with more advanced microfluidic solutions being introduced in later versions, once the architecture has been validated.

There are also development gaps in terms of software support. The proposed architecture will require extensions to the RISC-V instruction set to handle ternary arithmetic and photonic communication. Additionally, popular machine learning frameworks such as PyTorch and TensorFlow would need modifications to leverage the new hardware effectively. Developing these software extensions and ensuring compatibility with existing ecosystems is a non-trivial task that will require collaboration between hardware designers, software developers, and machine learning researchers.

Finally, the testing and prototyping phase presents challenges due to the lack of mature simulation tools for large-scale photonic networks. Traditional electronic simulation tools are not equipped to handle the unique properties of optical signals, such as wavelength-specific routing and interference. Developing accurate simulation environments that can model both electronic and photonic behavior will be crucial for ensuring the reliability and performance of the final product.

## **Optimizing RISC-V Instruction Sets for Ternary Computation**

The application of RISC-V instruction sets to the proposed architecture requires a reimagining of what traditional RISC-based instruction sets can achieve in a specialized ternary computation model. RISC-V, known for its flexibility and open-source nature, provides a robust starting point. However, significant changes and optimizations are necessary to maximize the performance of the 1.58-bit ternary architecture.

### Essential Extensions to RISC-V Instructions

To effectively support the ternary arithmetic of BitNet b1.58, the RISC-V instruction set must include specific extensions that facilitate ternary operations and optimized matrix multiplications. New instructions should be introduced for the following:

#### *Ternary Multiplication and Addition:*

Instructions such as TMUL (Ternary Multiply) and TADD (Ternary Add) can be added to handle the core operations involving ternary weights and activations. These instructions should be designed to leverage lookup tables embedded within the core hardware to enable single-cycle multiplications and additions, effectively removing the need for more computationally intensive arithmetic.

#### *Lookup Table Access:*

A specialized instruction, LUTLOAD, can be introduced to facilitate rapid access to precomputed values within the lookup tables. This instruction should allow the processor to index directly into

the lookup tables used in the Ternary Arithmetic Units (TAUs) for multiplication and activation functions.

#### *Photonically Accelerated Memory Access:*

Instructions like PHOTRD (Photon Read) and PHOTWR (Photon Write) should be introduced to handle data transmission via the photonic interconnect. These instructions would allow direct access to remote memory through the photonic RDMA capabilities, bypassing traditional memory hierarchy bottlenecks.

#### *SRAM Management:*

Given that each cluster has localized SRAM, instructions for efficient data movement between SRAM and compute units are necessary. Examples include SLOAD and SSTORE for fast loading and storing of intermediate values in localized SRAM.

#### Optimization Strategies for RISC-V in Ternary Models

Optimizing the RISC-V instruction set for this architecture involves not just adding new instructions but also modifying existing ones to reduce their complexity and better align with the streamlined computational pipeline. Traditional RISC-V features like branch prediction and speculative execution should be eliminated or significantly simplified. These features, while beneficial for general-purpose computing, introduce unnecessary overhead in a deterministic ternary model environment where the data flow is highly predictable.

The focus should be on pipelined parallelism and single-cycle execution of common operations. To achieve this, legacy instructions that support floating-point arithmetic, SIMD operations beyond what is needed for ternary manipulation, and other complex arithmetic operations can be excluded or replaced. Instead, instructions should be oriented toward maximizing throughput for matrix multiplications and additions, which are the backbone of LLM inference.

#### *Instructions to Be Excluded*

In the context of the ternary computation model, certain instructions commonly found in the RISC-V instruction set would be redundant or inefficient. These include:

- **Floating-Point Instructions:** All floating-point operations (FADD, FMUL, etc.) should be excluded, as they are not relevant for the ternary operations at the core of the BitNet b1.58 model.
- **Complex Branching and Speculative Execution:** Instructions related to speculative execution should be eliminated to streamline the pipeline and reduce power consumption. Given the deterministic nature of matrix operations in LLMs, branching can be minimized.

- **Hyper-Threading Support:** Instructions that facilitate hyper-threading can also be excluded, as the architecture prioritizes massive parallelism across many simple cores rather than deep multithreading within individual cores.

By excluding these instructions, the complexity of each core is reduced, which allows for more cores to be packed into a single chip, ultimately increasing the parallelism available for LLM inference.

## Conclusion

The proposed architecture for accelerating 1.58-bit ternary LLMs represents a significant departure from traditional GPU-based systems. By combining lightweight, specialized RISC-V compute cores with hybrid photonic-electrical networking, this design offers a path to achieving both high performance and energy efficiency. The application of RISC-V instruction sets to this architecture, with targeted extensions and optimizations, further enhances its capabilities by aligning the instruction set with the specific needs of ternary arithmetic and photonic data movement. While there are significant challenges in terms of manufacturing complexity, thermal management, and software ecosystem development, the potential benefits in terms of scalability, power efficiency, and performance make this an exciting direction for future AI hardware.

Moving forward, the focus should be on developing small-scale prototypes using FPGAs to validate the integration of photonic components, enhancing simulation tools for photonic networks, and working closely with software developers to build the necessary programming frameworks. By addressing these challenges, we can create an architecture that not only meets the demands of current LLMs but also sets the stage for future advancements in AI computation.

## Addendum: Potential Improvements for BitNet b1.58

### Outline of Potential Improvements for BitNet b1.58 Design

#### Improvements to Hardware Design for Reliability and Manufacturing Yield

##### Redundancy Mechanisms:

- Implement error correction codes (ECC) in the ternary computing units to enhance reliability in environments prone to thermal noise.
- Introduce spare computational units that can take over in the event of localized failures to increase fault tolerance.

##### Manufacturing Process Optimization:

- Adopt a FinFET or GAAFET-based process node optimized for ternary logic. Advanced nodes with high yields in digital circuits.
- Modularize the ternary processing units for easier testing and assembly, reducing the risk of failure during production.

##### Thermal Management:

- Use optimized thermal dissipation materials, such as graphene-based solutions, in the chip packaging to maintain operational temperatures.

##### Integration with Existing Ecosystems:

- Standardization of Ternary Logic Components:

Create a set of guidelines for designing ternary logic gates that ensure consistent manufacturing outcomes, similar to the CMOS standards.