# Nebula Submission

## Implementation of Scalable GPU Interconnects

## with AXI/CHI Protocols & Network-on-Chip Architecture

### Technical Abstract

*Proposed Solution for Scalable AI Computing Infrastructure*

**Project Category:** Hardware/Software Co-Design
**Target Scale:** 4-64 GPUs (2x2 to 8x8 grid configurations)
**Key Technologies:** ARM AMBA AXI4/CHI, Network-on-Chip, SystemVerilog
**Implementation:** RTL, SystemC TLM-2.0, Python Analysis Framework

**Team Name:** Team Bob

**Team Members:**

Pranav Chandra
Pramit Pal
Meghadri Ghosh

August 2025

# 1   Problem Statement and Motivation

Current GPU interconnect solutions face fundamental architectural limitations when scaling beyond 8-16 GPUs, creating critical bottlenecks for large-scale AI training and inference systems. Traditional PCIe-based connections exhibit **O(1) bandwidth scaling** characteristics, where aggregate bandwidth remains constant regardless of the number of GPUs, fundamentally limiting system performance. Point-to-point topologies fail to provide efficient routing for many-to-many communication patterns prevalent in modern AI workloads, while memory coherency protocols between GPUs require extensive software management, introducing significant latency overhead.

The growing demand for larger AI models necessitates distributed training across 32+ GPUs, yet existing interconnect solutions cannot efficiently support such scales. Software-managed coherency protocols add 100-1000+ cycle overhead for inter-GPU memory operations, while PCIe fabric contention creates unpredictable latency characteristics that severely impact training convergence times.

**Our objective** is to develop a standardized, hardware-based interconnect solution that achieves **O(N) bandwidth scaling** with the number of GPUs while maintaining full protocol compatibility with existing GPU architectures and providing deterministic, low-latency communication guarantees.

# 2   Technical Approach: Nebula Submission

## 2.1   System Architecture Overview

**Overview**

This section provides a step-by-step breakdown of the system architecture, from the highest-level system view down to the hardware modules and logic that enable scalable, high-performance GPU interconnects. The goal is to make the design and implementation process clear for engineers with a basic understanding of RAM, memory, and CPU concepts. Figure 1 shows a rough layout of what will further be discussed in this abstract.
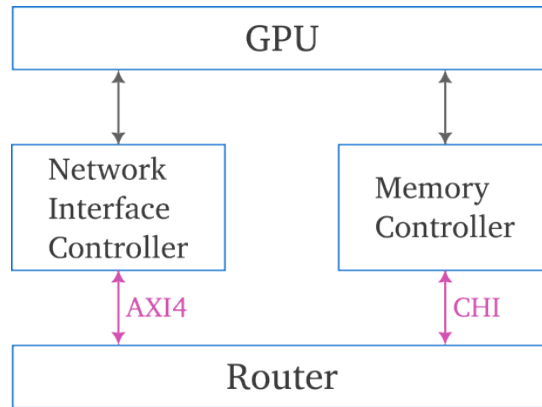


Figure 1: Top Level Architecture of the Various Components

### 2.1.1   Top-Level System Organization

At the highest level, this architecture is a hardware platform that connects multiple GPUs and memory controllers together to form a single, unified computing system. The system is organized as a grid (2D mesh) of GPUs, with each GPU able to communicate directly with any other GPU or memory controller in the grid. The system is designed to scale from 4 GPUs (2x2) up to 64 GPUs (8x8), with the ability to add or remove GPUs by changing configuration parameters.

**Key Components:**

- **GPUs**: The main processing units, each with its own local memory (RAM) and memory controller.

- **Memory Controllers**: Manage access to RAM for each GPU and handle requests from other GPUs.

- **Routers**: Hardware blocks that forward data between GPUs and memory controllers, forming the Network-on-Chip (NoC).

- **Network Interfaces**: Bridges that translate between GPU memory protocols (AXI4/CHI) and the NoC packet format.

- **NoC Links**: Physical connections (wires) that carry data between routers.

### 2.1.2 Grid and Mesh Topology

The GPUs and memory controllers are arranged in a 2D grid. Each GPU is connected to a router, and each router is connected to its neighboring routers (north, south, east, west) via NoC links. This forms a mesh topology, where data can travel from any GPU to any other by hopping through intermediate routers. Figure 2 demonstrates the NOC packet communication between routers.
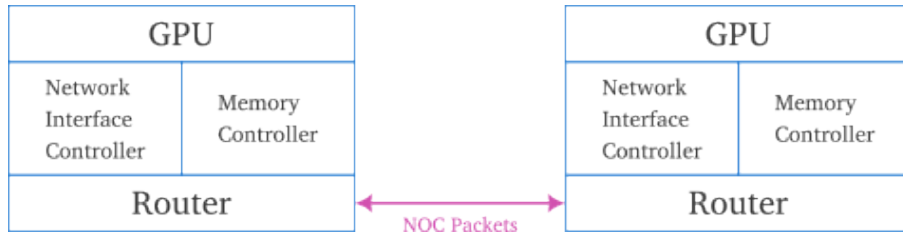


Figure 2: Example of interconnect between routers

**Implementation Steps:**

- Define the grid size (e.g., 4x4 for 16 GPUs) as a parameter in the hardware design.

- Instantiate a router for each grid position, connecting it to the local GPU and to its four neighbors (or fewer at the edges/corners).

- Connect each router to its local GPU via a network interface module.

- Connect routers to each other using point-to-point links (wires or PCB traces).

### 2.1.3 Hierarchical Clustering for Large Systems

For systems with more than 16 GPUs, the grid is divided into clusters (e.g., 4x4 clusters). Each cluster operates as a local mesh, and clusters are connected via higher-level inter-cluster links. This reduces congestion and makes routing more efficient at scale.

**Implementation Steps:**

- Partition the grid into clusters by grouping routers and GPUs into blocks (e.g., four 4x4 clusters in an 8x8 grid).

- Add inter-cluster routers and links that connect the clusters together.

- Implement address mapping logic that determines whether a packet should stay within a cluster or be forwarded to another cluster.

- Use configuration registers or tables to define cluster boundaries and routing rules.

### 2.1.4 Network Interface and Protocol Translation

Each GPU communicates using standard memory protocols (AXI4 for non-coherent, CHI for coherent transactions). The network interface module translates these memory transactions into NoC packets that can be routed through the mesh.
**Implementation Steps:**

- Instantiate an AXI4/CHI interface for each GPU, handling memory requests and responses.

- Implement a protocol translation engine that breaks up large memory transactions into smaller NoC packets, adds headers, and manages packet ordering.

- On the receive side, reassemble NoC packets into memory transactions and deliver them to the GPU or memory controller.

- Use hardware tables to track outstanding transactions and ensure correct ordering and completion.

### 2.1.5 Router Microarchitecture

Routers are the core of the NoC, responsible for forwarding packets between GPUs and memory controllers. Each router has multiple input and output ports, buffers for storing packets, and logic for deciding which direction to send each packet.
**Implementation Steps:**

- For each router, instantiate input and output ports for each direction (N, S, E, W, and local GPU).

- Implement input buffers (FIFOs) for each port to store incoming packets.

- Implement routing logic that examines the destination address in each packet and selects the appropriate output port.

- Use arbitration logic to resolve conflicts when multiple packets want to use the same output port.

- Implement credit-based flow control to prevent buffer overflows and ensure reliable delivery.

### 2.1.6 Address Mapping and Routing

To send data from one GPU to another, the system needs to know where each GPU is located in the grid and how to reach it. Address mapping logic translates memory addresses into NoC coordinates, and routing logic determines the path through the mesh.
**Implementation Steps:**

- Assign each GPU and memory controller a unique address range in the global system address space.

- Implement address decoders in hardware that extract the destination coordinates from the memory address.

- In each router, use the destination coordinates to select the next hop (direction) for each packet.

- For hierarchical systems, add logic to determine whether to route within a cluster or between clusters.

### 2.1.7 Putting It All Together: Data Flow Example

**Example: GPU 0 wants to read data from GPU 5**

1. GPU 0 issues a read request using its AXI4 interface.

2. The network interface translates the request into one or more NoC packets, adding headers with the destination coordinates for GPU 5.

3. The local router receives the packet and uses its routing logic to forward it toward GPU 5, hopping through intermediate routers as needed.

4. Each router along the path buffers the packet, checks for congestion, and forwards it to the next router.

5. When the packet reaches GPU 5's router, the network interface reassembles the packet and delivers the read request to GPU 5's memory controller.

6. The response follows the reverse path, with routers and network interfaces handling packetization, routing, and reassembly.

### 2.1.8 Hardware Implementation and Parameterization

All components (routers, network interfaces, protocol bridges) are implemented as parameterized SystemVerilog modules. This allows the system to be easily scaled up or down by changing configuration parameters (e.g., grid size, buffer depth, number of VCs).

**Implementation Steps:**

- Use SystemVerilog generate blocks to instantiate the required number of routers, network interfaces, and links based on the grid size.

- Set parameters for buffer sizes, data path widths, and protocol options at synthesis time.

- Use configuration files or registers to control address mapping, routing policies, and cluster organization.

- Synthesize and test the design on FPGA or in simulation, using testbenches that generate traffic patterns and check for correct operation.

| GPU 0,0 | GPU 1,0 | GPU 2,0 | GPU 3,0 |
|---------|---------|---------|---------|
| GPU 0,1 | GPU 1,1 | GPU 2,1 | GPU 3,1 |
| GPU 0,2 | GPU 1,2 | GPU 2,2 | GPU 3,2 |
| GPU 0,3 | GPU 1,3 | GPU 2,3 | GPU 3,3 |

Figure 3: Example 4x4 Grid Topology (16 GPUs)

## 2.2 Component 1: AXI4/CHI Protocol Implementation

### Overview

This section describes the practical implementation of the AXI4 and CHI protocol interfaces, which are responsible for enabling high-throughput, low-latency, and cache-coherent communication between GPUs and memory controllers in the NoC. The goal is to provide hardware-managed

transaction tracking, burst handling, and coherency state management, ensuring protocol compliance and scalability for large GPU grids. The following subsections detail the design and step-by-step implementation of the AXI4 interface state machines and the CHI coherency engine.

### 2.2.1 AXI4 Interface State Machine Design

The AXI4 interface is implemented as a set of five independent finite state machines (FSMs), one for each channel: Address Write (AW), Write Data (W), Write Response (B), Address Read (AR), and Read Data (R). Each FSM is responsible for protocol handshaking, data transfer, and error handling for its respective channel. The design is parameterized for scalability and protocol compliance.

**Implementation Details:**

- **Outstanding Transaction Management**:

  – A hardware table (implemented as a dual-port RAM or register array) is instantiated per interface to track all in-flight transactions. Each entry stores the transaction ID (TID), address, burst length, and status.

  – The table depth is parameterizable (default 16, can be increased for larger systems). The TID field width is set based on the maximum number of outstanding transactions required for the target grid size.

  – On receiving a new AW or AR command, a free entry is allocated and populated. On completion (B or R response), the entry is released.

  – The table is accessed by both the AW/AR and B/R FSMs, with arbitration logic to prevent hazards.

- **Burst Transaction Handling**:

  – The AXI4-to-NoC bridge segments large AXI bursts (up to 256 beats) into multiple NoC packets. Each packet contains a header with the TID, sequence number, and burst offset.

  – The bridge maintains a per-transaction counter to track the number of beats sent and received.

  – On the receive side, packets are reassembled into AXI bursts using a reorder buffer, ensuring correct ordering and data integrity.

  – Boundary protection logic ensures that bursts do not cross 4KB address boundaries, as required by the AXI4 specification.

- **Data Path Structure**:

  – The 512-bit data path is implemented as 8 parallel 64-bit lanes, each with its own byte-enable and data steering logic.

  – Multiplexers and demultiplexers are used to map AXI data to NoC flits and vice versa, supporting both aligned and unaligned accesses.

  – The number of lanes and lane width are set via SystemVerilog parameters, allowing easy scaling for different memory and NoC configurations.

  – The design supports both little-endian and big-endian data formats, with configurable byte swapping.

- **Credit-Based Flow Control**:

  – Each input buffer (per VC) in the router and interface maintains a credit counter, initialized to the buffer depth.

- When a flit is sent, the sender decrements its credit counter. When the receiver consumes a flit, it sends a credit return signal to the sender.
- If credits reach zero, the sender stalls until more credits are received, preventing buffer overflow and ensuring lossless operation.
- The credit return path is implemented as a dedicated signal or piggybacked on data flits, depending on timing constraints.

- **Error Handling and Protocol Compliance**:

  - Each FSM includes error detection logic for protocol violations (e.g., unexpected responses, burst length mismatches, invalid TIDs).
  - Detected errors are logged in a status register and can trigger interrupts or error responses on the AXI bus.
  - SystemVerilog assertions are used to verify protocol compliance during simulation and formal verification.

### 2.2.2 CHI Coherency Engine Architecture

### 2.2.3 CHI Coherency Engine Architecture

The CHI (Coherent Hub Interface) engine is responsible for hardware-managed cache coherency across all GPUs and memory controllers. It implements the five-state protocol (I, UC, UD, SC, SD) and manages all coherency requests, snoop filtering, and directory updates. The engine is split into Request Node (RN) and Home Node (HN) components, each with dedicated hardware logic.

**Implementation Details:**

- **Request Node (RN) Engine**:

  - Each GPU instantiates an RN FSM that handles all outgoing cache line requests (read, write, atomic, etc.).
  - The RN maintains a snoop filter, implemented as a hardware directory (RAM), which tracks the ownership and sharing status of each cache line.
  - On a cache miss, the RN issues a CHI request to the HN, including the line address and request type. The RN also initiates snoop requests to other GPUs as needed.
  - The RN FSM manages state transitions for each cache line, updating the snoop filter and directory based on incoming responses and acknowledgments.
  - Atomic updates are handled by locking the relevant directory entry until all snoop responses are received.

- **Home Node (HN) Engine**:

  - Each memory controller instantiates an HN FSM, which acts as the central directory for all cache lines mapped to its address range.
  - The HN maintains a directory RAM, with one entry per cache line, storing the current state (I, UC, UD, SC, SD) and a list of sharers/owners.
  - On receiving a CHI request, the HN serializes the transaction, updates the directory state, and issues snoop requests to relevant RNs as needed.
  - The HN tracks outstanding transactions in a queue, ensuring correct ordering and completion semantics.
  - The HN generates completion responses to the requesting RN once all snoop responses are collected and the directory is updated.

- **Five-State Protocol**:
  - Each cache line is assigned a 3-bit state field: Invalid (I), Unique Clean (UC), Unique Dirty (UD), Shared Clean (SC), Shared Dirty (SD).
  - State transitions are triggered by incoming requests (read, write, invalidate, etc.) and are managed by the RN and HN FSMs.
  - The protocol ensures that only one node can hold a dirty (modified) copy at a time, and that all sharers are kept coherent via snoop and invalidate messages.
  - Message ordering is enforced by hardware queues and sequence numbers, preventing race conditions and ensuring protocol correctness.

- **Snoop Broadcast and Acknowledgment**:
  - Invalidation and snoop requests are broadcast using multicast logic in the NoC. The HN or RN generates a multicast packet addressed to all relevant nodes.
  - Each node receiving a snoop request responds with an acknowledgment, which is tracked by the sender using a bitmap or counter.
  - The transaction is only completed once all required acknowledgments are received, ensuring global coherency.
  - Outstanding snoops and acknowledgments are tracked in hardware using per-transaction state machines and counters.

- **Error Handling and Debug**:
  - Protocol violations (e.g., unexpected state transitions, missing acknowledgments) are detected by assertions and logged in error registers.
  - Debug interfaces allow inspection of directory state, outstanding transactions, and snoop filter contents for verification and troubleshooting.

## 2.3 Component 2: Network-on-Chip Router Microarchitecture

### Overview

This section details the practical implementation of the Network-on-Chip (NoC) router microarchitecture, which is central to achieving scalable, low-latency, and deadlock-free communication between GPUs in the system. The router is designed as a highly parameterized SystemVerilog module, supporting flexible grid sizes, buffer depths, and virtual channel (VC) counts. The following subsections describe the five-stage pipeline, buffer architecture, and routing algorithms, with explicit implementation steps and design choices.

### 2.3.1 Five-Stage Router Pipeline

Each router is implemented as a parameterized SystemVerilog module, with the pipeline split into five distinct stages. Each stage is realized as a separate always_ff block, with pipeline registers between stages for timing closure and modularity.
  **Implementation Details:**

1. **Buffer Write (BW)**:
   - Incoming flits are received on each input port and validated for protocol compliance (e.g., correct header, CRC check).
   - Flits are written into the appropriate VC FIFO buffer, selected based on the traffic class and VC allocation policy.

- If the buffer is full, backpressure is asserted to the upstream router or interface.
- Error detection logic flags malformed or unexpected flits, which are dropped or logged for debug.

2. **Route Computation (RC)**:

   - For each head flit in the VC buffer, the router compares the destination coordinates (extracted from the flit header) with its own coordinates.
   - The default routing algorithm is dimension-ordered XY routing: if the x-destination differs, route east/west; otherwise, route north/south.
   - If the preferred output port is congested (as indicated by buffer occupancy or congestion signals), adaptive logic selects an alternative output direction.
   - The selected output port and next VC are stored in a per-flit routing table for use in later pipeline stages.

3. **Virtual Channel Allocation (VA)**:

   - Arbitration logic selects which VC(s) are eligible to advance to the switch allocation stage, based on round-robin, priority, or age-based policies.
   - For each eligible flit, the router requests allocation of a downstream VC on the selected output port.
   - VC state machines (per VC) track the state of each channel: Idle, Routing, Active, Waiting. State transitions are triggered by flit arrival, allocation, and credit availability.
   - If no downstream VC is available, the flit remains in the buffer and is retried in the next cycle.

4. **Switch Allocation (SA)**:

   - The router contains a crossbar switch connecting all input VCs to all output ports.
   - Arbitration is performed among competing VCs for each output port, using round-robin or weighted priority to ensure fairness and prevent starvation.
   - Grant history is tracked to ensure that all VCs eventually make progress, even under heavy load.
   - The winning flit is selected for transmission in the next stage.

5. **Switch Traversal (ST)**:

   - The selected flit is transmitted across the crossbar to the output port, with the appropriate VC and output buffer selected.
   - Credit counters for the downstream buffer are decremented, and credit return signals are generated as flits are consumed downstream.
   - Error monitoring logic checks for protocol violations (e.g., credit underflow, invalid VC assignment) and logs or flags errors as needed.
   - The pipeline is fully stallable: if any downstream resource is unavailable, the affected pipeline stage stalls, and backpressure propagates upstream.

### 2.3.2 Advanced Buffer Architecture

### 2.3.3 Advanced Buffer Architecture

The buffer architecture is designed for high throughput, deadlock avoidance, and efficient resource utilization. Each input port is equipped with multiple VCs, each with its own FIFO buffer and state machine.

**Implementation Details:**

- **Multi-Level Buffering**:

  - Each input port has a parameterizable number of VCs (default 4), each with a parameterizable FIFO depth (default 16 flits).
  - The FIFOs are implemented as dual-port RAMs or register arrays, with separate read/write pointers for efficient access.
  - Optionally, a shared buffer pool can be enabled, allowing dynamic allocation of buffer space among VCs to maximize utilization under bursty traffic.

- **VC State Machines**:

  - Each VC is managed by a dedicated FSM with four states: Idle (waiting for flit), Routing (route computation in progress), Active (allocated and transmitting), Waiting (blocked on downstream resource).
  - State transitions are triggered by flit arrival, successful VC/switch allocation, and credit availability.
  - The FSMs are implemented as always_ff blocks, with state encoded in registers for synthesis efficiency.

- **Credit-Based Flow Control**:

  - Each VC maintains a credit counter, initialized to the buffer depth. Credits are decremented on flit send and incremented on credit return from downstream.
  - If credits reach zero, the VC stalls and backpressure is asserted to the upstream router or interface.
  - Credit return signals are implemented as dedicated wires or piggybacked on data flits, depending on timing and area constraints.
  - Credit underflow/overflow is detected and flagged as an error for debug and verification.

- **Congestion Monitoring**:

  - Buffer occupancy is tracked in hardware for each VC and input port.
  - If occupancy exceeds a programmable threshold, a congestion signal is asserted to the route computation stage and to upstream routers.
  - Congestion information is used by the adaptive routing logic to select less congested paths, improving load balancing and throughput.
  - Buffer occupancy statistics are exported to the performance monitoring infrastructure for analysis.

### 2.3.4  Routing Algorithm Implementation

### 2.3.5  Routing Algorithm Implementation

The router supports both deterministic and adaptive routing algorithms, with hardware support for deadlock avoidance, load balancing, and congestion mitigation.
**Implementation Details:**

- **XY Routing**:

  - The default routing algorithm is dimension-ordered XY routing, implemented as combinational logic in the route computation stage.
  - The router compares its own (x, y) coordinates with the destination (x, y) in the flit header. If x-destination > x, the flit is routed east; if x-destination < x, west; otherwise, y-direction is used.
  - This approach guarantees deadlock freedom by strictly ordering traversal in the x and then y dimensions.
  - The routing logic is parameterized to support different mesh sizes and topologies.

- **Adaptive Routing**:

  - If the preferred output port is congested (as indicated by buffer occupancy or explicit congestion signals), the router evaluates alternative output directions.
  - The adaptive logic selects the least congested available output, using buffer occupancy and recent grant history as inputs.
  - Congestion status is propagated to upstream routers via dedicated signals, enabling global load balancing.
  - The adaptive routing policy is configurable (e.g., minimal adaptive, fully adaptive, or turn model).

- **Load Balancing**:

  - The router tracks per-port utilization and buffer occupancy statistics in hardware.
  - Weighted random selection is used among non-congested outputs to distribute traffic and avoid hotspots.
  - Hotspot detection logic flags ports with sustained high occupancy, triggering adaptive rerouting.
  - Load balancing parameters are programmable via configuration registers.

- **VC Ordering and Deadlock Avoidance**:

  - Message classes (e.g., requests, responses, data) are mapped to specific VCs to prevent cyclic dependencies and ensure deadlock freedom.
  - Dependency tracking logic monitors VC usage and enforces ordering constraints, blocking flits that would create cyclic waits.
  - The VC allocation policy is designed to guarantee progress for all message classes under all traffic conditions.
  - Deadlock and livelock detection logic is included for verification and debug.

## 2.4 Component 3: Network Interface Protocol Translation

**Overview**

This section describes the practical implementation of the protocol translation logic that bridges the AXI4/CHI interfaces and the packet-switched NoC. The goal is to enable seamless, high-throughput conversion between burst-oriented memory protocols and the flit-based NoC, while preserving ordering, coherency, and error detection. The translation engines are implemented as dedicated hardware modules, with parameterizable tables and logic for scalability.

### 2.4.1 AXI-to-NoC Translation Engine

The AXI-to-NoC translation engine converts AXI4 memory transactions into NoC packets and vice versa, handling burst decomposition, address mapping, and packet assembly.

**Implementation Details:**

- **Transaction Decomposition**:

  - Large AXI bursts (up to 256 beats) are segmented into multiple NoC packets, each carrying a subset of the burst data.
  - Each packet includes a header with the transaction ID (TID), sequence number, burst offset, and length.
  - The translation engine maintains a per-transaction counter and a reorder buffer to track outstanding packets and ensure in-order reassembly.
  - On the receive side, packets are reassembled into AXI bursts, with error checking for missing or out-of-order packets.

- **State Tracking**:

  - A hardware table (RAM) tracks up to 64 outstanding operations, with fields for TID, address, burst length, sequence number, and timeout counters.
  - Each entry is allocated on a new AXI command and released when all corresponding NoC packets are acknowledged and reassembled.
  - Timeout logic detects stalled or lost packets and triggers error handling or retransmission as needed.

- **Address Mapping**:

  - AXI addresses are mapped to NoC destination coordinates using a configurable address decoder, implemented as combinational logic or a small lookup table.
  - The mapping supports both static partitioning (fixed address ranges per GPU) and dynamic schemes (hashing, modulo, etc.) for load balancing.
  - The decoder logic is synthesized from a configuration file, allowing easy adaptation to different system topologies.

- **Packet Assembly**:

  - Each NoC packet includes a header (destination, sequence, class), payload (data or control), and CRC for error detection.
  - Header and payload are generated in hardware, with fields populated from the AXI transaction and address mapping logic.
  - CRC is computed using a parallel LFSR (Linear Feedback Shift Register) for high-speed operation.

- The translation engine supports both single-flit and multi-flit packets, with config-
  urable payload sizes.

- **Error Handling and Debug**:

  - Protocol violations (e.g., missing packets, CRC errors, timeout) are detected and
    logged in error registers.
  - Debug interfaces allow inspection of outstanding transactions, reorder buffers, and
    error counters for verification and troubleshooting.

### 2.4.2 CHI-to-NoC Protocol Mapping

### 2.4.3 CHI-to-NoC Protocol Mapping

The CHI-to-NoC protocol mapping engine ensures that cache-coherent CHI transactions are
correctly translated into NoC packets, preserving message ordering, coherency semantics, and
error detection.
**Implementation Details:**

- **Message Classification**:

  - CHI messages (requests, responses, data, snoops) are classified in hardware and
    mapped to specific VCs based on message type and priority.
  - The mapping logic is implemented as a combinational decoder, with configuration
    registers for traffic class assignment.
  - Priority is assigned by message class, with critical coherency messages given higher
    priority to minimize latency.

- **Coherency State Management**:

  - Each node maintains a directory RAM for cache line state, updated by the CHI
    engine and protocol mapping logic.
  - Snoop traffic is optimized by only sending requests to relevant nodes, as determined
    by the directory and snoop filter.
  - The protocol mapping engine tracks outstanding snoop requests and acknowledgments,
    ensuring completion before transaction finalization.

- **Cache Line Handling**:

  - 64-byte cache lines are segmented into multiple flits, each with a checksum for error
    detection.
  - Reassembly logic at the destination checks integrity and reconstructs the full cache
    line before passing it to the CHI engine.
  - The segmentation and reassembly logic is parameterized for different cache line and
    flit sizes.

- **Ordering Preservation**:

  - Sequence numbers and dependency tracking are used to ensure that CHI completion
    semantics are maintained across the NoC.
  - The protocol mapping engine enforces ordering constraints for requests, responses,
    and data, blocking or reordering packets as needed.
  - Dependency tracking logic is implemented as a small hardware table, indexed by
    transaction ID and message class.

- **Error Handling and Debug**:

  - Protocol violations (e.g., missing acknowledgments, out-of-order completions, checksum errors) are detected and logged in error registers.
  - Debug interfaces allow inspection of directory state, outstanding transactions, and flit buffers for verification and troubleshooting.

## 2.5 Component 4: System Integration Architecture

### Overview

This section describes the practical implementation of the system integration architecture, which ties together the NoC, protocol bridges, and GPU/memory interfaces into a unified, scalable system. The focus is on address mapping, routing, collective operations, and memory consistency, with hardware support for performance optimization and system-level features.

### 2.5.1 Hierarchical Addressing Scheme

The hierarchical addressing scheme enables scalable, efficient routing and memory access across the entire GPU grid, supporting both local and global operations.

**Implementation Details:**

- **Global Address Space**:

  - The system address space is partitioned into fixed ranges, with each GPU and memory controller assigned a unique region.
  - Address translation logic in hardware maps global addresses to NoC destination coordinates, using a combination of bit extraction and lookup tables.
  - The mapping supports both flat and hierarchical topologies, with configuration registers for flexible assignment.

- **Route Computation**:

  - Routers extract destination coordinates from the address field in the packet header, using bit masks and shifts.
  - For hierarchical topologies, additional fields in the header indicate cluster and local node IDs, enabling multi-level routing.
  - The route computation logic is parameterized for different mesh and cluster sizes.

- **Collective Operations**:

  - Broadcast and multicast are implemented using hardware support for multi-destination packet replication.
  - Special address ranges or packet types trigger collective logic, which generates and routes copies of the packet to all relevant nodes.
  - The collective engine tracks outstanding responses and ensures completion before releasing resources.
  - Hardware barriers and reduction operations can be supported by extending the collective logic.

- **Memory Consistency**:

  - Ordering is enforced by hardware barriers, sequence numbers, and dependency tracking in the protocol bridges and routers.

– The system supports both strong and relaxed consistency models, selectable via configuration registers.

– Consistency violations (e.g., out-of-order completions) are detected and flagged for debug and verification.

- **Error Handling and Debug**:

  – Address mapping errors, routing failures, and collective operation timeouts are detected and logged in error registers.

  – Debug interfaces allow inspection of address translation tables, route computation logic, and collective engine state.

### 2.5.2 Performance Optimization Mechanisms

### 2.5.3 Performance Optimization Mechanisms

The system includes a comprehensive set of hardware mechanisms to optimize performance for AI workloads, focusing on routing, QoS, transaction handling, and VC management.
**Implementation Details:**

- **Adaptive Routing**:

  – Real-time congestion monitoring is implemented in each router, with congestion signals propagated to upstream nodes.

  – The routing logic dynamically selects the least congested path, balancing load and minimizing latency.

  – Congestion thresholds and routing policies are programmable via configuration registers.

- **Quality of Service (QoS)**:

  – Traffic is classified into multiple classes (e.g., control, data, collective) at the protocol bridge.

  – Priority queuing and bandwidth allocation are implemented in the routers, with per-class counters and schedulers.

  – QoS parameters are configurable, allowing tuning for different workload requirements.

- **Transaction Optimization**:

  – Request combining logic merges adjacent or overlapping memory requests to reduce packet count and improve efficiency.

  – Response batching groups multiple responses into a single packet when possible, reducing overhead.

  – Predictive prefetching logic anticipates future requests based on access patterns, issuing prefetches to hide memory latency.

- **Virtual Channel Management**:

  – Traffic classes are mapped to separate VCs, preventing head-of-line blocking and ensuring isolation.

  – VC allocation and arbitration policies are programmable, supporting both static and dynamic schemes.

– VC usage statistics are exported to the performance monitoring infrastructure for analysis and tuning.

- **Error Handling and Debug**:

  – Performance anomalies (e.g., excessive congestion, QoS violations) are detected and logged in error registers.

  – Debug interfaces allow inspection of routing tables, QoS schedulers, and transaction optimizers for verification and tuning.

hello world!