

SpinNaker

**a universal
Spiking Neural Network
architecture**

version 0.0 - DRAFT
21 February 2006

SpiNNaker - a chip multiprocessor for neural network simulation

Features

- 20 ARM968 processors, each with:
 - 128 Kbytes of tightly-coupled data memory;
 - 64 Kbytes of tightly-coupled instruction memory;
 - DMA controller;
 - communications controller;
 - interrupt controller;
 - low-power 'wait for interrupt' mode.
- Multicast communications router
 - 6 serial inter-chip receive interfaces;
 - 6 serial inter-chip transmit interfaces;
 - 1024 associative routing entries.
- Interface to external SDRAM
 - over 1 Gbyte/s sustained block transfer rate.
- Fault-tolerant architecture
 - defect detection, isolation, and function migration.
- Boot, test and debug interfaces (to be determined).

Introduction

SpiNNaker is a chip multiprocessor designed specifically for the real-time simulation of large-scale spiking neural networks. Each chip (along with its associated SDRAM chip) forms one node in a scalable parallel system, interconnected to the other nodes through self-timed links.

The processing power is provided through the multiple ARM cores on each chip. Each ARM models multiple (up to 1,000) neurons, with each neuron being a coupled pair of differential equations modelled in continuous 'real' time. Neurons communicate through atomic 'spike' events, and these are communicated as discrete packets through the on- and inter-chip communications fabric. The packet contains a routing key that is defined at its source and is used to implement multicast routing through an associative router in each chip.

One processor on each SpiNNaker chip will perform system management functions; the communications fabric supports point-to-point packets to enable co-ordinated system management across local regions and across the entire system.

Background

SpiNNaker was designed at the University of Manchester within an EPSRC-funded project in collaboration with the University of Southampton, ARM Limited and Siiistix Limited. The work would not have been possible without EPSRC funding, and the support of the EPSRC and the industrial partners is gratefully acknowledged.

Intellectual Property rights

All rights to the SpiNNaker design are the property of the University of Manchester with the exception of those rights that accrue to the project partners in accordance with the contract terms.

Disclaimer

The details in this datasheet are presented in good faith but no liability can be accepted for errors or inaccuracies. The design of a complex chip multiprocessor is a research activity where there are many uncertainties to be faced, and there is no guarantee that an SpiNNaker system will perform in accordance with the specifications presented here.

The APT group in the School of Computer Science at the University of Manchester was responsible for all of the architectural and logic design of the SpiNNaker chip, with the exception of synthesizable components supplied by ARM Limited. All design verification was also carried out by the APT group. As such the industrial project partners bear no responsibility for the correct functioning of the device.

Change history

version	date	changes
0.0	27/12/05	first draft

Contents

1. Chip organization	5
1.1 Block diagram	5
1.2 System-on-Chip hierarchy	6
2. System architecture	7
2.1 Routing	7
2.2 System-level address spaces	8
3. ARM968 processing subsystem	9
3.1 Features	9
3.2 ARM968 subsystem organisation	9
3.3 Fault-tolerance	9
3.4 Test	9
4. ARM 968	10
4.1 Features	10
4.2 Organization	10
4.3 Fault-tolerance	10
4.4 Test	10
5. Counter/timer and interrupt controller	11
5.1 Features	11
5.2 Register summary	11
5.3 Register details	11
5.4 Fault-tolerance	11
5.5 Test	12
5.6 Notes	12
6. DMA controller	13
6.1 Features	13
6.2 Register summary	13
6.3 Register details	13
6.4 Fault-tolerance	13
6.5 Test	13
7. Communications controller	14
7.1 Features	14
7.2 Packet formats	14
7.3 Control byte summary	15
7.4 Register summary	16
7.5 Register details	16
7.6 Fault-tolerance	18
7.7 Test	18
7.8 Notes	18
8. Communications NoC	19
8.1 Features	19
8.2 Block diagram	19
8.3 Arbiter structure	19
8.4 Fault-tolerance	20
8.5 Test	20
9. Communications Router	21
9.1 Features	21

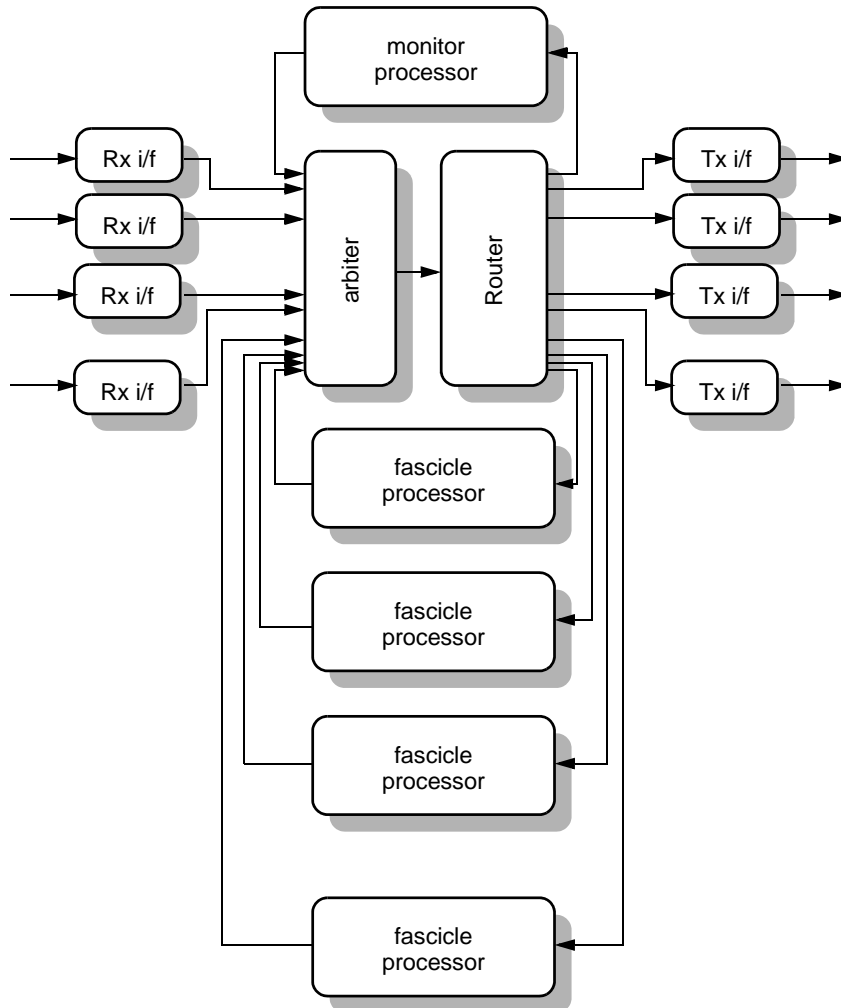
9.2 Description	21
9.3 Internal organization	21
9.4 Multicast (MC) router	22
9.5 The algorithmic (ALG) router	25
9.6 Packet error handler	25
9.7 Emergency routing	25
9.8 Errant packets	25
9.9 Fault-tolerance	25
9.10 Test	26
9.11 Notes	26
10. Inter-chip transmit and receive interfaces	27
10.1 Features	27
10.2 Programmer view	27
10.3 Fault-tolerance	27
10.4 Test	27
11. System NoC	29
11.1 Features	29
11.2 Fault-tolerance	29
11.3 Test	29
12. SDRAM interface	30
12.1 Features	30
12.2 Fault-tolerance	30
12.3 Test	30
13. System Controller	31
13.1 Features	31
13.2 Register summary	31
13.3 Register details	31
13.4 Fault-tolerance	31
13.5 Test	31
14. Router configuration registers	32
14.1 Features	32
14.2 Register summary	32
14.3 Register details	32
14.4 Fault-tolerance	32
14.5 Test	32
15. Boot ROM	33
15.1 Fault-tolerance	33
15.2 Test	33
16. System RAM	34
16.1 Features	34
16.2 Fault-tolerance	34
16.3 Test	34
17. Boot, test and debug support	35
17.1 Fault-tolerance	35
17.2 Test	35
18. Input and Output signals	36
18.1 External SDRAM interface	36
18.2	36

19. Area estimates	37
20. Power estimates	38

1. Chip organization

1.1 Block diagram

The primary functional components of SpiNNaker are illustrated in the figure below, which shows the Communications NoC and its clients.



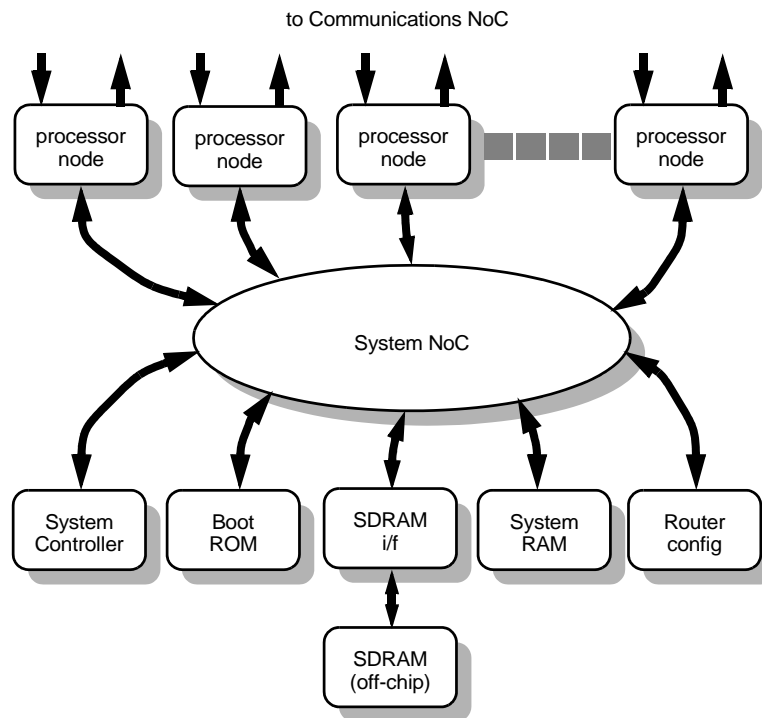
Each chip contains 20 identical processing subsystems each of which is responsible for modelling a number of neurons with associated inputs and outputs - a fascicle.

Following self-test, at start-up one of the processors is nominated as the Monitor Processor and thereafter performs system management tasks.

The router is responsible for routing neural event packets both between the on-chip fascicle processors and from and to other SpiNNaker chips. The Tx and Rx interface components are used to extend the on-chip communications NoC across to other SpiNNaker chips. The arbiter assembles inputs from the various on- and off-chip sources into a single serial stream which is then passed to the Router.

In addition to the primary function, there are additional resources accessible from the processor systems via the System NoC. Each of the fascicle processors has access to the shared off-chip SDRAM, and various system components also connect through the System NoC in order that, whichever processor is Monitor Processor, it will have access to these components.

The sharing of the SDRAM is an implementation convenience rather than a functional requirement, although it may facilitate function migration in support of fault-tolerant operation.



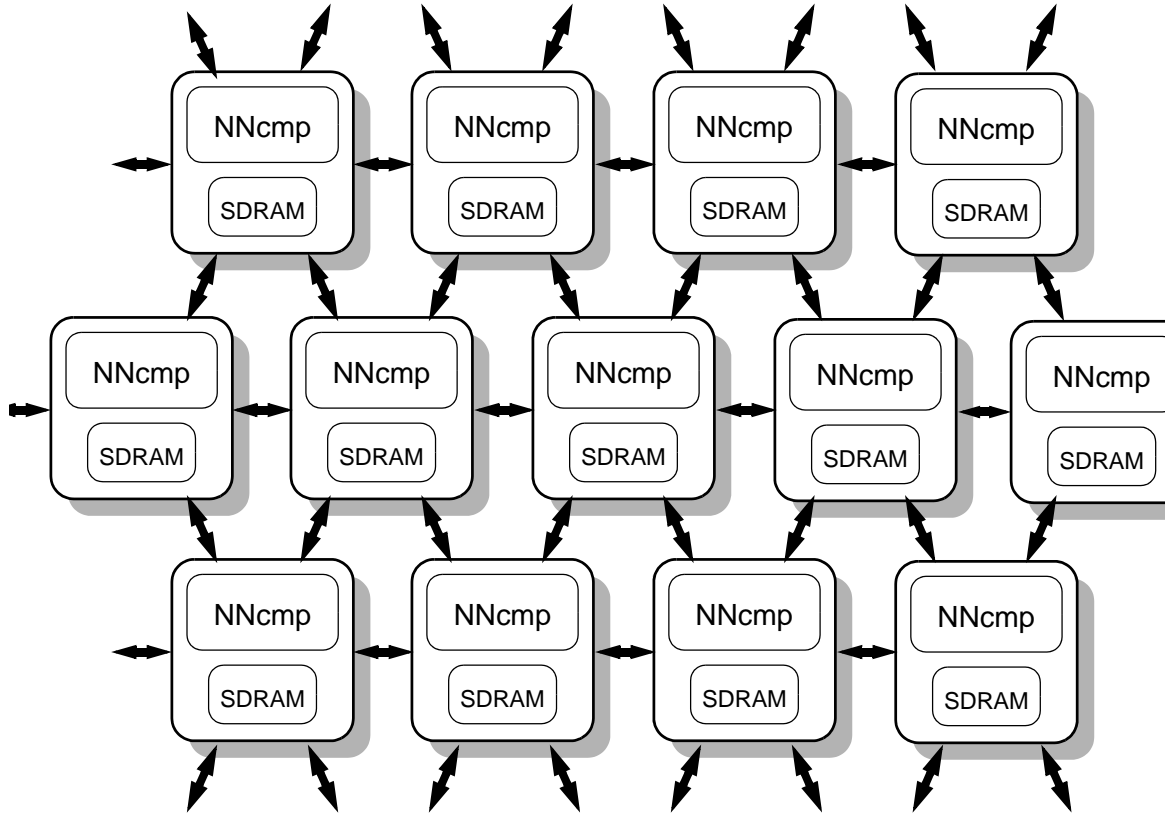
1.2 System-on-Chip hierarchy

The SpiNNaker chip is viewed as having the following structural hierarchy, which is reflected throughout the organisation of this datasheet:

- ARM968 processor subsystem
 - the ARM968, with its tightly-coupled instruction and data memories
 - Timer/counter and interrupt controller
 - DMA controller
 - communications controller, including communications NoC interface
 - System NoC interface
- Communications NoC
 - Router, including multicast, algorithmic, default and emergency routing functions
 - 6 inter-chip transmit interfaces
 - 6 inter-chip receive interfaces
 - communications NoC arbiter and fabric
- System NoC
 - SDRAM interface
 - System Controller
 - Router configuration registers
 - Boot ROM
 - System RAM
 - System NoC arbiter and fabric
- Boot, test and debug
 - central controller for ARM968 JTAG functions

2. System architecture

SpiNNaker is designed to form (with its associated SDRAM chip) a node of a massively parallel system. The system architecture is illustrated below:



2.1 Routing

The nodes are arranged in a *hexagonal* mesh with bidirectional links to 6 neighbours. The system supports both multicast packets (to carry neural event information, routed by the associative Multicast Router), and point-to-point packets (to carry system management and control information, routed algorithmically).

Emergency routing

In the event of a link failing or congesting, traffic that would normally use that link is redirected in hardware around two adjacent links that form a triangle with the failed link. This “emergency routing” is intended to be temporary, and the operating system will identify a more permanent resolution of the problem. The local Monitor Processor is informed of all uses of emergency routing.

Deadlock avoidance

The communications system has potential deadlock scenarios because of the possibility of circular dependencies between links. The policy used here to prevent deadlocks occurring is:

- *no Router can ever be prevented from issuing its output.*

The mechanisms used to ensure this are the following:

- outputs have sufficient buffering and capacity detection so that the Router knows whether or not an output has the capacity to accept a packet;
- emergency routing is used, where possible, to avoid overloading a blocked output;

- where emergency routing fails (because, for example, the alternative output is also blocked) the packet is 'dropped' to the local Monitor Processor;
- the local Monitor Processor is guaranteed to accept the dropped packet (eventually).

The expectation is that the communications fabric will be lightly-loaded so that blocked links are very rare. Where the operating system detects that this is not the case it will take measures to correct the problem by modifying routing tables or migrating functionality to a different part of the system.

Errant packet trap

Packets that get mis-routed could continue in the system for ever, following cyclic paths. To prevent this all packets are time stamped and a coarse global time phase signal is used to trap old packets. To minimize overhead the time stamp is 2 bits, cycling 00 -> 01 -> 11 -> 10, and when the packet is two time phases old (time sent XOR time now = 0b11) it is dropped to the local Monitor Processor and an error flagged. The length of a time phase can be adapted dynamically to the state of the system; normally timed-out packets should be very rare so the time phase can be conservatively long to minimise the risk of packets being dropped due to congestion.

2.2 System-level address spaces

The system incorporates a number of different levels of component that must be enumerated in some way:

- Each Node (where a Node is an SpiNNaker chip plus SDRAM) must have a unique, fixed address which is used as the destination ID for a point-to-point packet, and the addresses must be organised logically for algorithmic routing to function efficiently.
- Processors will be addressed relative to their host Node address, but this mapping will not be fixed as an individual Processor's role can change over time. Point-to-point packets addressed to a Node will be delivered to the local Monitor Processor, whichever Processor is serving that function. Internal to a Node there will be some hard-wired addressing of each Processor for system diagnosis purposes, but this mapping will be hidden outside the Node.
- Neurons occupy an address space that identifies each Neuron uniquely within the domain of its multicast routing path (where this domain must include alternative links that may be taken during emergency routing). Where these domains do not overlap it is possible to reuse the same address, though this must be done with considerable care. Neuron addresses can be assigned arbitrarily, and this flexibility can be exploited to optimize Router utilization (for example by giving Neurons with the same routing requirements related addresses so that they can all be routed by the same Router entries).

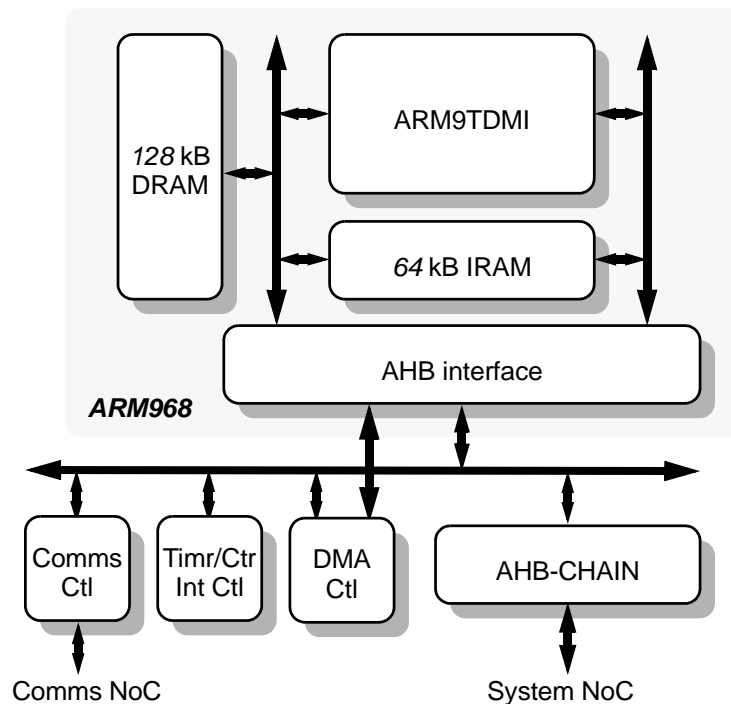
3. ARM968 processing subsystem

SpiNNaker incorporates 20 ARM968 processing subsystems which provide the computational capability of the device. Each of these subsystems is capable of generating and processing neural events communicated via the Communications NoC and, alternatively, of fulfilling the role of Monitor Processor.

3.1 Features

- a synthesized ARM968 module with
 - a 200 MIPS ARM9 processor
 - 64 kB tightly-coupled instruction memory
 - 128 kB tightly-coupled data memory
- a local AHB with
 - communications controller connected to Communications NoC
 - wrapper to interface to the System NoC
 - DMA controller
 - timer/counter and interrupt controller

3.2 ARM968 subsystem organisation



3.3 Fault-tolerance

The fault-tolerance of the ARM968 subsystem is defined in terms of its component parts, described below.

3.4 Test

The test strategies for the ARM968 subsystem are likewise defined in terms of its component parts.

4. ARM 968

The ARM968 (with its associated tightly-coupled instruction and data memories) forms the core processing resource in SpiNNaker. It is a standard synthesizable IP component from ARM Ltd, and as such there is limited scope for customizing it for this application.

4.1 Features

- 200 MIPS ARM9TDMI processor.
- 64 kB tightly-coupled instruction memory (I-RAM).
- 128 kB tightly-coupled data memory (D-RAM).
- AHB interface to external system.

4.2 Organization

See ARM DDI 0311C – the ARM968E-S datasheet.

4.3 Fault-tolerance

Fault insertion

- ARM9TDMI can be disabled.
- Software can corrupt I-RAM and D-RAM to model soft errors. (Can these be detected?)

Fault detection

- The I-RAM and D-RAM are protected by parity bits?
- A watchdog timer can catch runaway software.
- Self-test routines, run at start-up and during normal operation, can detect faults.

Fault isolation

- The ARM968 unit can be disabled from the System Controller.
- defective locations in the I-RAM and D-RAM can be mapped out of use by software.

Reconfiguration

- software will avoid using defective I-RAM and D-RAM locations.
- functionality will migrate to an alternative Processor in the case of permanent faults that go beyond the failure of one or two memory locations.

4.4 Test

production test

start-up test

run-time test

5. Counter/timer and interrupt controller

Each processor node on an SpiNNaker chip has a local counter/timer and interrupt controller that is used to enable and disable interrupts from various sources, and to wake the processor from sleep mode when required. The controller provides centralised management of IRQ and FIQ sources, and offers an efficient indication of the active sources for vectoring purposes.

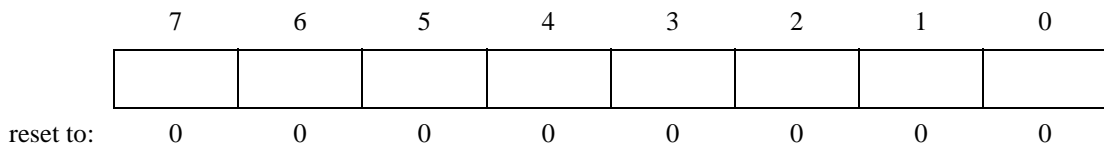
5.1 Features

- manages the various interrupt sources to each local processor:
 - arriving multicast packet without payload
 - arriving multicast packet with payload
 - arriving point-to-point packet without payload
 - arriving point-to-point packet with payload
 - DMA complete
 - timer interrupt
 - interrupt from another processor on the chip
 - system fault interrupt
- the counter/timer unit provides two independent counters, for example for:
 - millisecond interrupts for real-time dynamics
 - watchdog timer for fault trapping

5.2 Register summary

Name	Offset	R/W	Function

5.3 Register details



5.4 Fault-tolerance

Fault insertion

Fault detection

Fault isolation

Reconfiguration

5.5 Test

production test

start-up test

run-time test

5.6 Notes

- millisecond interrupt could be provided from a centralised C/T unit? But will all processors want to receive the same time interrupts?
- extra interrupt(s) for packet parity failure?

6. DMA controller

Each ARM968 processing subsystem includes a DMA controller. The DMA controller is primarily used for transferring inter-neural connection data from the SDRAM in large blocks in response to an input event arriving at a fascicle processor, and for returning updated connection data during learning.

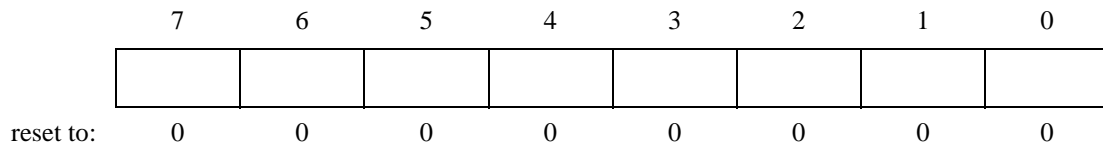
6.1 Features

-

6.2 Register summary

Name	Offset	R/W	Function

6.3 Register details



6.4 Fault-tolerance

Fault insertion

Fault detection

Fault isolation

Reconfiguration

6.5 Test

production test

start-up test

run-time test

7. Communications controller

Each processor node on SpiNNaker includes a communications controller which is responsible for generating and receiving packets to and from the communications network.

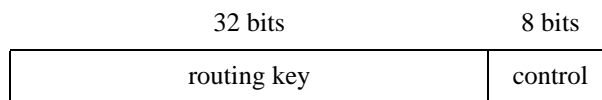
7.1 Features

- Support for 2 packet types:
 - 40-bit multicast neural event packets routed by a key provided at the source;
 - 40-bit point-to-point packets routed algorithmically by destination address.
- Packets may optionally carry an additional 32-bit payload.
- 2-bit time stamp (used by Routers to trap errant packets).
- Parity (to detect corrupt packets).

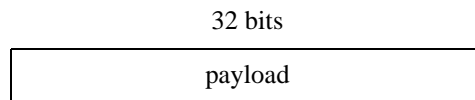
7.2 Packet formats

Neural event multicast (mc) packets (type 1)

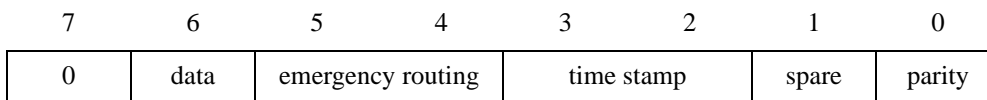
Neural event packets include a 32-bit routing key inserted by the source and a control byte:



In addition they may include an optional (not normally used) 32-bit payload:

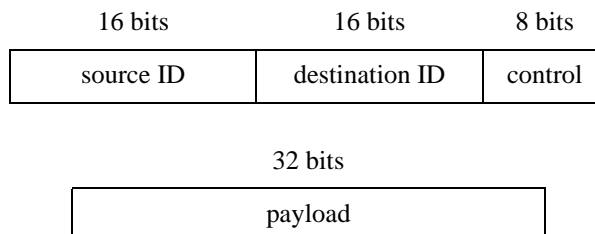


The 8-bit control field includes packet type (= 0 for multicast packets), a data payload indicator, emergency routing, time stamp and error detection (parity) information, plus a spare bit for software use:



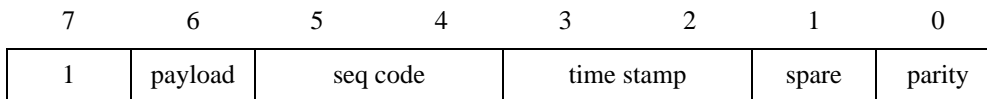
Point-to-point (p2p) packets (type 0)

Point-to-point packets include 16-bit source and destination chip IDs, plus a control byte and an optional (normally used) 32-bit payload:



Here the 8-bit control field includes packet type (=1 for p2p packets), a data payload indicator, a sequence code, time stamp and error detection (parity) information, plus a spare bit for software

use:



7.3 Control byte summary

Field Name	bits	Function
parity	0	parity of complete packet (including payload when used)
spare	1	available for software use
time stamp	3:2	phase marker indicating time packet was launched
seq code	5:4	p2p only: start, middle odd/even, end of payload
emergency routing	5:4	mc only: used to control routing around a failed link
data	6	data payload attached
packet type	7	= 0 for mc; = 1 for p2p

parity

The complete packet (including the data payload where used) will have odd parity.

spare

This bit is set at packet launch to the value defined in the control register, and the receiver(s) may inspect it. It may be used for any purpose by the software.

time stamp

The system has a global time phase that cycles through 00 -> 01 -> 11 -> 10 -> 00. Global synchronisation must be accurate to within less than one time phase (the duration of which is programmable and may be dynamically variable). A packet is launched with a time stamp equal to the current time phase, and if a Router finds a packet that is two time phases old (time now XOR time launched = 11) it will drop it to the local Monitor Processor.

seq code

p2p packets use these bits to indicate the sequence of data payloads:

- 11 -> start packet: the first packet in a sequence (of >1 packets)
- 10 -> middle even: the second, fourth, sixth, ... packet in a sequence
- 01 -> middle odd: the third, fifth, seventh, ... packet in a sequence
- 00 -> end: the last (or only) packet in a sequence

emergency routing

mc packets use these bits to control emergency routing around a failed or congested link:

- 00 -> normal mc packet;
- 01 -> the packet has been redirected by the previous Router through an emergency route along with a normal copy of the packet. The receiving Router should treat this as a combined normal plus emergency packet.

- 10 -> the packet has been redirected by the previous Router through an emergency route which would not be used for a normal packet.
- 11 -> this emergency packet is reverting to its normal route.

data

- indicates whether the packet has a 32-bit data payload (=1) or not (=0).

packet type

- this bit indicates whether the packet is a multicast (0) or point-to-point (1) packet.

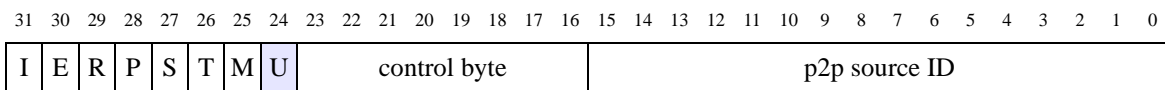
7.4 Register summary

Name	Offset	R/W	Function
r0: status/config	0x0	R/W	Indicates the current status of the controller
r1: unused	0x4	-	-
r2: send data	0x8	W	32-bit data for transmission
r3: send key	0xC	W	Send mc key/p2p dest ID & seq code
r4: receive data	0x10	R	32-bit received data
r5: receive key	0x14	R	Received mc key/p2p source ID & seq code

A packet will contain data if r4 is written before r5; this can be performed using an ARM STM instruction.

7.5 Register details

r0: status & control



The functions of these fields are described in the table below:

Name	bits	R/W	Function
p2p source ID	15:0	W	16-bit chip source ID for p2p packets
control byte	23:16	R	control byte of last received packet
U: unused	24	-	-
M: packet type	25	W	send multicast (=1) or point-to-point (=0) packet
T: transmit full	26	R	transmit buffer full
S: spare	27	W	outgoing packet 'spare' bit value
P: parity	28	R	received packet parity error (=1)

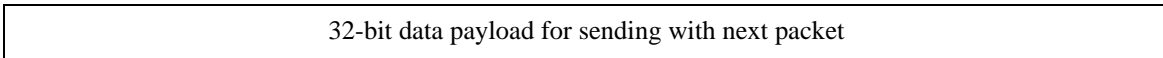
Name	bits	R/W	Function
R: received	29	R	packet received
E: int. enable	30	W	enable interrupt by received packet
I: int. status	31	R	enabled interrupt caused by received packet

The p2p source ID is expected to be configured once at start-up.

r1: unused

r2: send data

3 3 2 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 9 8 7 6 5 4 3 2 1 0
1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0



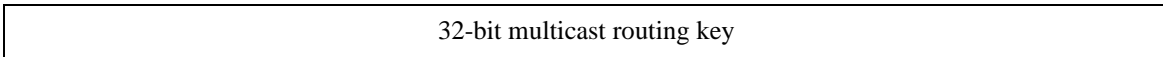
If data is written into r2 before a send key or dest ID is written into r3, the packet initiated by writing to r3 will include the contents of r2 as its data payload. If no data is written into r2 before a send key or dest ID is written into r3 the packet will carry no data payload.

r3: send mc key/p2p dest ID & sequence code

Writing to r3 will cause a packet to be issued (with a data payload if r2 was written previously).

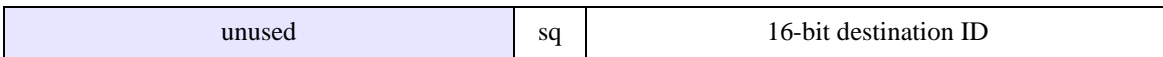
If bit[25] of the control register is set the Communication Controller is set to send multicast packets and a 32-bit routing key should be written into r3. The 32-bit routing key is used by the associative multicast Routers to deliver the packet to the appropriate destination(s).

3 3 2 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 9 8 7 6 5 4 3 2 1 0
1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0



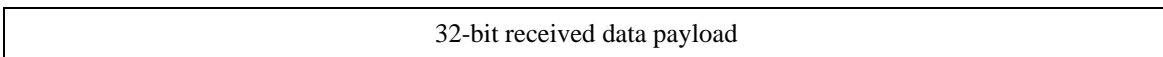
If bit[25] of the control register is clear the Communication Controller is set to send point-to-point packets and the value written into r3 should include the 16-bit address of the destination chip in bits[15:0] and a sequence code in bits[17:16]. (See 'seq code' on page 15.)

3 3 2 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 9 8 7 6 5 4 3 2 1 0
1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0



r4: received data

3 3 2 2 2 2 2 2 2 2 2 2 1 1 1 1 1 1 1 1 1 1 9 8 7 6 5 4 3 2 1 0
1 0 9 8 7 6 5 4 3 2 1 0 9 8 7 6 5 4 3 2 1 0



If a received packet carries a data payload the payload will be delivered here and will remain valid until r5 is read.

r5: received mc key/p2p source ID & sequence code

A received packet will deliver its routing key or source ID and sequence code to r5. This will be the exact value that the sender placed into its r3 for transmission. The register is read sensitive - once

read it will change as soon as the next packet arrives.

7.6 Fault-tolerance

Fault insertion

Software can cause the Communications Controller to misbehave in several ways including inserting dodgy routing keys, source IDs, destination IDs.

Do we need to be able to force parity errors in transmit packets?

Fault detection

Parity of received packet?

Fault isolation

The Communications Controller is mission-critical to the local processing subsystem, so if it fails the subsystem should be disabled and isolated.

Reconfiguration

The local processing subsystem is shut down and its functions migrated to another subsystem on this or another chip. It should be possible to recover all of the subsystem state and to migrate it, via the SDRAM, to a functional alternative.

7.7 Test

production test

start-up test

run-time test

7.8 Notes

- time phase accuracy: if we assume that the system time phase is F and the skew is K (that is, all parts of the system transition from one phase to its successor within a time K), then a packet has at least $F-K$ to reach its destination and will be killed after at most $2F+K$. Thus, if we want to allow for a maximum packet transit time of $F-K = T$ and can achieve a minimum phase skew of K , then T and K are both system constants and we should choose $F = T+K$. The longest packet life is then $2T+3K$.

8. Communications NoC

The communications NoC has the primary role of carrying neural event packets between Fascicle Processors on the same or different chips.

8.1 Features

- On- and inter-chip links
- Router with associative multicast, algorithmic, default and emergency routing functions.
- Arbiter to merge all sources into a sequential packet stream into the Router.
- Individual links can be reset to clear blockages and deadlocks.

8.2 Block diagram

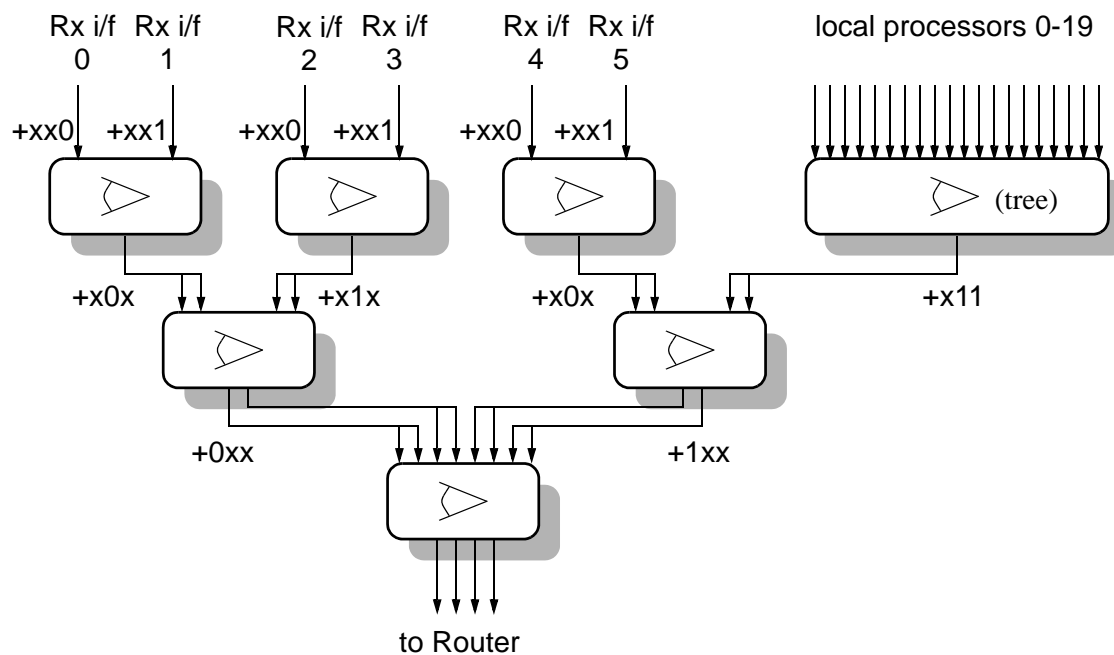
A block diagram of the Communications NoC was given in section 1.1 on page 5.

8.3 Arbiter structure

As the input links converge on the Router they must merge through 2-way CHAIN arbiters, and the number of parallel links must increase to absorb the bandwidth. The following hierarchy is proposed:

- the local processor links can all be merged through a single-link arbiter tree as the local bandwidth is low, e.g. at most 20 processors x 1,000 neurons x 100Hz x 40 bits = 80 Mbit/s.
- the Rx interfaces can each carry up to 1 Gbit/s, about half the on-chip single-link bandwidth, so the first layer of arbiters can be single-link, the 2nd layer dual-link and the 3rd layer quad-link (i.e. 8-bits or 48 wires wide).
- at each arbiter merging Rx interfaces the packet must pick up 1 bit to indicate its source, for default routing.

The Arbiter structure is illustrated below. Each doubling of the wires represents a doubling of the number of CHAIN links. The numbers indicate source tagging of the packets.



8.4 Fault-tolerance

Fault insertion

There is little direct control of the Communications NoC fabric except at the periphery as noted in the sections below.

Fault detection

Most failures will cause local asynchronous deadlock, which is readily detected at both the transmitting and receiving ends of the link.

Fault isolation

If links fail their clients will have to be disabled and their functions migrated.

Reconfiguration

Client functional migration is required.

8.5 Test

production test

start-up test

run-time test

9. Communications Router

The Communications Router is responsible for routing all packets that arrive at its input to one or more of its outputs. Its primary function is to route multicast neural event packets, which it does through an associative multicast router subsystem. But it is also responsible for routing point-to-point packets (for which it uses algorithmic routing), for default routing (when a multicast packet does not match any entry in the multicast router) and for emergency routing (when an output link is blocked due to congestion or hardware failure).

Various error conditions are identified and handled by the Communications Router, for example packet parity errors, time-out, and output link failure.

9.1 Features

- 1024 programmable associative multicast routing entries.
 - associative routing based on source ‘key’.
 - with flexible “don’t care” masking.
- algorithmic routing of point-to-point packets.
- support for 40- and 72-bit multicast and point-to-point packets.
- default routing of unmatched multicast packets.
- automatic re-routing around failed links.
- Failure detection and handling:
 - packet parity error
 - time-expired packet
 - output link failure

9.2 Description

We assume that messages arrive from other nodes via the link receiver interfaces and from internal clients and are presented to the router one-at-a-time. The Arbiter is responsible for determining the order of presentation of the messages, but as each message is handled independently the order is unimportant (though it is desirable for packets following the same route to stay in order).

Each message contains an identifier that is used by the Communications Router to determine which of the outputs the message is sent to. These outputs may include any subset of the output links, where the message may be sent via the respective link transmitter interface, and/or any subset of the internal processor nodes, where the message is sent to the respective Communications Controller.

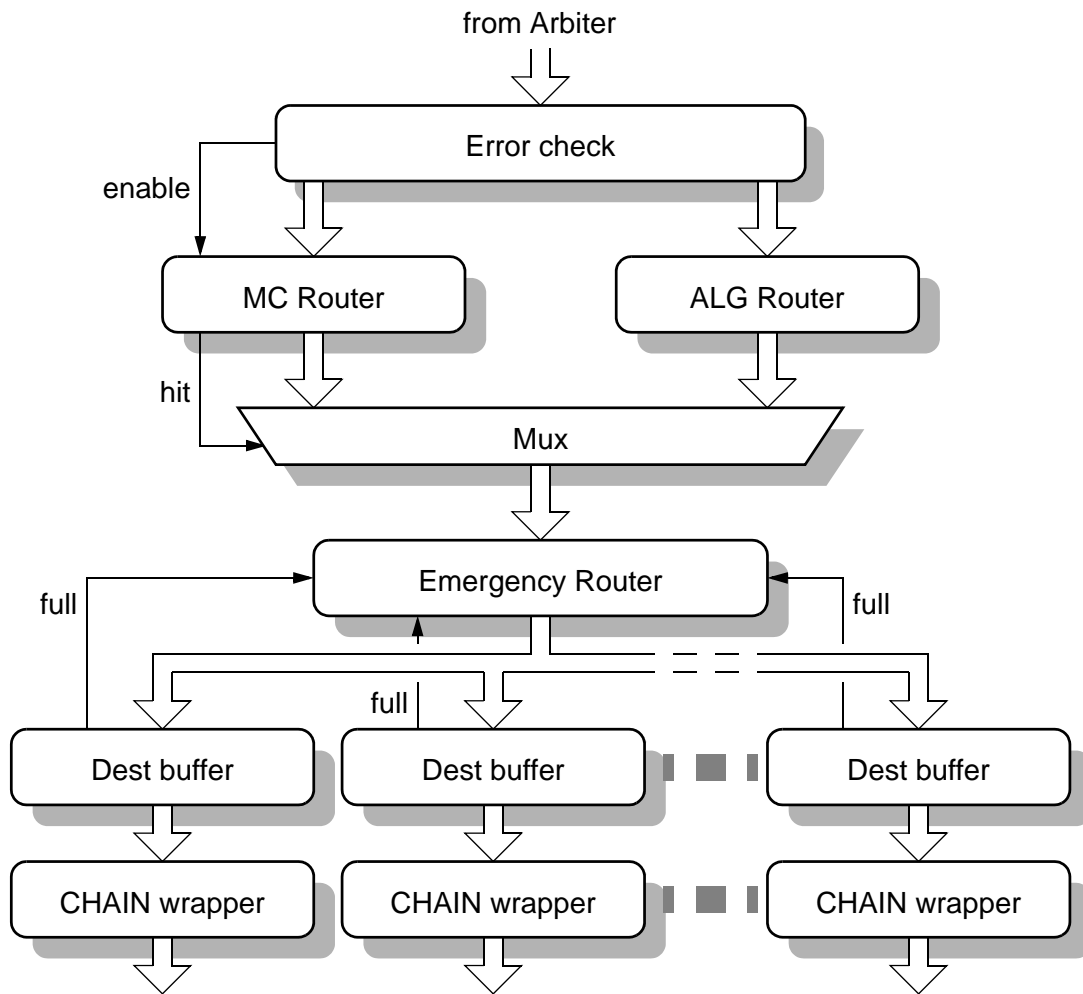
For the neural network application the identifier can be simply a number that uniquely identifies the source of the message – the neuron that generated the message by firing. This is ‘source address routing’. In this case the message need contain only this identifier, as a neural spike is an “event” where the only information is that the neuron has fired.

The Router then functions simply as a look-up table where for each identifier it looks up a routing word, where each routing word contains 1 bit for each destination (each link transmitter interface and each local processor) to indicate whether or not the message should be passed to that destination.

9.3 Internal organization

The internal organization of the Communications Router is illustrated in the figure opposite.

Packets are passed as complete 40- or 72-bit units from the Arbiter, together with an identifier of the Rx interface that the packet arrived through (for default routing). The first stage of processing



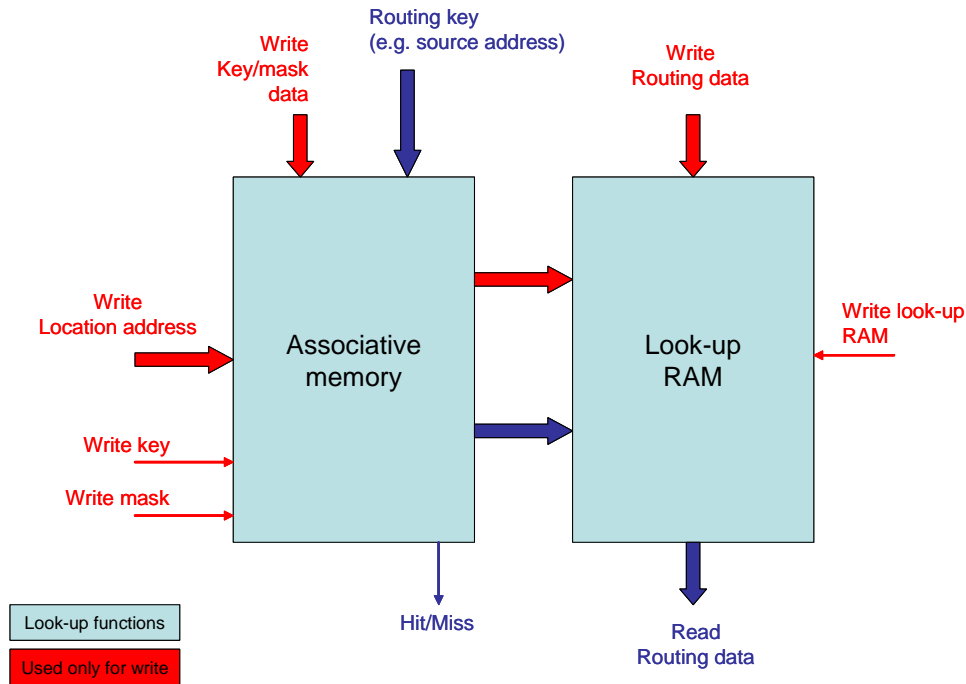
here is to identify errors. The second stage passes the packet to the appropriate routing engines – the multicast (MC) router is activated only if the packet is error-free and of multicast type, but the algorithmic (ALG) router is activated for every packet as it handles default routing in addition to point-to-point algorithmic and error routing. The output of the router stage is a vector of destinations to which the packet should be relayed. The third stage is the emergency routing mechanism for handling failed or congested links, which it detects using ‘full’ signals fed back from the individual destination output buffers.

Notes

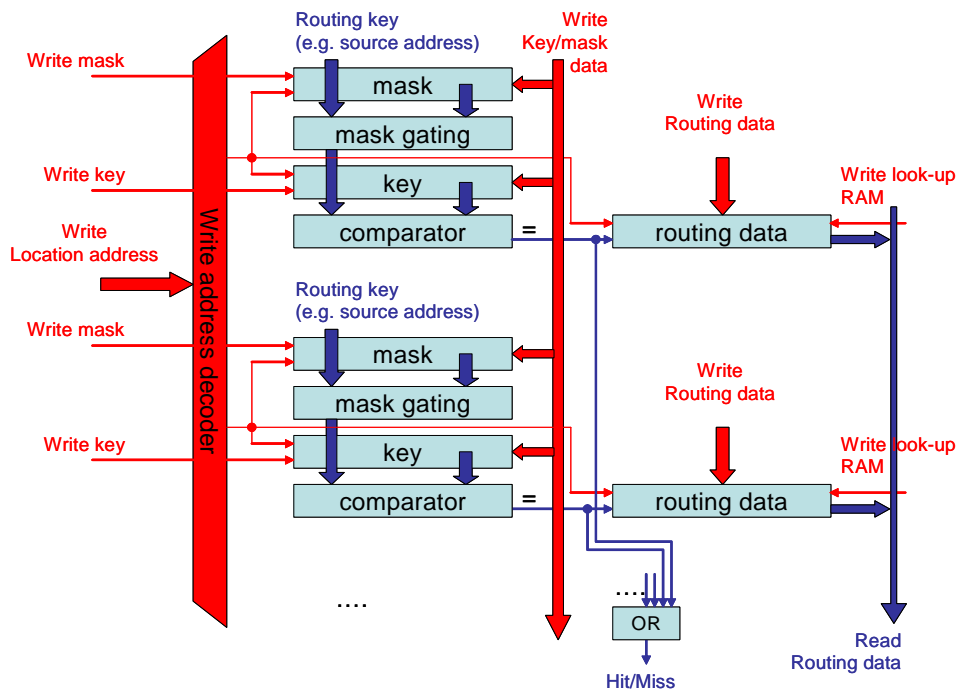
- the Router needs to know which is the Monitor Processor for routing terminating p2p, error, and dropped packets
- how are details of errors communicated to the Monitor Processor?
- Emergency routing may cause two packets to be issued onto the same output link (one for the MC routed data and the second for the alternative route for the blocked link). These are merged into a single packet with a different emergency-routing type to indicate its dual purpose.

9.4 Multicast (MC) router

The internal organisation of the multicast router is illustrated in the figure below.



Implementation



Multicast router optimisations

The simple look-up table as described above works in principle but in practice would be too large to fit onto a chip. However, several optimisations address this problem, reducing the table to a practical size:

- The table within a particular node need contain entries only for message identifiers whose routes from their source to all of their destinations pass through, are generated in, or end in that node.

- Default routing can be supported that passes messages from a link receiver interface to the diametrically opposite link transmitter interface when the message identifier is not found in the look-up table.
- Groups of message identifiers can be routed using the same look-up table entry by making some of the identifier bits “don’t care” as far as the look-up process is concerned.

The logical structure of the Router is then an associative (content-addressed) memory, with programmable masking on a per-entry basis to support the “don’t care” optimization, connected to a conventional memory that holds the per-entry output routing word. The associative memory can be implemented using any of the usual techniques such as VLSI CAM cells or hash-addressed RAM.

Additional mechanisms can support the default routing mentioned above when there is no match for the message identifier in the associative memory and, through partitioning the identifier address space, provision can be made for conventional 1-to-1 destination address routing and broadcast mechanisms.

Illustrative example

By way of illustration, let us assume a neural modelling system where each neuron has a unique 14-bit number, and when it fires it transmits this number prepended by two zero bits as the message identifier (and the message has no other content). Further, we assume that the neurons are handled in groups of 256, where each group of 256 is assigned to a particular processor on a particular node and each neuron in a group has the same set of inputs and the same output destinations as every other neuron in the same group.

The message identifiers can then be processed by the Routers as the 16-bit binary number: 00nnnnnnXXXXXXXXXX, where:

- 00 indicates that this is a source address message identifier that should be routed according to the routing table;
- nnnnnn is the neuron group identifier, so at most $26 = 64$ routing look-up table entries are required;
- XXXXXXXX indicates that the bottom 8 bits of the message identifier can be treated as “don’t care” and play no role in the routing. They will be used only by the destination processor to identify the neuron that fired within the group.

Messages beginning with 01, 10, and 11 may be used for 1-to-1 destination address routing, broadcast, and some other purpose respectively, using conventional routing algorithms.

It can be seen from this example that the optimisations are a vital aspect of the invention, reducing the size of the look-up table in this very small example from $16,384 (= 2^{14})$ to at most 64 entries. For larger systems the benefits of the optimisations are likely to be even more significant.

The “don’t care” bits are programmable independently for each look-up table entry on each node, and they can be distributed anywhere across the message identifier for maximum flexibility, so a single look-up table entry may route several groups together at one node, whereas at the next node they may be routed independently to different destinations.

Route set up

The routing look-up tables may be configured using external software that takes a neural “netlist”, describing the way the neurons interconnect, and then maps the neurons onto processors and determines the routing table values using algorithms similar to those used to configure an FPGA. As with FPGA configuration, resource constraints such as the routing table size and link bandwidth limitations must be taken into account during the mapping process.

The routing table configuration is then loaded into the local Router by a local processor that follows instructions from a control system using the 1-to-1 message routing mechanism.

For static neural network modelling the routing is fixed after initialisation. It is possible to allow local processors to modify the routing tables while the system is running, if this is required to model

developmental processes (for example), provided this is done with due care.

9.5 The algorithmic (ALG) router

The algorithmic router uses the 16-bit destination ID in a point-to-point packet to determine which (single) output the packet should be routed to. A 64K entry SRAM lookup table directs the p2p packet to:

- the local Monitor Processor, or
- an adjacent chip via the appropriate link.

In addition, the algorithmic router performs default and error routing functions.

9.6 Packet error handler

The packet error handler is a routing engine that simply flags the packet for routing to the local Monitor Processor if it detects any of the following:

- a packet parity error;
- a packet that is two time phases old.

9.7 Emergency routing

If a link fails (temporarily, due to congestion, or permanently, due to component failure) action will be taken at two levels:

- the blocked link will be detected in hardware and subsequent packets rerouted via the other two sides of one of the routing triangles of which the suspect link was an edge.
- the monitor processor will be informed. It will assess the problem, and take appropriate action:
 - if the problem was due to transient congestion, it will note the congestion but do nothing further;
 - if the problem was due to recurring congestion, it will negotiate and establish a new route for some of the traffic using this link;
 - if the problem appears permanent, it will reset the link (incurring some packet loss) and then, if this does not clear the problem, negotiate and establish new routes for all of the traffic using this link.

The hardware support for these processes include:

- default routing processes in adjacent nodes that are invoked by flagging the packet as an emergency type;
- mechanisms to inform the monitor processor of the problem;
- mechanisms the monitor processor can use to reset the link;
- means of inducing the various types of fault for testing purposes.

Emergency rerouting around the triangle requires additional emergency packet types for mc packets. p2p packets will find their own way to their destination following emergency routing.

9.8 Errant packets

In order to ensure that packets cannot circulate for ever within the system each packet includes a time phase field. This is set when the packet is launched, and if a packet arrives at a Router two time phases after it was launched it will be routed directly (and only) to the local monitor processor for error-handling purposes.

9.9 Fault-tolerance

The Communications Router has limited fault-tolerance capacity, mainly coming down to mapping

out a failed multicast router entry. This is a useful mechanism as the multicast router dominates the silicon area of the Communications Router.

Fault insertion

- enable Router to flip packet parity bits?

Fault detection

- packet parity errors
- packet time-phase errors
- packet unroutable errors (e.g. a locally-sourced multicast packet which doesn't match any entry in the multicast router).

Fault isolation

- a mechanism is required to disable a multicast router entry if it fails. Possible just an 'entry valid' bit?

Reconfiguration

- since all multicast router entries are identical the function of any entry can be relocated to a spare entry (within the same segment of the router if segmentation is used to save power).
- if a router (segment) becomes full a global reallocation of resources can move functionality to a different router (segment)

9.10 Test

production test

start-up test

run-time test

9.11 Notes

- The Router will require a number of traffic monitor features, e.g. packet counters, congestion indicators, count packet under match & mask, dropped packet count, emergency routing count, count on each output link, ...

10. Inter-chip transmit and receive interfaces

Inter-chip communication is implemented by extending CHAIN links from chip to chip. In order to sustain CHAIN link throughput, there is a protocol conversion at each chip boundary from standard CHAIN 1-of-5 (including EOP) return-to-zero to 2-of-7 non-return-to-zero. Each conversion maps two 2-bit CHAIN symbols to a single 4-bit 2-of-7 symbol.

10.1 Features

- transmit (Tx) interface:
 - converts two on-chip 1-of-5 RTZ symbols into one off-chip 2-of-7 NRZ symbol;
 - control input to induce a fault;
 - failure detection output.
 - fault reset input.
- receive (Rx) interface:
 - converts one off-chip 2-of-7 NRZ symbol into two on-chip 1-of-5 RTZ symbols;
 - control input to induce a fault;
 - failure detection output.
 - fault reset input.

10.2 Programmer view

There are no programmer-accessible features implemented in these interfaces. In normal operation these interfaces provide transparent connectivity between the routing network on one chip and those on its neighbours.

10.3 Fault-tolerance

The fault inducing, detecting and resetting functions are controlled from the System Controller (see 'System Controller' on page 31).

Fault insertion

- an input controlled by the System Controller causes the interface to deadlock

Fault detection

- an output to the System Controller indicates deadlock

Fault isolation

- the interface can be disabled to isolate the chip-to-chip link. This may be the same input from the System Controller that is used to insert a fault.

Reconfiguration

- the link interface can be reset by the System Controller to attempt recovery from a fault
- the link interface can be isolated and an alternative route used

10.4 Test

production test

start-up test

run-time test

S
p
i
n
n
a
k
e
r

11. System NoC

The System NoC has a primary function of connecting the Fascicle Processors to the SDRAM interface. It is also used to connect the Monitor Processor to system control and test functions, and for a variety of other purposes.

11.1 Features

- supports full bandwidth block transfers between the SDRAM and the Fascicle Processors.
- can be reset (in subsections?) to clear deadlocks.

11.2 Fault-tolerance

Fault insertion

Fault detection

Fault isolation

Reconfiguration

11.3 Test

production test

start-up test

run-time test

12. SDRAM interface

The SDRAM interface connects the System NoC to an off-chip SDRAM device. It will be the ARM xxx.

12.1 Features

- lots of bandwidth, please!

12.2 Fault-tolerance

Fault insertion

Fault detection

Fault isolation

Reconfiguration

12.3 Test

production test

start-up test

run-time test

13. System Controller

The System Controller incorporates a number of functions used by the Monitor Processor for system start-up, fault-tolerance testing (invoking, detecting and resetting faults), general performance monitoring, and such like.

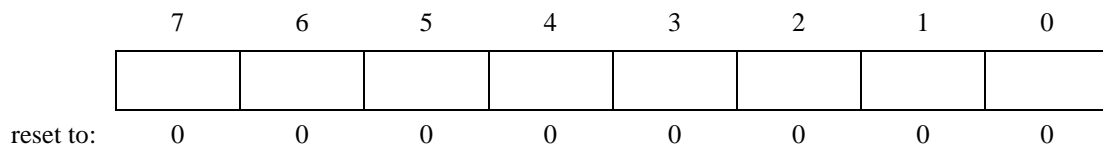
13.1 Features

- lots of fiddly bits
- input & output ports to connect to the communications NoC Tx & Rx interface fault invocation, detection and reset functions.

13.2 Register summary

Name	Offset	R/W	Function

13.3 Register details



13.4 Fault-tolerance

Fault insertion

Fault detection

Fault isolation

Reconfiguration

13.5 Test

production test

start-up test

run-time test

14. Router configuration registers

The Router is highly configurable, and the Monitor Processor is responsible for initialising it and updating it when necessary. The Router configuration registers are accessed via the system NoC.

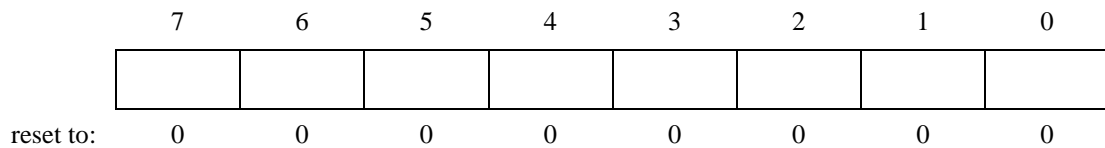
14.1 Features

- used to set up the associative routing tables.
- give read/write access to the Router tables for test purposes.

14.2 Register summary

Name	Offset	R/W	Function

14.3 Register details



14.4 Fault-tolerance

Fault insertion

Fault detection

Fault isolation

Reconfiguration

14.5 Test

production test

start-up test

run-time test

15. Boot ROM

15.1 Fault-tolerance

Fault insertion

Fault detection

Fault isolation

Reconfiguration

15.2 Test

production test

start-up test

run-time test

16. System RAM

The System RAM is an additional 256 kByte block of on-chip RAM used primarily by the Monitor Processor to enhance its program and data memory resources as it will be running more complex (though less time-critical) algorithms than the Fascicle Processors.

As the choice of Monitor Processor is made at start-up (and may change during run-time for fault-tolerance purposes) the System RAM is made available to whichever processor is Monitor Processor via the System NoC. It is probably important that accesses by the Monitor Processor to the System RAM are non-blocking as far as SDRAM accesses by the Fascicle Processors are concerned, so the System NoC should ensure this is the case.

The System RAM may also be used by the Fascicle Processors to communicate with the Monitor Processor and with each other, should the need arise.

16.1 Features

- 256 kB of SRAM, available via the System NoC.
- can be disabled to model complete failure for fault-tolerance testing.
- can we include parity or ECC to improve fault-tolerance?

16.2 Fault-tolerance

Fault insertion

- It is straightforward to corrupt the contents of the System RAM to model a soft error – any processor can do this. It is not clear how this would be detected.
- The System RAM can be disabled to model a total failure.

Fault detection

- The Monitor Processor may perform a System RAM test at start-up, and periodically thereafter.
- It is not clear how soft errors can be detected without some sort of parity or ECC system.

Fault isolation

- Faulty words in the System SRAM can be mapped out of use.

Reconfiguration

- For hard failure of a single bit, avoid using the word containing the failed bit.
- If the System RAM fails completely the only option is to use the SDRAM instead, which will probably result in compromised performance for the Fascicle Processors due to loss of SDRAM bandwidth. An option then would be to relocate some of the Fascicle Processors' workload to another chip.

16.3 Test

production test

- run standard memory test patterns from one of the processing subsystems.

start-up test

run-time test

17. Boot, test and debug support

17.1 Fault-tolerance

Fault insertion

Fault detection

Fault isolation

Reconfiguration

17.2 Test

production test

start-up test

run-time test

18. Input and Output signals

18.1 External SDRAM interface

Signal	Type	Function

18.2

19. Area estimates

We are targeting a UMC 130nm process 10mm x 10mm die. (Europractice runs on this process are multiples of 5mm x 5mm. The test chips will be 5mm x 5mm.)

Assumptions

- RAM is around $1 \mu\text{m}^2/\text{bit}$ = 6M T/mm².
- logic is 0.1 x the density of RAM = 100k gates/mm².
- The pad ring occupies 0.25 mm all round the chip, so the core is 9.5 x 9.5 = 90.25 mm².

Using these assumptions we total up the core logic area as follows:

- The processor nodes = 20 x 3.8 = 76 mm².
 - An ARM968 with 64kByte I-RAM and 128kByte D-RAM is 3.5 mm².
 - DMA, interrupt, counter/timer, communications controllers: 20 k gates = 0.2 mm².
 - Communications and Systems NoC interfaces = 0.1 mm².
- The Communications NoC = 9.7 mm².
 - The associative router with 1024 associative entries is ~ 9mm².
 - The algorithmic router with 64k x 3 entries is 0.2 mm².
 - The Arbiter is small ~ 0.1 mm².
 - The Tx and Rx interfaces are small: altogether ~ 0.2 mm².
 - Communications network fabric ~ 0.2 mm².
- The Systems NoC = 4 mm².
 - The 256 kByte System RAM is 2 mm².
 - The Boot ROM is small ~ 0.2 mm².
 - The System Controller with 20k gates is 0.2 mm².
 - The SDRAM controller with 60k gates is 0.6 mm².
 - The network fabric is ~ 100kgates = 1 mm².
- Boot, test and debug = 0.5 mm².

Total area

The total core logic area is thus $76 + 9.7 + 4 + 0.5 = 90.2 \text{ mm}^2$.

Notes

Associative Router = $1024 \times 96 \text{ latches} + 96 \text{ gates} = 500\text{k gates} = 5 \text{ mm}^2$?

20. Power estimates

Processor

ARM968 (from ARM web site) consumes 0.12 to 0.23 mW/MHz on a 130 nm process, and delivers 1.1 dhrystone MIPS/MHz. Thus, to a good approximation, its power-efficiency is 5,000 to 10,000 MIPS/W and it uses 100-200 pJ/instruction.

neuron dynamics

30 instructions at 1 kHz = 30 kIPS = 3-6 μ W.

connection processing

1,000 inputs at 10 Hz (ave.) and 10 instructions/input = 100 kIPS = 10-20 μ W.

SDRAM access

assume SDRAM uses 250mW at 1 Gbyte/s; accessing 4 bytes costs 1 nJ.

1,000 inputs at 10 Hz (ave.) = 40 kByte/s = 10 μ W.

communications link

2.5V I/Os, 10 pF/wire = 30 pJ/transition

3 transitions/4 bits + EOP = 33 transitions/spike = 1nJ/spike/link.

Router

assume power budget at full throughput of 200 MHz is 200 mW, so 1 nJ/route.

neuron total

at 10 Hz (ave.), with H hops, power = 3-6 + 10-20 + 10 + (1 + 2H)10⁻³ μ W
= 23-36 μ W (routing & inter-chip hops are negligible).

Chip

20 processors x 1,000 neurons/processor x 13-26 μ W = 260-520 mW.

Node

chip + SDRAM = 460-720 mW.

System

1 billion neurons = 50,000 nodes = 23-36 kW.