# Network Architecture for FPGA-Based Event Filtering at $\overline{\text{P}}$ANDA/FAIR

Sören Fleischer

December 10, 2013

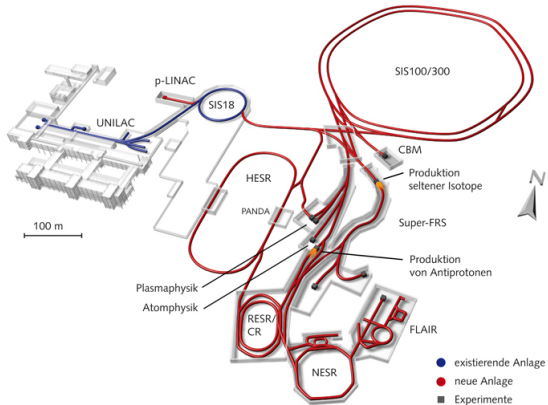# Inhaltsverzeichnis

## A Detector Experiment



Figure: Topology of $\overline{P}$ANDA at FAIR

# A Detector Experiment



Figure: The $\overline{\text{P}}$ANDA Detector

## A Detector Experiment



Figure: $\overline{P}$ANDA's Electromagnetic Calorimeter

## A Detector Experiment



Figure: $\overline{P}$ANDA's Micro Vertex Detector

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

**General Idea**
Compute Node v3
Compute Node Network Topology

# Online Data Reduction



Figure: Online Data Reduction

Motivation
Online Data Reduction
Data Transport
Implementation
To Do

General Idea
Compute Node v3
Compute Node Network Topology

# Compute Node v3



Figure: CN v3 Mother Board

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
**Compute Node v3**
Compute Node Network Topology

## Compute Node v3



Figure: CN v3 Daughter Board, no RAM installed

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
**Compute Node v3**
Compute Node Network Topology

# Compute Node v3



Figure: CN v3 Daughter Board, connectors

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
Compute Node v3
**Compute Node Network Topology**

# Compute Node Network Topology



Figure: Fixed data paths on each CN

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
Compute Node v3
**Compute Node Network Topology**

# Compute Node Network Topology



Figure: Possible topology in a shelf of 6 CNs

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
Compute Node v3
**Compute Node Network Topology**

# Example of a Tree network



Figure: Example of a Tree network

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
Compute Node v3
**Compute Node Network Topology**

# Network Topologies



Figure: Network topologies (example graphs)

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
Compute Node v3
**Compute Node Network Topology**

# Network Topologies



(a) A partial mesh network

(b) Path 1

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
Compute Node v3
**Compute Node Network Topology**

# Network Topologies



(a) A partial mesh network

(b) Path 1

(c) Path 2

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
Compute Node v3
**Compute Node Network Topology**

# Network Topologies



(a) A partial mesh network

(b) Path 1

(c) Path 2

(d) Path 3

Motivation
**Online Data Reduction**
Data Transport
Implementation
To Do

General Idea
Compute Node v3
**Compute Node Network Topology**

# Network Topologies



(a) A partial mesh network

(b) Path 1

(c) Path 2

(d) Path 3

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

**Requirements, Practical Considerations**
Circuit Switching or Packet Switching
Routing
Overview

## Requirements

- Scalability
  Some 100 FPGAs

- Flexibility
  Changing topology without changing any bitstream

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

**Requirements, Practical Considerations**
Circuit Switching or Packet Switching
Routing
Overview

## Requirements

- Scalability
  Some 100 FPGAs

- Flexibility
  Changing topology without changing any bitstream

- Being lightweight
  Don't consume too much precious Logic Cells on the FPGAs

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

**Requirements, Practical Considerations**
Circuit Switching or Packet Switching
Routing
Overview

## Requirements

- Scalability
  Some 100 FPGAs

- Flexibility
  Changing topology without changing any bitstream

- Being lightweight
  Don't consume too much precious Logic Cells on the FPGAs

- Speed/Throughput
  The maximum speed of Xilinx' RocketIO of $6.5\frac{Gb}{sec}$ should be
  achieved

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

**Requirements, Practical Considerations**
Circuit Switching or Packet Switching
Routing
Overview

## Requirements

- Scalability
  Some 100 FPGAs

- Flexibility
  Changing topology without changing any bitstream

- Being lightweight
  Don't consume too much precious Logic Cells on the FPGAs

- Speed/Throughput
  The maximum speed of Xilinx' RocketIO of $6.5\frac{Gb}{\sec}$ should be achieved

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
**Circuit Switching or Packet Switching**
Routing
Overview

# Circuit Switching or Packet Switching

- Circuit Switching
  - Easier to implement

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
**Circuit Switching or Packet Switching**
Routing
Overview

# Circuit Switching or Packet Switching

- Circuit Switching
  - Easier to implement
  - Guaranteed speed for a given circuit

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
**Circuit Switching or Packet Switching**
Routing
Overview

# Circuit Switching or Packet Switching

- Circuit Switching
    - Easier to implement
    - Guaranteed speed for a given circuit
    - "Line Busy" error condition

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
**Circuit Switching or Packet Switching**
Routing
Overview

# Circuit Switching or Packet Switching

- Circuit Switching
  - Easier to implement
  - Guaranteed speed for a given circuit
  - "Line Busy" error condition

- Packet Switching

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
**Circuit Switching or Packet Switching**
Routing
Overview

## Circuit Switching or Packet Switching

- Circuit Switching

  - Easier to implement
  - Guaranteed speed for a given circuit
  - "Line Busy" error condition

- Packet Switching

  - Harder to implement

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
Circuit Switching or Packet Switching
Routing
Overview

# Circuit Switching or Packet Switching

- Circuit Switching
    - Easier to implement
    - Guaranteed speed for a given circuit
    - "Line Busy" error condition
- Packet Switching
    - Harder to implement
    - No guaranteed speed

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
**Circuit Switching or Packet Switching**
Routing
Overview

# Circuit Switching or Packet Switching

- Circuit Switching
    - Easier to implement
    - Guaranteed speed for a given circuit
    - "Line Busy" error condition

- Packet Switching
    - Harder to implement
    - No guaranteed speed
    - Might require more buffers (Block RAM) on the FPGAs

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
Circuit Switching or Packet Switching
Routing
Overview

# Circuit Switching or Packet Switching

- Circuit Switching
    - Easier to implement
    - Guaranteed speed for a given circuit
    - "Line Busy" error condition

- Packet Switching
    - Harder to implement
    - No guaranteed speed
    - Might require more buffers (Block RAM) on the FPGAs
    - Throttling mechanism

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
Circuit Switching or Packet Switching
Routing
Overview

# Circuit Switching or Packet Switching

- Circuit Switching
    - Easier to implement
    - Guaranteed speed for a given circuit
    - "Line Busy" error condition

- Packet Switching
    - Harder to implement
    - No guaranteed speed
    - Might require more buffers (Block RAM) on the FPGAs
    - Throttling mechanism
    - Packet packing/unpacking overhead

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
**Circuit Switching or Packet Switching**
Routing
Overview

# Circuit Switching or Packet Switching

- Circuit Switching
    - Easier to implement
    - Guaranteed speed for a given circuit
    - "Line Busy" error condition

- Packet Switching
    - Harder to implement
    - No guaranteed speed
    - Might require more buffers (Block RAM) on the FPGAs
    - Throttling mechanism
    - Packet packing/unpacking overhead

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
**Circuit Switching or Packet Switching**
Routing
Overview

# Circuit Switching or Packet Switching



Figure: Circuit switching in the olden days

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
Circuit Switching or Packet Switching
**Routing**
Overview

## Routing

- Where should the routing take place?
- Not on the precious FPGAs. Instead, on a PC which sends the routes to the FPGAs.

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
Circuit Switching or Packet Switching
**Routing**
Overview

## Routing

- Where should the routing take place?
- Not on the precious FPGAs. Instead, on a PC which sends the routes to the FPGAs.
- Dijkstra's Algorithm for route calculation

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
Circuit Switching or Packet Switching
**Routing**
Overview

## Routing

- Where should the routing take place?
- Not on the precious FPGAs. Instead, on a PC which sends the routes to the FPGAs.
- Dijkstra's Algorithm for route calculation

Motivation
Online Data Reduction
**Data Transport**
Implementation
To Do

Requirements, Practical Considerations
Circuit Switching or Packet Switching
Routing
**Overview**

# Complete Reduction System



Supervisor PC

Ethernet Switch

shelf manager    shelf manager    ...

Figure: High level diagram of the reduction system

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

**Overview**
Crossbar Switch
Dialer
Resource Consumption

## One Compute Node



Figure: One Compute Node

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

# One FPGA#0



Figure: One Switching FPGA (#0)

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

# Crossbar Switch

- Has $n$ interfaces
- For each interface, it has

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

# Crossbar Switch

- Has *n* interfaces
- For each interface, it has
  - A status (st_idle, st_a_connected, st_b_connected, st_lo_connected)

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

## Crossbar Switch

- Has $n$ interfaces
- For each interface, it has
  - A status (st_idle, st_a_connected, st_b_connected, st_lo_connected)
  - A target interface number

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

## Crossbar Switch

- Has $n$ interfaces
- For each interface, it has
    - A status (st_idle, st_a_connected, st_b_connected, st_lo_connected)
    - A target interface number
- For connected interfaces: Copies data from source interface to target interface

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

## Crossbar Switch

- Has $n$ interfaces
- For each interface, it has
    - A status (st_idle, st_a_connected, st_b_connected, st_lo_connected)
    - A target interface number
- For connected interfaces: Copies data from source interface to target interface
- For idle interfaces: Waits for requests to establish a connection or to identify oneself

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

## Crossbar Switch

- Has *n* interfaces
- For each interface, it has
    - A status (st_idle, st_a_connected, st_b_connected, st_lo_connected)
    - A target interface number
- For connected interfaces: Copies data from source interface to target interface
- For idle interfaces: Waits for requests to establish a connection or to identify oneself

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

# Example for Crossbar Switch state signals

| Interface number | Interface status | Target interface |
|:---:|:---|:---:|
| 1 | st_idle | 0 |
| 2 | st_a_connected | 4 |
| 3 | st_lo_connected | 3 |
| 4 | st_b_connected | 2 |

Table: Example of the content of target and interface status information within a 4-interface crossbar switch

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

# State signals and 2 state machines



Figure: Schematic data flow inside a crossbar switch

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
**Crossbar Switch**
Dialer
Resource Consumption

# State Diagram of the Crossbar Switch

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

# One FPGA#1-#4



Figure: One Algorithm FPGA (#1-#4)

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

# Data flow between algorithms



Figure: Schematic data flow from sender algorithm to receiver algorithm

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

## Dialer

- Establishes a connection to a receiver algorithm, as requested by the associated sender algorithm.
- Has a "phone book" of receiver algorithms who the associated sender algorithm wants to send data to. Consisting of:

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

## Dialer

- Establishes a connection to a receiver algorithm, as requested by the associated sender algorithm.
- Has a "phone book" of receiver algorithms who the associated sender algorithm wants to send data to. Consisting of:
  - ID/Class of the receiver algorithm

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

## Dialer

- Establishes a connection to a receiver algorithm, as requested by the associated sender algorithm.
- Has a "phone book" of receiver algorithms who the associated sender algorithm wants to send data to. Consisting of:
    - ID/Class of the receiver algorithm
    - The route to the receiver algorithm

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

# Dialer

- Establishes a connection to a receiver algorithm, as requested by the associated sender algorithm.
- Has a "phone book" of receiver algorithms who the associated sender algorithm wants to send data to. Consisting of:
    - ID/Class of the receiver algorithm
    - The route to the receiver algorithm
    - Error counters for problems which can arise while establishing that connection:

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

## Dialer

- Establishes a connection to a receiver algorithm, as requested by the associated sender algorithm.
- Has a "phone book" of receiver algorithms who the associated sender algorithm wants to send data to. Consisting of:
    - ID/Class of the receiver algorithm
    - The route to the receiver algorithm
    - Error counters for problems which can arise while establishing that connection:
        - Line busy

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

# Dialer

- Establishes a connection to a receiver algorithm, as requested by the associated sender algorithm.
- Has a "phone book" of receiver algorithms who the associated sender algorithm wants to send data to. Consisting of:
    - ID/Class of the receiver algorithm
    - The route to the receiver algorithm
    - Error counters for problems which can arise while establishing that connection:
        - Line busy
        - Out of range

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

## Dialer

- Establishes a connection to a receiver algorithm, as requested by the associated sender algorithm.
- Has a "phone book" of receiver algorithms who the associated sender algorithm wants to send data to. Consisting of:
    - ID/Class of the receiver algorithm
    - The route to the receiver algorithm
    - Error counters for problems which can arise while establishing that connection:
        - Line busy
        - Out of range
        - Timeout

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

## Dialer

- Establishes a connection to a receiver algorithm, as requested by the associated sender algorithm.
- Has a "phone book" of receiver algorithms who the associated sender algorithm wants to send data to. Consisting of:
  - ID/Class of the receiver algorithm
  - The route to the receiver algorithm
  - Error counters for problems which can arise while establishing that connection:
    - Line busy
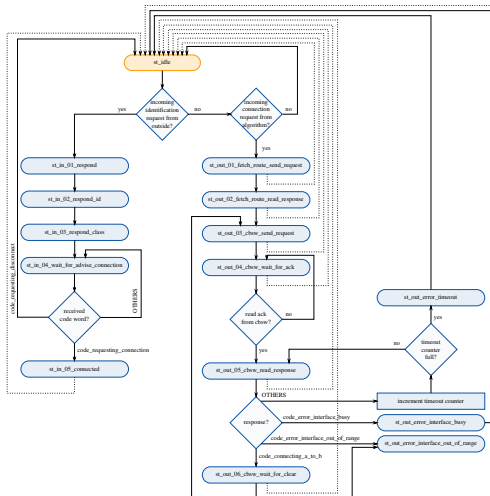    - Out of range
    - Timeout

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

## Example for a Dialer Memory



Figure: Example of a dialer memory with g_num_routes=4

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

# Dialer State Diagram

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
**Dialer**
Resource Consumption

# Dialer State Diagram



Figure: Dialer Flow Chart (Part 2)

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
Dialer
**Resource Consumption**

# Test project with 9 CBSWs

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
Dialer
**Resource Consumption**

# Test project with 9 CBSWs

| Project File: | final.xise | Parser Errors: | No Errors |
|---|---|---|---|
| Module Name: | topmodule_9x4_with_broad_sc | Implementation State: | Programming File Generated |
| Target Device: | xc4vfx60-10ff672 | • Errors: | No Errors |
| Product Version: | ISE 14.1 | • Warnings: | 9923 Warnings (9901 new) |
| Design Goal: | Balanced | • Routing Results: | All Signals Completely Routed |
| Design Strategy: | Xilinx Default (unlocked) | • Timing Constraints: | All Constraints Met |
| Environment: | System Settings | • Final Timing Score: | 0 (Timing Report) |

| Device Utilization Summary | | | | [-] |
|---|---|---|---|---|
| **Logic Utilization** | **Used** | **Available** | **Utilization** | **Note(s)** |
| Total Number Slice Registers | 8,786 | 50,560 | 17% | |
| Number used as Flip Flops | 8,591 | | | |
| Number used as Latches | 195 | | | |
| Number of 4 input LUTs | 18,124 | 50,560 | 35% | |
| Number of occupied Slices | 12,374 | 25,280 | 48% | |
| Number of Slices containing only related logic | 12,374 | 12,374 | 100% | |
| Number of Slices containing unrelated logic | 0 | 12,374 | 0% | |
| Total Number of 4 input LUTs | 18,487 | 50,560 | 36% | |
| Number used as logic | 18,111 | | | |
| Number used as a route-thru | 363 | | | |
| Number used as Shift registers | 13 | | | |
| Number of bonded IOBs | 35 | 352 | 9% | |
| Number of BUFG/BUFGCTRLs | 2 | 32 | 6% | |
| Number used as BUFGs | 2 | | | |
| Number of FIFO16/RAMB16s | 12 | 232 | 5% | |
| Number used as RAMB16s | 12 | | | |
| Average Fanout of Non-Clock Nets | 3.93 | | | |

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
Dialer
**Resource Consumption**

# Test project for an FPGA #0

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
Dialer
**Resource Consumption**

# Test project for an FPGA #0

| Project File: | final.xise | Parser Errors: | No Errors |
|---|---|---|---|
| Module Name: | topmodule_cn_fpga0_broad_sc | Implementation State: | Programming File Generated |
| Target Device: | xc4vfx60-10ff672 | • Errors: | No Errors |
| Product Version: | ISE 14.1 | • Warnings: | 15943 Warnings (3098 new) |
| Design Goal: | Timing Performance | • Routing Results: | All Signals Completely Routed |
| Design Strategy: | Performance with IOB Packing | • Timing Constraints: | All Constraints Met |
| Environment: | System Settings | • Final Timing Score: | 0 (Timing Report) |

| Device Utilization Summary | | | | [-] |
|---|---|---|---|---|
| **Logic Utilization** | **Used** | **Available** | **Utilization** | **Note(s)** |
| Number of Slice Flip Flops | 3,092 | 50,560 | 6% | |
| Number of 4 input LUTs | 10,717 | 50,560 | 21% | |
| Number of occupied Slices | 6,371 | 25,280 | 25% | |
| Number of Slices containing only related logic | 6,371 | 6,371 | 100% | |
| Number of Slices containing unrelated logic | 0 | 6,371 | 0% | |
| Total Number of 4 input LUTs | 10,999 | 50,560 | 21% | |
| Number used as logic | 10,713 | | | |
| Number used as a route-thru | 282 | | | |
| Number used as Shift registers | 4 | | | |
| Number of bonded IOBs | 35 | 352 | 9% | |
| IOB Flip Flops | 23 | | | |
| Number of BUFG/BUFGCTRLs | 1 | 32 | 3% | |
| Number used as BUFGs | 1 | | | |
| Number of FIFO16/RAMB16s | 6 | 232 | 2% | |
| Number used as RAMB16s | 6 | | | |
| Average Fanout of Non-Clock Nets | 4.24 | | | |

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
Dialer
**Resource Consumption**

# Test project for an algorithm FPGA

- Crossbar Switch with 10 interfaces
- 1 dummy sender
- 9 dummy receivers

Motivation
Online Data Reduction
Data Transport
**Implementation**
To Do

Overview
Crossbar Switch
Dialer
**Resource Consumption**

# Test project for an algorithm FPGA

| Project File: | final.xise | Parser Errors: | No Errors |
|---|---|---|---|
| Module Name: | topmodule_cn_fpga1234 | Implementation State: | Placed and Routed |
| Target Device: | xc4vfx60-10ff672 | • Errors: | No Errors |
| Product Version: | ISE 14.1 | • Warnings: | 21004 Warnings (0 new) |
| Design Goal: | Timing Performance | • Routing Results: | All Signals Completely Routed |
| Design Strategy: | Performance with IOB Packing | • Timing Constraints: | All Constraints Met |
| Environment: | System Settings | • Final Timing Score: | 0  (Timing Report) |

| Device Utilization Summary | | | | [-] |
|---|---|---|---|---|
| **Logic Utilization** | **Used** | **Available** | **Utilization** | **Note(s)** |
| Number of Slice Flip Flops | 2,209 | 50,560 | 4% | |
| Number of 4 input LUTs | 6,075 | 50,560 | 12% | |
| Number of occupied Slices | 4,050 | 25,280 | 16% | |
| Number of Slices containing only related logic | 4,050 | 4,050 | 100% | |
| Number of Slices containing unrelated logic | 0 | 4,050 | 0% | |
| Total Number of 4 input LUTs | 6,292 | 50,560 | 12% | |
| Number used as logic | 6,073 | | | |
| Number used as a route-thru | 217 | | | |
| Number used as Shift registers | 2 | | | |
| Number of bonded IOBs | 35 | 352 | 9% | |
| IOB Flip Flops | 7 | | | |
| Number of BUFG/BUFGCTRLs | 1 | 32 | 3% | |
| Number used as BUFGs | 1 | | | |
| Number of FIFO16/RAMB16s | 6 | 232 | 2% | |
| Number used as RAMB16s | 6 | | | |
| Average Fanout of Non-Clock Nets | 3.59 | | | |

Motivation
Online Data Reduction
Data Transport
Implementation
**To Do**

**RocketIO Instantiation**
Slow Control
Supervisor PC software
Thank you

## RocketIO Instantiation

- So far, only internal signals have been used for data transport
- Real life scenarios require the use of RocketIO/Aurora cores, which need to be instantiated

Motivation
Online Data Reduction
Data Transport
Implementation
To Do

**RocketIO Instantiation**
Slow Control
Supervisor PC software
Thank you

## RocketIO Instantiation

- So far, only internal signals have been used for data transport
- Real life scenarios require the use of RocketIO/Aurora cores, which need to be instantiated
- Possible issues: The high latency (some 100 clock cycles) will at least require tuning of the timeouts.

Motivation
Online Data Reduction
Data Transport
Implementation
To Do

**RocketIO Instantiation**
Slow Control
Supervisor PC software
Thank you

## RocketIO Instantiation

- So far, only internal signals have been used for data transport
- Real life scenarios require the use of RocketIO/Aurora cores, which need to be instantiated
- Possible issues: The high latency (some 100 clock cycles) will at least require tuning of the timeouts.

Motivation
Online Data Reduction
Data Transport
Implementation
**To Do**

RocketIO Instantiation
**Slow Control**
Supervisor PC software
Thank you

## Slow Control

- The supervisor PC needs to know the network of the FPGAs
- Therefore, it needs to contact each FPGA and ask it for its "inventory" (Crossbar Switches and Dialers)

Motivation
Online Data Reduction
Data Transport
Implementation
To Do

RocketIO Instantiation
**Slow Control**
Supervisor PC software
Thank you

## Slow Control

- The supervisor PC needs to know the network of the FPGAs
- Therefore, it needs to contact each FPGA and ask it for its "inventory" (Crossbar Switches and Dialers)
- Each Crossbar Switch must be asked for their respective neighbors

Motivation
Online Data Reduction
Data Transport
Implementation
**To Do**

RocketIO Instantiation
**Slow Control**
Supervisor PC software
Thank you

## Slow Control

- The supervisor PC needs to know the network of the FPGAs
- Therefore, it needs to contact each FPGA and ask it for its "inventory" (Crossbar Switches and Dialers)
- Each Crossbar Switch must be asked for their respective neighbors
- Requires some mechanism of communication between PC and FPGAs (Ethernet? IPMI?)

Motivation
Online Data Reduction
Data Transport
Implementation
**To Do**

RocketIO Instantiation
**Slow Control**
Supervisor PC software
Thank you

## Slow Control

- The supervisor PC needs to know the network of the FPGAs
- Therefore, it needs to contact each FPGA and ask it for its "inventory" (Crossbar Switches and Dialers)
- Each Crossbar Switch must be asked for their respective neighbors
- Requires some mechanism of communication between PC and FPGAs (Ethernet? IPMI?)

Motivation
Online Data Reduction
Data Transport
Implementation
**To Do**

RocketIO Instantiation
Slow Control
**Supervisor PC software**
Thank you

## Supervisor PC software

- Implementation of Dijkstra's algorithm
- Or any other suitable mechanism to determine routes (e.g. manual definition)

Motivation
Online Data Reduction
Data Transport
Implementation
**To Do**

RocketIO Instantiation
Slow Control
**Supervisor PC software**
Thank you

## Supervisor PC software

- Implementation of Dijkstra's algorithm
- Or any other suitable mechanism to determine routes (e.g. manual definition)

Motivation
Online Data Reduction
Data Transport
Implementation
To Do

RocketIO Instantiation
Slow Control
Supervisor PC software
Thank you

Thank you for your attention