

UNIVERSITY OF MANCHESTER

SCHOOL OF COMPUTER SCIENCE

SOFTWARE ENGINEERING

Databases on SpiNNaker

Author:

Arthur Ceccotti

Supervisor:

David Lester

April 28, 2016

Abstract

A Database Management System (DBMS) is a suite of programs which manage large structured sets of data [1]. Thus it performs the organization, storage, security and integrity of the user's data.

This report contains details of an approach for implementing a DBMS, namely SpiDB (SpiNNaker Databases), on the Spiking Neural Network Architecture (SpiNNaker), an emerging, highly distributed hardware design optimised for Neural Network simulations. The open-source project covers different implementations of a Key-Value store and a Relational DBMS under the architecture constraints.

As a research project, it has a strong focus on the evaluation of results, which allows exploration of how the SpiNNaker hardware performs under a non-neuromorphic environment. The conclusions gathered show limitations and strengths of the architecture under this general purpose application, with suggestions on possible modifications. This can serve as feedback for improvements on the ongoing SpiNNaker development.

Acknowledgements

I would firstly like to thank my supervisor, David Lester, for the useful feedback and motivation during the development of my project. Massive gratitude goes also to Alan Bary Stokes, who became a second supervisor to me, always raising interesting insights and introducing me to the SpiNNaker lifestyle, alongside Andrew Rowley and Thomas Heinis.

As always, friends and family will always be in my heart. I love you all.

Contents

1	Introduction	8
1.1	SpiNNaker Overview	8
1.2	Project Aim	8
2	Background	9
2.1	SpiNNaker Architecture	9
2.2	SpiNNaker Communication fabric	10
3	Development	12
3.1	Planning	12
3.1.1	Requirements Analysis	12
3.1.2	Technologies	13
3.2	Design	14
3.3	Implementation	15
3.3.1	Key-Value Store	15
3.3.2	Relational Database	19
3.3.3	User Interface	21
3.4	Testing and Debugging	21
3.5	Challenges	23
3.5.1	Out-of-order execution	23
3.5.2	Unreliable communication	24
3.5.3	API Bugs	25
4	Evaluation	26
4.1	Reliability & transmission rate	26
4.2	Performance benchmark	29
4.3	SDRAM bandwidth utilisation	30
5	Conclusion	32
5.1	Achievements	32
5.2	Future work	33
	Bibliography	34
A	Queries	36
A.1	Query formats	36
A.2	Key-Value Database	37
A.2.1	Definitions	37

A.2.2	PUT	38
A.2.3	PULL	38
A.3	Relational Database	40
A.3.1	Definitions	40
A.3.2	Queries	42
B	Experiments	43
B.1	Specification	43
B.2	Communication Reliability	43
B.2.1	Sending & receiving SDP packets	43
B.2.2	Storing packets	44
B.2.3	Plot data	46

List of Figures

2.1	SpiNNaker chip layout [2]	9
2.2	SpiNNaker CMP and SDRAM [3]	9
2.3	SpiNN-3 [4]	10
2.4	SpiNN-4 [5]	10
2.5	SpiNNaker communication protocols	11
3.1	Development Plan	13
3.2	SpiDB tree structure	15
3.3	Role assignments per chip	15
3.4	Query from host	17
3.5	Root to leaf SDP	17
3.6	Response to host	17
3.7	Graphical User Interface	21
3.8	Debugging	22
3.9	Testing	22
4.1	PUT performance with variable delay	27
4.2	HASH PULL performance with variable delay	28
4.3	Latency scalability under uniform and local traffic distributions [3]	29
4.4	PUT operation performance benchmark	30
4.5	SDRAM bandwidth utilisation [4]	31

List of Definitions

Branch SpiDB application processor in charge of merging, aggregating database entries and maintaining a tree structure of the query plan

Data Tightly Coupled Memory (DTCM) 62-KB of memory dedicated for data storage and program stack, private to each processor

Database Management System (DBMS) Suite of programs which manage large sets of data, organising the structure and retrieval of database entries [1].

Fixed-Route (FR) Communication protocol with custom routing through chips

Host User's machine connected to SpiNNaker board over Ethernet communication

Instruction Tightly Coupled Memory (ITCM) 32-KB of memory dedicated for machine instructions, private to each processor

Leaf SpiDB application processor in charge of processing storage and retrieval of database entries in SDRAM

Million-core machine System composed of ~1200 SpiNN-4 machines (1,000,000 ARM processors)

Multicast (MC) Communication protocol of one source, multiple destinations

Nearest Neighbour (NN) Communication protocol of one source and destination physically adjacent to it

Point-to-Point (P2P) Communication protocol of one source, one destination

Root SpiDB application processor in charge of communicating with host and forwarding database queries to other processors

SpiNNaker Datagram Protocol (SDP) Communication protocol built on top of P2P to allow larger transfers of data between two cores in the system

SpiDB (SpiNNaker Database) Experimental Database Management System on SpiNNaker, developed by me as part of the third year project

SpiNNaker (Spiking Neural Network Architecture) Massively parallel computer system composed of a large grid of ARM cores designed to model billions of neurons in biological real time [4].

SpiNNaker Chip Multiprocessor (CMP) Custom chip composed of 18 ARM processors and SDRAM

Synchronous Dynamic Random Access Memory (SDRAM) 128-MB of dynamic memory, stitch bounded to each SpiNNaker chip, accessible from all cores within it

SpiNN-3 System composed of 4 SpiNNaker chips (72 ARM processors)

SpiNN-4 System composed of 48 SpiNNaker chips (864 ARM processors)

Chapter 1

Introduction

This chapter describes a high level view of the SpiNNaker platform and its main uses. It highlights also the motivation and aims of my project and how it may impact on the improvement of a large scale international research.

1.1 SpiNNaker Overview

”SpiNNaker (Spiking Neural Network Architecture) is a biologically inspired, massively parallel computing engine designed to facilitate the modelling and simulation of large-scale spiking neural networks of up to a billion neurons and trillion synapses (inter-neuron connections) in biological real time.” [4] The SpiNNaker project, inspired by the fundamental structure and function of the human brain, began in 2005 and it is a collaboration between several universities and industrial partners: University of Manchester, University of Southampton, University of Cambridge, University of Sheffield, ARM Ltd, Silistix Ltd, Thales [6]. A single SpiNNaker board is composed of hundreds of processing cores, allowing it to efficiently compute the interaction between populations of neurons.

1.2 Project Aim

This research project involves making use of the SpiNNaker software stack and hardware infrastructure, both optimised for neural network simulations, in order to explore and evaluate its usability and performance as a general purpose platform. This has been achieved through the development of a distributed Key-Value store and a simple Relational Database Management System (DBMS). This open-source project has grown to be the *largest non-neuromorphic application in the SpiNNaker API* (SpiNNaker Software Stack), available widely for further research.

SpiNNaker was designed from the outset to support parallel execution of application code while ensuring energy efficiency [7]. This appealed as an extraordinary opportunity to store and retrieve data in a fast, distributed way, under a DBMS. This allows exploration of a broad range of ideas outside of the initial scope of SpiNNaker, testing some of its capabilities and limitations against a non-neuromorphic application, bringing useful feedback to the team.

Given such opportunities, the main aims of this project are:

- **Testing usability** of SpiNNaker as a general purpose platform, exploring its strengths and weaknesses under an environment it was not originally designed for.
- Performing **application benchmarks**, allowing analysis which can provide insights for design decisions to the current architecture and possibly influence on the structure of the next generations of SpiNNaker chips.

Chapter 2

Background

This chapter describes an architectural overview of the SpiNNaker research, hardware specifications and multicore communication used on my database project.

2.1 SpiNNaker Architecture

The basic building block of the SpiNNaker machine is the SpiNNaker Chip Multiprocessor (CMP), a custom designed globally asynchronous locally synchronous (GALS) system [4] with 18 ARM968E-S processor nodes [8] (figure 2.1). The SpiNNaker chip contains two silicon dies: the SpiNNaker die itself and a 128-MByte SDRAM (Synchronous Dynamic Random Access Memory) die clocked at 133-MHz [9], which is physically mounted on top of the SpiNNaker die and stitch-bonded to it (figure 2.2) [10]. The SDRAM serves as local shared memory for the 18 cores within the chip, also utilised for memory based communication [11].

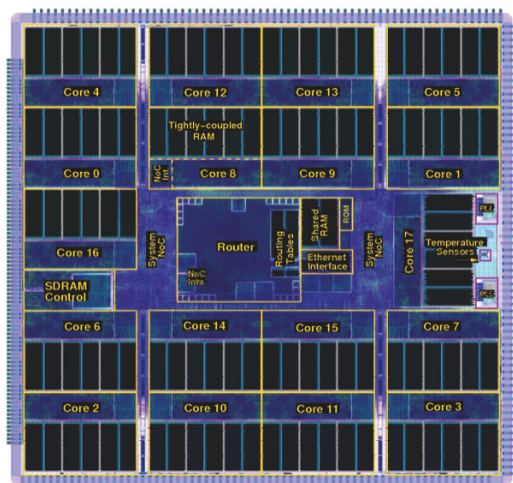


Fig. 2.1: SpiNNaker chip layout [2]

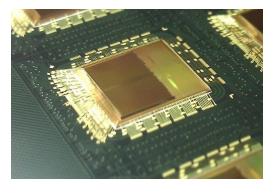


Fig. 2.2: SpiNNaker CMP and SDRAM [3]

Each ARM core within the chip follows a 32-bit Harvard Architecture [12], holding a private 32-KB



Fig. 2.3: SpiNN-3 [4]

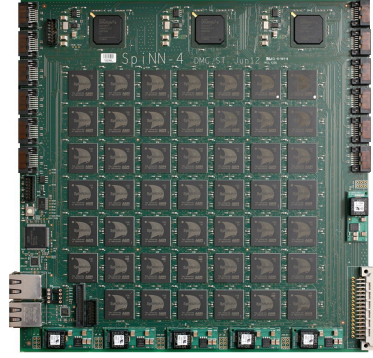


Fig. 2.4: SpiNN-4 [5]

Instruction Tightly Coupled Memory (ITCM) and 64-KB Data Tightly Coupled Memory (DTCM) [4]. It has a small peak power consumption of 1-W at the nominal clock frequency of 180-MHz [13].

Figure 2.3 shows SpiNN-3, a 4 chip SpiNNaker board with 72 processors, next to figure 2.4 (SpiNN-4), composed of 48 chips for a total of 864 processing cores.

2.2 SpiNNaker Communication fabric

Cores on a SpiNNaker board exchange packets through wired connections over a large communication fabric. Each SpiNNaker chip is surrounded by a lightweight, packet-switched asynchronous communication infrastructure [10]. Packet exchange can be used to transmit information initially private to a core and it is managed under the API's even driven architecture, thus incoming packets issue an interrupt on the receiving core, which is woken up from low-power mode. Data transfer and path routing between different chips is handled by an external processor, named the *router*.

There are currently 4 different communication protocols in the system, visualised in figure 2.5. It is worth noting that **none of these protocols guarantee successful delivery of data** and this effect is worsened if there is increased traffic in the communication fabric. Sending a large amount of packets simultaneously is likely to result on packet drops, which is a consequence of the early design decisions of SpiNNaker. In a neuromorphic environment missing spikes (electric signals between neurons in the brain) is common, thus not being a major concern of the hardware architecture, designed for neural network simulations. This results on a large issue when developing applications where high reliability is necessary, as later discussed in chapter 4.

- **Multicast (MC):** The MC protocol, originally designed to simulate asynchronous neural spikes [14], is used when a single source issues information to multiple destinations (one-to-many), as a fan-out topology [2]. The packet contains a 32 bit routing key, used by an external router process to carry out the delivery, and an optional 32 bit payload, both provided at the source.
- **Point-to-Point (P2P):** P2P packets have a single source and a single destination core (one-to-one). Each packet contains a 16-bit source ID, destination ID and an optional 32-bit payload [11]. On top of this layer, the SpiNNaker Datagram Protocol (SDP) was designed to allow transfers of up to 256-bytes of data between two cores, by sending a sequence of P2P packets with payloads [15]. SDP can be used to communicate to the *host* machine, wrapped around a larger UDP packet.

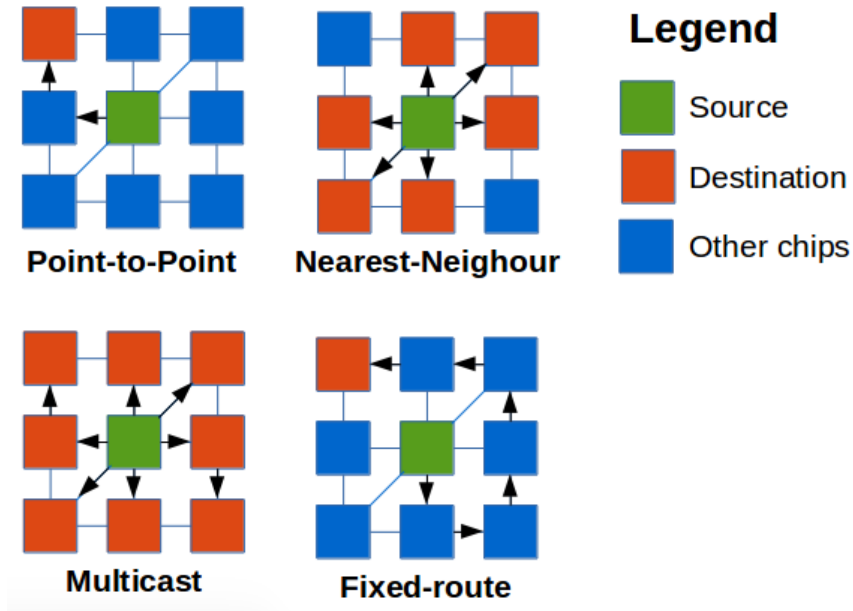


Fig. 2.5: SpiNNaker communication protocols

- **Nearest-neighbour (NN):** Under the NN protocol, packets issued at a chip will only be delivered to the monitor core on a physically adjacent node [2]. NN packets contain a 32-bit payload and a 32-bit address/operation field [11].
- **Fixed-route (FR):** FR packets use the same mechanism as MC, without a key field. Packets being routed by the contents of a register instead of a routing key.

My database application makes extensive use of the SDP protocol, alongside MC.

Chapter 3

Development

This chapter describes the development of my database application, SpiDB, from September 2015 to March 2016, involving planning, testing and implementation. The project is open-source, which can be found at <https://github.com/SpiNNakerManchester/SpiNNakerGraphFrontEnd>.

By being part of the SpiNNaker team in Manchester, I was directly exposed to the ongoing research, frequently receiving feedback on my work. The collaboration was effectively bi-directional, as I was able to constantly find, evaluate and fix inconsistencies and bugs in the Software API not known to the team.

3.1 Planning

The first phase of my project involved making a detailed plan of approach, taking into consideration the time resources and learning curve. Such a plan was useful to set myself deadlines for deliverables in the cycle of short iterations and keep track of progress. A high level chart with the weekly delivery plan can be seen on figure 3.1.

The overall project was divided into two phases: developing a No-SQL Key-Value store for insertion and retrieval of non-structured data and an SQL based Relational Database for creation and manipulation of table structures.

3.1.1 Requirements Analysis

The project plan involved analysing the importance of different requirements based on their relative difficulty and scope within the project aim. Modern database management systems have a broad range of complex requirements. Given limited resources, the following requirements were prioritised for this application:

- **Reliability:** User's queries must complete in a reasonable way. This means any internal errors, inconsistencies or communication failures should be handled from within the database system, avoiding unpredicted failures to the user. This is a difficult task given the unreliability of the SpiNNaker communication fabric, as discussed in section 2.2.
- **Durability:** User's queries with the aim of modifying the database state should persist, being internally stored until removal. The SpiNNaker hardware does not contain permanent storage components, reducing this constraint to "insert operations must persist until shutdown". It would be possible to externally connect the board to a mass-storage device, but it is outside of the scope of this project.

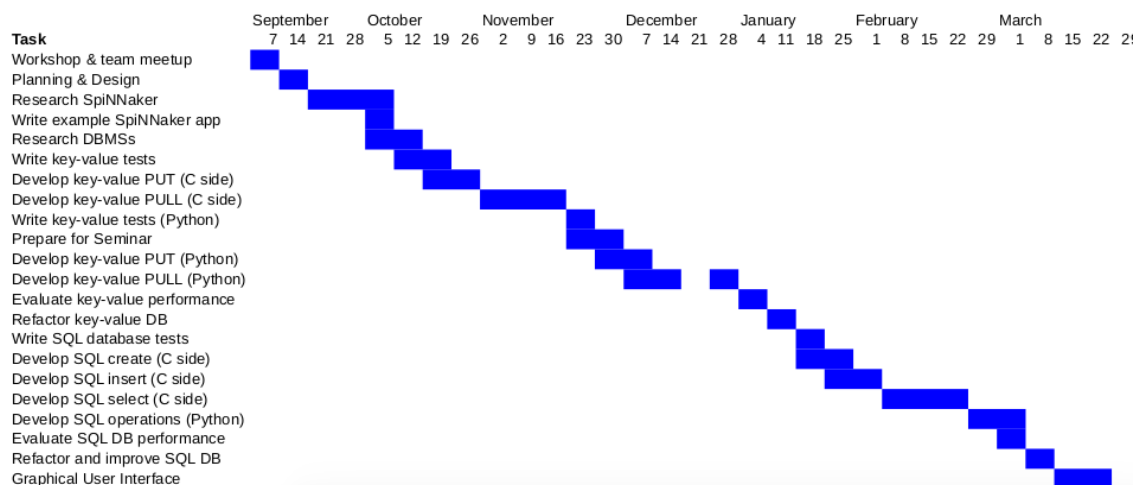


Fig. 3.1: Development Plan

- **Isolation:** Executing transactions concurrently must have the same effect as doing so sequentially. Concurrency control is an important part of this parallel database system, as it aims to handle thousands of concurrent queries distributed among thousands of processing cores. An approach to ensure isolation in the system is further discussed in section 3.5.1.
- **Scalability:** The system must be able to handle many parallel queries and have enhanced performance when given more hardware resources, in this case processor count. This is arguably the most important of all requirements, as the SpiNNaker team is currently focusing on building a large scale machine composed of 1,000,000 processing cores. If this application scales well, it will quickly be able to show the strengths and weaknesses of the machine.

These main requirements do not cover two of the four ACID (Atomicity, Consistency, Durability, Isolation) properties of a database: atomicity and consistency. Atomicity, although very important on a large commercial application, is extremely hard to achieve in a decentralised system, with the use of complex multi-core rollbacks, and falls out of the scope of this experimental project. Consistency is significant when ensuring data written to the database is valid according to constraints, triggers and cascades, which are originally non-existent in a reduced instruction data store and do not contribute to the project research.

Lastly another outstanding requirement not included in the plan was a strong security protocol. Data encryption and authorisation have many advantages, but present themselves as unnecessary complexity for a small, private, experimental project.

3.1.2 Technologies

Part of the project plan involves research on the SpiNNaker hardware and software stack, extensively used on my project. Given the steep learning curve of a low level distributed system and event driven platform, on the 7th of September 2015, I attended the 5th SpiNNaker Workshop, where tens of researchers from around the globe gathered for a one week course on the SpiNNaker Project in Manch-

ester. I was officially introduced to the team, whom I would learn from and work with for the rest of the year.

The SpiNNaker Software Stack is split into two: the Python toolchain, running on the *host* machine (user's computer), and C code, compiled to run on the board. The full software and documentation can be found at <https://github.com/SpiNNakerManchester>, and is composed of over 70,000 lines of code, a lot of which I had to carefully read and learn to use.

These technologies were used to develop the following deliverables:

- **Python:** (2000 lines of code) uploading binaries to the board, status checking, query parsing, Graphical User Interface, data analytics and socket communication.
- **C:** (2500 lines of code) event driven communication between processors, low level memory management, distributed query processing.

3.2 Design

Internally, I designed SpiDB to have a hierarchical tree structure architecture, composed of *root*, *branch* and *leaf* nodes/cores. This structure allows a divide-and-conquer approach to the database query plan, having the aim of improving load balancing, scalability and allow aggregation of data entries.

The following binaries run on different cores on the chip, each with a specific role:

- **root:** Each chip contains a single *root* node, which handles incoming packets from *host* by redirecting them to *branch* and *leaf* nodes, in an intelligent way, for parallel processing. Thus it reads user queries and manages load balancing throughout the board.
- **branch:** The middle layer is composed of 4 *branch* cores in charge of aggregating data returned by *leaf* nodes, serving also as "capacitors", slowing down excessive queries which may overload a destination core.
- **leaf:** The user queries are finally processed by 12 *leaf* nodes per chip, with the aim of storing and retrieving database entries (Key-Value pairs or table rows) on shared memory, unlike other nodes which handle mostly communications.

These roles can be visualised on figures 3.2 and 3.3.

The two remaining cores on the SpiNNaker chip are used for internal use of the system (thus omitted from the diagrams), as they run an instance of *sark*, low level SpiNNaker Operating System, and *reinjector*, in charge of re-issuing dropped Multicast packets.

There are two important advantages, in favour of scalability, which lead me to make this design choice. Firstly, SpiNNaker communication is unreliable, meaning that a lot of traffic and a central point of failure cause large packet drop rate. This hierarchical structure strongly reduces the problem, and it assures cores will never receive packets from more than 4 other cores, distributing queries in an efficient way and protecting cores from excessive incoming packets. Secondly this approach is inheritably beneficial for merges and aggregation (eg. sql keywords COUNT, MAX, MIN, AVG, SUM, etc.), as these can be done on different layers over iterations.

A disadvantage of this design is that less cores perform the actual storage and retrieval of data. Out of 16 application cores, 4 are used only for distribution of the query plan (*root* and *branches*). This consequently impacts performance, as each *leaf* node will be assigned more memory to read and will be kept engaged with query processing.

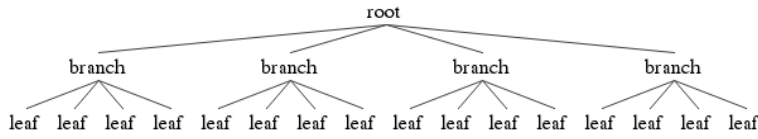


Fig. 3.2: SpiDB tree structure

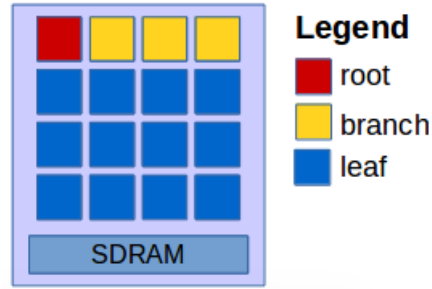


Fig. 3.3: Role assignments per chip

3.3 Implementation

This section describes the final implementation achieved during the development of the project. It outlines the supported database operations, their internal functionality, expected results and performance.

3.3.1 Key-Value Store

The first half of the project involved developing a distributed Key-Value store, in which users can insert and retrieve data entries in a dictionary form. This means the user must be able to set keys mapping to values, both of type *int* or *string*, and retrieve such values when given the same keys. This section describes in detail the internal processing of the insert (*put*) and retrieve (*pull*) queries.

3.3.1.1 PUT

The *put* operation is the main way of inserting data onto the SpiDB distributed database. Upon completion, such operation will store the specified key, mapping to its value, on the memory of an arbitrary chip in the Spinnaker board (as chosen by the *root* core). This operation expects an acknowledgement from the core which stored the entry. It has complexity linear to the size of the input Key-Value and constant to database size.

Example usage:

```

1 put "hello" "world"
2 put 123 "foo"
3 put "life" 42
4 put 1 2

```

Internally the following steps occur:

1. User issues query of format *put* "key" "value" on host machine.
2. Query and metadata are converted into a byte array with the following format, transferred via UDP over Ethernet to the board (figure 3.4).

```

1 typedef struct putQuery {
2     spiDBcommand    cmd;
3     uint32_t        id;
4
5     uint32_t        info;
6     uchar           k_v[256];
7 } putQuery;

```

In this case *SpiDBCommands* is set to the constant representing the put operation, *id* identifies every query uniquely, *info* contains bit encoding of the size and type of both key and value, *k_v* contains key and value appended.

3. Packet arrival triggers interrupt on *root* core, via the SDP protocol, and is decoded.
4. *root* selects a *leaf* core to store the data entries in one of the following ways, specified by the user:
 - **Naive:** as packets arrive, they are assigned to a different *leaf* node in a Round-Robin manner.
 - **Hashing:** the given key is used to produce a 32-bit hash value, which is decoded to assign the query to a specific *leaf* node.
5. *root* communicates to chosen *leaf* node with *put* contents via SDP (figure 3.5).
6. *leaf* node triggers an interrupt and stores Key-Value entry into its dedicated region of SDRAM.
7. *leaf* sends an acknowledgement message over UDP back to host (figure 3.6).
8. User receives timing information and query result.

3.3.1.2 PULL

The *pull* operation is the main way of retrieving data from the SpiDB distributed database. Upon completion, such operation will return the value mapped by a given key, from an arbitrary chip in the SpiN-Naker board, or not respond if such key was not found. This operation expects a response only if the key is present on the database, thus it is an undecidable problem. The reason for this is further explained in section 3.5.2. This operation has complexity linear to the size of the input and linear to database size, although highly improved with a large number of processors.

Example usage:

```

1 pull "hello"
2 pull 123

```

Internally the following steps occur:

1. User issues query of format *pull* "key" on host machine.

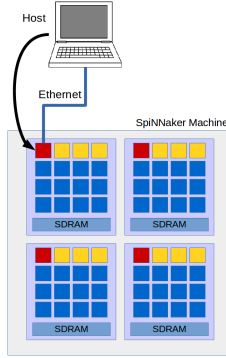


Fig. 3.4: Query from host

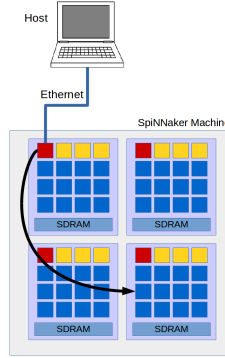


Fig. 3.5: Root to leaf SDP

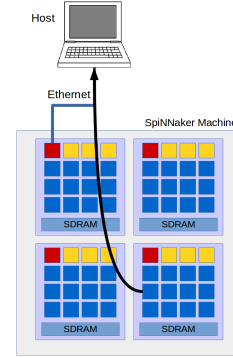


Fig. 3.6: Response to host

2. Query and metadata are converted into a byte array with the following format, transferred via UDP over Ethernet to the board.

```

1 typedef struct pullQuery{
2     spiDBcommand    cmd;
3     uint32_t        id;
4
5     uint32_t        info;
6     uchar           k[256];
7 } pullQuery;

```

Where *SpiDBCommands* represents the pull operation constant, *id* is the query, *info* contains bit encoding of the size and type of the given key, *k* contains the key itself encoded as a byte-array.

3. Packet triggers interrupt on *root* core, via SDP, and is decoded.
4. *root* selects one of the following search strategies, based on database type:
 - **Naive**: there is no knowledge of which, if any, chip contains the given key. Therefore *root* issues a multicast packet to all *leaf* nodes on the board, requesting them to linearly scan their regions of shared memory, searching for the entry.
 - **Hashing**: the key is used to produce a 32-bit hash value, which, if existent on the database, must be present at the memory of a specific core, pointed by the decoding of such hash value. Therefore *root* node sends a single SDP packet to chosen *leaf*, requesting it to search for the given key.
5. Each *leaf* node that received a search request triggers an interrupt and searches for Key-Value mapping in SDRAM memory. If key is found, value is returned over UDP back to host.
6. User receives timing information and query result.

Algorithms 1 and 2 show a high level simplification of code running in the *root* and *leaf* nodes, as described above, through the event driven application, implementing *put* and *pull* operations.

Algorithm 1 Root core

```
1: procedure ONRECEIVE(sdp)
2:   if sdp.command is PUT then
3:     if DB_TYPE is NAIVE then
4:       forward PUT query to next core (Round-robin)
5:     else if DB_TYPE is HASH then
6:        $h \leftarrow \text{hash}(sdp.key)$ 
7:        $chipx \leftarrow h[0 : 7]$ 
8:        $chipy \leftarrow h[8 : 15]$ 
9:        $core \leftarrow h[16 : 24]$ 
10:      forward sdp PUT query to (chipx, chipy, core)
11:   if sdp.command is PULL then
12:     if DB_TYPE is NAIVE then
13:       issue multicast PULL to all cores in the system
14:     else if DB_TYPE is HASH then
15:        $h \leftarrow \text{hash}(sdp.key)$ 
16:        $chipx \leftarrow h[0 : 7]$ 
17:        $chipy \leftarrow h[8 : 15]$ 
18:        $core \leftarrow h[16 : 24]$ 
19:       forward sdp PULL query to (chipx, chipy, core)
```

Algorithm 2 Leaf core

```
1: procedure ONRECEIVE(sdp)
2:   if sdp.command is PUT then
3:     store sdp.key and sdp.value in SDRAM
4:   if sdp.command is PULL then
5:      $entry \leftarrow SDRAM[0]$ 
6:      $i \leftarrow 0$ 
7:     while entry is not null do
8:       if entry.key is sdp.key then
9:         send response to host with sdp.value
10:      return
11:      $entry \leftarrow SDRAM[i]$ 
12:      $i \leftarrow i + 1$ 
```

3.3.2 Relational Database

The second half of the project involves developing an SQL-based distributed database on top of the Key-Value store layer. Entries are stored in a structured way, bounded by the definition of tables. In this RDMS users can create tables with different fields, insert values and retrieve them with given conditions. This section describes in detail the internal processing of the *create*, *insert* and *select* queries.

3.3.2.1 CREATE

In SpiDB, the *create* operation is used to generate a table definition in the SpiNNaker hardware. Data can only be inserted into the database if a corresponding table exists. This query has as parameters the table name, field names and their types (*int* or *string*). This operation expects an acknowledgement from the *root* core, which sets the failure flag if the table definition already exists.

Example usage:

```
1 CREATE TABLE Dogs(name varchar(10), owner varchar(35), age integer);  
2 CREATE TABLE People(name varchar(35), lastname varchar(20));
```

Internally the following steps occur:

1. User issues query of format *CREATE TABLE name(column1 type(size), ...)* on host machine.
2. Query and metadata are converted into a byte array and sent to the board (format can be found on appendix section A.1).
3. *root* core receives and decodes SDP packet.
4. If table does not exist on the database yet, *root* core stores table definition and metadata in its region of shared SDRAM, accessible by other cores for insert/retrieve operations.
5. *root* core sends acknowledgement back and information is displayed to the user.

Complexity: linear to the size of the input, constant to database size.

3.3.2.2 INSERT

New values can be added to the database management system through the *INSERT* query. A single entry is an *int* or *string* value assigned to a column on a given table at a new row. Multiple entries can be safely inserted concurrently, as they are distributed across different cores. This operation expects an acknowledgement from the *leaf* node in charge of storing the value and metadata, with the failure flag set if memory is full.

Example usage:

```
1 INSERT INTO Dogs(name, owner, age) VALUES ("Toddy", "Arthur", 8);  
2 INSERT INTO Dogs(name, owner, age) VALUES ("Guto", "Arthur", 10);  
3 INSERT INTO People(name, lastname) VALUES ("Arthur", "Ceccotti");  
4 INSERT INTO People(name, lastname) VALUES ("Canny", "Dattlin");
```

Internally the following steps occur:

1. User issues query of format *INSERT INTO table(column1, ...) VALUES (value1, ...)* on host machine.

2. Query is broken down into column-value pairs, containing also value type, size and specified table (eg. ["name":("Arthur", type: string, size: 6)]). This step is necessary because SDP packets have a limit size of 256-bytes, thus not being able to carry one single large packet containing all the assignments for a table row. Streaming smaller packets also decreases the need for a large storage buffer at the destination [15].
3. Column-value pairs (entries) are converted into a byte array, sent to the board, being received and decoded by the *root* node.
4. The column-value pair query is forwarded to a *leaf* node in a Round-Robin way.
5. *leaf* node receives query, checks if specified table exists in shared memory SDRAM and reads its definition.
6. *leaf* node stores column-value pair into a new row in reserved region for given table. The SpiDB RDMS is row-based, thus entire rows are stored consecutively in a structured way at all cores.
7. *leaf* node acknowledges query, returned to the user.

Complexity: linear to the size of the input, constant to database size.

3.3.2.3 SELECT

In the SpiDB system, multiple entries from a table can be retrieved using the *SELECT* query, following the standard SQL syntax. Selection criteria can be specified by the user and matching results are streamed until a time-out is reached. If no value in the table matches such criteria, there will be no response from the board, which is a consequence of not making use of internal acknowledgements, as discussed in section 3.5.2. This is a similar behaviour to *pull* requests.

The select query has the capability of producing a very large amount of traffic in the system, as it is the only query with a variable number of packets returned. All other queries have at most one acknowledgement per incoming packet.

Example usage:

```

1 SELECT name FROM Dogs WHERE owner = "Arthur";
2 SELECT name , owner FROM Dogs;
3 SELECT * FROM People WHERE age > 5 and age < 20;
4 SELECT lastname FROM People WHERE name = lastname;
```

Internally the following steps occur:

1. User issues query of format *SELECT field1, field2, ... FROM table WHERE condition*, which is converted and sent as a single packet to the SpiNNaker board.
2. Upon receipt, *root* node issues a multicast packet to every *leaf* node on the board, requesting search for all rows which match the *WHERE* criteria.
3. Every *leaf* node linearly checks appropriate condition and forwards matching column-value pairs on packets to their appropriately assigned *branch* nodes.
4. *branch* nodes aggregate fields if requested, controls speed of execution and sends packet to host with such entry definition.
5. User receives an arbitrary amount of entries, displayed as a table of results.

Complexity: linear to the size of the input, linear to database size.

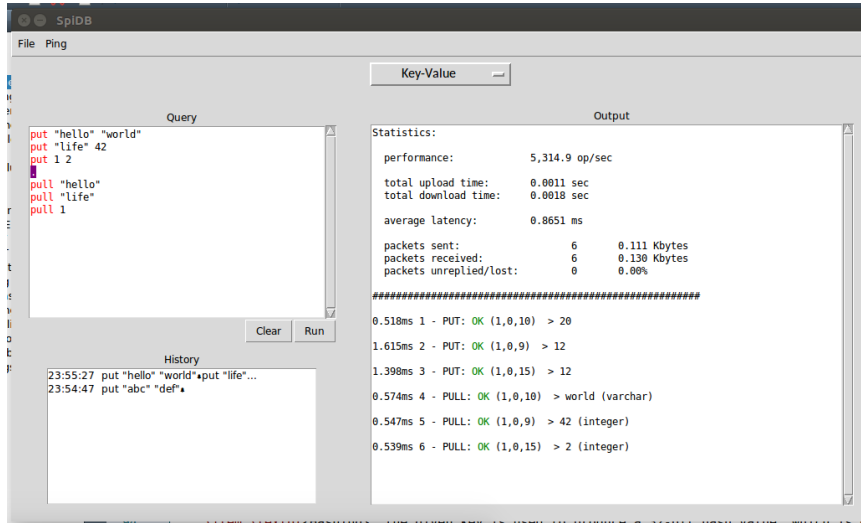


Fig. 3.7: Graphical User Interface

3.3.3 User Interface

Although a graphical user interface was not an essential part of the project plan, a simple one was created for the purposes of demonstration, data plotting and ease of visualisation. By making use of the SpiDB API, the UI includes features to import *sql* files, ping the board, issue concurrent queries, visualise previous queries, display result statistics and graphs. An example instance can be seen in figure 3.7.

3.4 Testing and Debugging

I started the project with a Test-Driven Development approach, writing Python tests with the *unittest* module and C assertions with part of the SpiNNaker API module *debug.h* and my own code. This allowed high reliability from the start. Running tests can be seen in figure 3.9. Testing and development was done from the bottom up, starting from internal memory management, core communication and finally the user interface and host communication. Testing involved initially a number of insert queries with different types and sizes. Given the SpiNNaker architecture is word aligned, it was important to test boundaries along data sizes along multiples of 4 bytes. Examples of tests run can be seen on table 3.1. Both the Key-Value store and the SQL management systems were also tested for their performance and scalability, running over 100,000 simultaneous incoming queries with real data pulled from a modern English vocabulary. The results from these experiments can be read on chapter 4.

Real-time debugging on the SpiNNaker board is relatively hard, but luckily the API provides ways to log messages in each core's private data memory, which can be read by an external process upon execution of the program. Debugging was performed with a tool developed by the team, named *ybug*, available at <https://github.com/SpiNNakerManchester/ybug>, which allows core status checking, uploading binaries, reading logs and memory content, among other functionality [16]. On the *host* Python code, debugging was done using the PyCharm IDE runtime tools.

```

700009f0 00000000 00000001 10041004 00000001 ffffffff 10041004 00000001 00000000
70000a10 10041004 00000001 00000001 10041004 b2bc347d 1a19cb9f 10041004 11b514d4
70000a30 d997c7f3 10041004 10041004 c575452e 10041004 86e811f9 9f0bb90a 10041004
70000a50 74257bc7 6daccb39 10042005 00000003 6c6c6548 0000006f 20051004 6c6c6548
70000a70 0000006f 00000003 20042003 74736554 00676e69 200c202c 7320794d 74726f68
70000a90 79656b20 696b2061 6f20646e 65722066 6974616c 796c6576 6e6f6c20 74732067
192.168.240.253:0,0,0 > lobuf 2
[INFO] (slave/slave.c: 89): Initializing Slave...
[INFO] (slave/slave.c: 27): DataSpec data address is 70000950
[INFO] (src/data_specification.c: 64): magic = ad130ad6, version = 1.0
[INFO] (slave/slave.c: 38): System region: 70000998
[INFO] (slave/slave.c: 39): Data region: 700009b0
[INFO] (slave/slave.c: 47): Initialization completed successfully!
[ASSERT] (slave/slave.c: 30): Failed putting 0x00000000 (s: ) (type: 2) -> 0x00000000 (s: ) (type: 2)
[ASSERT] (slave/slave.c: 30): Failed putting 0x626f6f46 (s: Foobar) (type: 2) -> 0x00000000 (s: ) (type: 2)
192.168.240.253:0,0,0 > snewm 700009b0
700009b0 10041004 ffffffff ffffffff 10041004 ffffffff ffffffff 10041004 00000000
700009d0 00000000 10041004 00000001 00000001 10041004 00000002 00000002 10042005
700009f0 00000003 6c6c6548 0000006f 20051004 6c6c6548 0000006f 00000003 20042003
70000a10 74736554 00676e69 200c202c 7320794d 74726f68 79656b20 696b2061 6f20646e
70000a30 65722066 6974616c 796c6576 6e6f6c20 74732067 676e6972 726f6e20 73657420
70000a50 676e6974 20002000 20062000 626f6f46 00007261 00000000 00000000 00000000
70000a70 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000000
70000a90 00000000 00000000 00000000 00000000 00000000 00000000 00000000 00000000
192.168.240.253:0,0,0 > lobuf 2
[INFO] (slave/slave.c: 89): Initializing Slave...
[INFO] (slave/slave.c: 27): DataSpec data address is 70000950
[INFO] (src/data_specification.c: 64): magic = ad130ad6, version = 1.0
[INFO] (slave/slave.c: 38): System region: 70000998
[INFO] (slave/slave.c: 39): Data region: 700009b0
[INFO] (slave/slave.c: 47): Initialization completed successfully!
[INFO] (slave/slave.c: 96): Starting pull tests.
[ASSERT] (slave/slave.c: 30): Failed putting 0x00000000 (s: ) (type: 2) -> 0x00000000 (s: ) (type: 2)
[ASSERT] (slave/slave.c: 30): Failed putting 0x626f6f46 (s: Foobar) (type: 2) -> 0x00000000 (s: ) (type: 2)
[INFO] (slave/slave.c: 105): Finished pull tests.
192.168.240.253:0,0,0 >

```

Fig. 3.8: Debugging

```

DataSpecification - [~/Spinnaker/Repositories/DataSpecification] - ~/Spinnaker/Repositories/SpinnakerGraphFrontEnd/examples/spiDB/core_root/unit_test
File Edit View Navigate Code Refactor Run Tools VCS Window Help
SpinnakerGraphFrontEnd examples spiDB core_root unit_tests root_put_tests.c
Project Pull h x root_put_tests.c root.c x message_queue.h x cluster_slave.c x
SpinnMachine ~/Spinnaker/Repositories
print
Run spiDB
0:0: 1: [DEBUG] (core_root/root.c: 280): =====
0:0: 1: [DEBUG] (core_root/root.c: 281): ===== TESTS FINISHED =====
0:0: 1: [DEBUG] (core_root/root.c: 282):
0:0: 1: [DEBUG] (core_root/root.c: 283):
0:0: 1: [DEBUG] (core_root/root.c: 284): PUT - Received 14/14
0:0: 1: [DEBUG] (core_root/root.c: 285): - Passes: 7/14
0:0: 1: [DEBUG] (core_root/root.c: 286):
0:0: 1: [DEBUG] (core_root/root.c: 287): PULL - Received 15/14
0:0: 1: [DEBUG] (core_root/root.c: 288): - Passes: 7/15
0:0: 1: [DEBUG] (core_root/root.c: 289):
0:0: 1: [DEBUG] (core_root/root.c: 290):
0:0: 1: [DEBUG] (core_root/root.c: 291):
0:0: 1: [DEBUG] (core_root/root.c: 60): PUT PULL
0:0: 1: [DEBUG] (core_root/root.c: 64): 0: - -> HelloWorld
0:0: 1: [DEBUG] (core_root/root.c: 64): 1: - -> Hellofoo
0:0: 1: [DEBUG] (core_root/root.c: 64): 2: - -> Hellbar
0:0: 1: [DEBUG] (core_root/root.c: 64): 3: FAILED - FAILED -> HelloIgg
0:0: 1: [DEBUG] (core_root/root.c: 64): 4: - -> World
0:0: 1: [DEBUG] (core_root/root.c: 64): 5: FAILED - FAILED -> 12345678
0:0: 1: [DEBUG] (core_root/root.c: 64): 6: - -> !
0:0: 1: [DEBUG] (core_root/root.c: 64): 7: FAILED - FAILED -> A1B2C3ABCD
0:0: 1: [DEBUG] (core_root/root.c: 64): 8: FAILED - FAILED -> A1Test
0:0: 1: [DEBUG] (core_root/root.c: 64): 9: FAILED - FAILED -> 1ATest
0:0: 1: [DEBUG] (core_root/root.c: 64): 10: FAILED - FAILED -> 456Test
0:0: 1: [DEBUG] (core_root/root.c: 64): 11: - -> TestNumberValue159
0:0: 1: [DEBUG] (core_root/root.c: 64): 12: FAILED - FAILED -> A relatively long string... with spaces and other characters in it!uhul
0:0: 1: [DEBUG] (core_root/root.c: 64): 13: - -> uhulA relatively long string... with spaces and other characters in it!
0:0: 1: [DEBUG] (core_root/root.c: 293):
0:0: 1: [DEBUG] (core_root/root.c: 294): =====
0:0: 1:
>

```

Fig. 3.9: Testing

Operation	arg1	arg2	expected result	reason
PUT	"hello"	"world"	success	
PUT	"foo"	"bar"	success	
PUT	"A"	"B"	success	
PUT	42	"life"	success	
PUT	"abc"	123	success	
PUT	1	2	success	
PUT	1.5	2.5	failure	no floating point support
PUT	"hello"	"earth"	failure	key already exists
PUT	(very long string)	"value"	failure	256-byte limit for now
PUT	(very long int)	"value"	failure	256-byte limit for now
PUT	(empty)	"string"	failure	both key and value must be set
PUT	"你"	"好"	failure	ASCII support only
PULL	"hello"	"world"		
PULL	42	"life"		
PULL	"abc"	123		
PULL	"HELLO"	failure		case-sensitive
PULL	(very long string)	failure		256-byte limit for now
PULL	"a new key"	failure		key does not exist

Table 3.1: Examples of Key-Value store test cases

3.5 Challenges

This section outlines a list of different challenges and problems I had to face during the development of the application and how they influenced decision making.

3.5.1 Out-of-order execution

A multi-threaded system must account for the fact that queries may not be executed in the order they are issued, which can be a concern depending on the application. As SpiNNaker is a distributed architecture, it causes problems and inconsistencies. For example, using SpiDB, if we try to execute the following code sequence in order:

Listing 3.1: Non-blocking execution

```

1 put "hello" "world"
2 pull "hello"

```

We are not guaranteed that the *pull* query will retrieve the value "world". As both queries are executed simultaneously on the SpiNNaker board, there is a chance that the *pull* operation will terminate before the *put* does, thus failing to return a value.

As a solution to this dependency constraint, my SpiDB database API includes the syntax symbol "." (dot) which blocks execution until all preceding operations terminate. This allows the programmer to control execution flow, choosing what is allowed to compute out-of-order (in parallel) and what should be sequentialised at a cost of performance. This is a similar concept to the Verilog blocking and non-blocking assignments.

Listing 3.2: Blocking execution

```
1 put "hello" "world"  
2 .  
3 pull "hello"
```

The above code will assure sequential code, meaning the *pull* operation will always return "world". It is usually a good idea to block execution when given a dependency constraint. This can also be done with larger code fragments.

Listing 3.3: Blocking execution

```
1 put "hello" "world"  
2 put "life" 42  
3 put 123 456  
4 .  
5 pull "hello"  
6 pull "life"  
7 pull 123
```

It is worth noting that although non-block operations can cause out-of-order execution, it does not occur very frequently. This is because, when transmitting data to the board, queries are serialised over Ethernet in order of appearance and in addition there is a small forced transmission delay between these packets. This will be further discussed on chapter 4.

3.5.2 Unreliable communication

The SpiNNaker communication protocol can be unreliable, as packets are not guaranteed to arrive at destination, as discussed earlier in section 2.2. This effect is worsened when there is large traffic in the system, so the more packets are being sent, the more packets are lost. While developing the application I had to consider this cost, assuring reduced communication when possible.

Queries such as *put* issue only one packet at a time, allowing an end-to-end acknowledgement without significant packet drops or performance costs. For this reason the *put* query, alongside the SQL commands *create* and *insert* always expect a response, which contains timing information and a flag whether it was successful or not. The steps of running SpiDB queries can be revised in section 3.3.

On the other hand the *pull* and *select* commands issue multicast packets (one-to-many). This means for each query, a packet is sent to up to hundreds of codes, which SpiNNaker is particularly optimised for, but the opposite (many-to-one) is more difficult. Multiple cores sending *SDP* packets to a single core results on a very large packet drop. Only a total of 4 of these packets arrive successfully, as it is the limit of *SDP* channels per core. In the worst case, too many packets have also the capability of crashing the destination core. This means when a packet has multiple destinations, acknowledgements would worsen the situation, as most of these would be dropped and it would highly increase traffic in the system, which itself is a reason for packet loss.

This analysis has lead me to make the design decision of not using internal acknowledgements for every packet sent. In *pull* queries, *leaf* nodes only respond when they have found the specified key, which means there will be at most one acknowledgement. If the key is not found on the database, not a single *leaf* core will respond, which can only be known externally with the use of a time-out.

The advantage of this approach is that it increases speed of execution for successful queries and avoids excessive traffic in the system, increasing reliability. This comes at the cost of slow execution of

instances of *pull* requests where the key does not exist on the database. Using this approach performance is poor when requested entries do not exist. Assuming most of the time users will try to retrieve entries previously inserted on the database, I evaluate this decision to be wise.

3.5.3 API Bugs

I have been one of the few users of a large amount of the SpiNNaker API, currently under construction. This means some of it has not been fully tested, resulting in strange behaviour at times, making debugging of my own application difficult. Besides aiming to gather benchmarks for SpiNNaker, my project itself has been very useful when exposing unexpected errors or inconsistencies in the team's codebase.

Facing these issues was certainly a challenge, as they belonged to domains I had little knowledge of. I was a good opportunity for me to learn and improve code quality of a large, collaborated project, by testing its usability and evaluating its outputs. This means my project involved not only developing the database application, but also collaborating to improve SpiNNaker itself.

The main API bugs I exposed and helped resolve were:

- **Duplicate packets:** when multiple distinct SDP packets were sent from different sources to the same destination, strangely the receiving core would read the same packet duplicated a number of times. This behaviour was highly unexpected to the team, so upon extensive testing, I was able to point the issue to a running processor named *reinjector*, in charge of re-issuing lost packets, and resolve the problem.
- **Data Specification Error:** when uploading binaries to the board, sometimes cores would crash with an internal error state (SWERR), upon evaluating and providing the team with detailed feedback on the issue, I found this to be caused by the SpiNNaker Data Specification framework, which handles allocation of data in shared RAM.

Chapter 4

Evaluation

This research project has among its main aims the evaluation of results, described in this chapter. It is important to analyse how this general purpose application performs under the SpiNNaker environments and what architectural limitations it faces, which is useful feedback to the team. The chapter outlines also performance benchmarks of SpiDB in comparison with other database systems and how it can be improved.

4.1 Reliability & transmission rate

According to my experiments, sending SDP packets immediately between two cores results on a very large packet drop ration. Without explicit delay/wait between the transmission of packets, the number of successful deliveries is only about 10% of those sent for a large number of packets. This effect is strongly reduced by the addition of a small delay of 5-10 μ s, which guarantees over 99% successful deliveries, as can be seen on table 4.1. The table also shows that long intervals, greater than 10 μ s, are mostly redundant, as they do not decrease packet loss significantly, but tragically reduce throughput.

Upon experimenting, I found that sending a large amount of consecutive, immediate packets (at least 1,000,000) has a chance of consequently crashing the destination core or causing it to reach the erroneous watchdog state (WDOG), meaning it cannot respond to commands.

Note that the experiment, made on SpiNN-3, involved sending a single SDP packet multiple times to the same destination on the same chip, containing only the 8 byte SDP header. The destination did not store or read the packet contents, only incrementing a counter upon packet receipt. This was used to show the maximum possible transmission rate allowed by the SDP protocol under the hardware communication fabric. Code snippets and more information can be found on the appendix under the appendix section B.2.

Ideally these packets would send useful information, to be read upon arrival, which would keep the destination busy. From my experience and input from the team, the best way to achieve this is by

Packets sent	Successful deliveries (%)				
	no delay	2 μ s delay	5 μ s delay	10 μ s delay	100 μ s delay
50,000	9.56%	57.73%	95.95%	98.36%	99.85%
100,000	12.15%	54.97%	97.99%	99.13%	99.92%
200,000	13.07%	50.55%	99.01%	99.33%	99.96%
500,000	12.97%	50.08%	99.49%	99.80%	99.99%
1,000,000	13.05%	45.06%	98.84%	99.88%	99.99%

Table 4.1: Successful SDP deliveries with delay between each packet

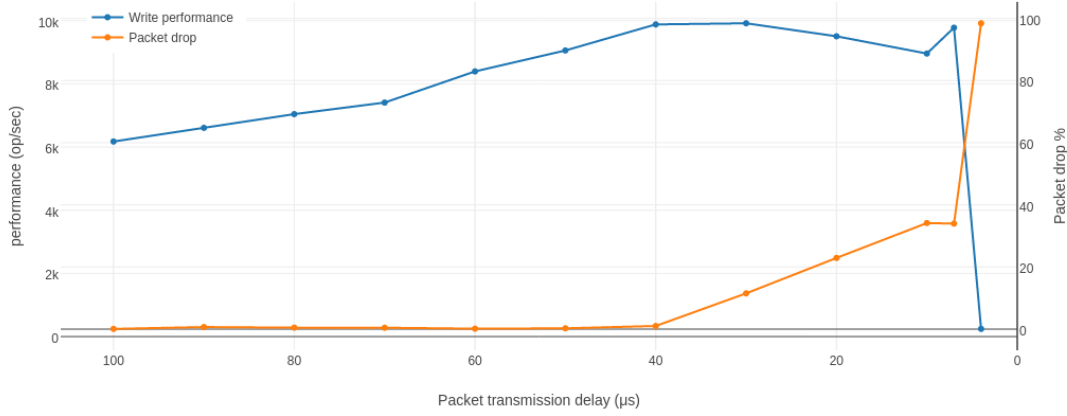


Fig. 4.1: PUT performance with variable delay

immediately storing the SDP packet contents into a buffer when it is received and then handling it at a different point in time (listing B.3 under appendix). The reason for this is that if another packet arrives as we are processing the current packet with same or higher priority, the incoming packet will drop.

High priority should be assigned to storing incoming packets into a buffer and the actual processing should have lower priority, as it can be handled at a later point in time. It is important to note that this can cause starvation and buffer overflow if there are too many immediate packets being received. For instance, if our SDP packet is of size 12-byte (8-byte header + 32-bit data) and stored into a buffer in private DTCM memory upon receipt, we would only ever be able to hold up to about 5,000 messages at once (64-Kbytes DTCM size / 12-byte packet size). Realistically a large part of DTCM contains the stack, local and global variables, so that number will be drastically reduced. In my application, SpiDB, insert and retrieve database queries have a size of 256-bytes, which means a limit of 250 entries in the buffer if memory were empty.

A measured test, transmitting data from a host machine to an Ethernet attached SpiNNaker chip and then to an Application Processor on that node via shared memory, achieved speeds in excess of 22 Mb/s [3].

This evaluation is important because it allows finding the optimal transmission delay and latency for an application. A developer on SpiNNaker using the SDP protocol, which is experimental and yet to be optimised [3], needs to find the balance between transmission rate and reliability, which are inversely proportional.

On SpiDB, I experimented with different time delays when transmitting packets from host over Ethernet to a core on SpiNNaker, in order to find the best evaluation. The speed of one operation is calculated as the time of successful reply minus the time that operation was sent from host, thus being a round-trip-time plus internal processing time. The performance is calculated as the amount of operations successfully executed in one second.

Figure 4.1 plots the performance (left Y axis) and packet drop percentage (right Y axis) of PUT

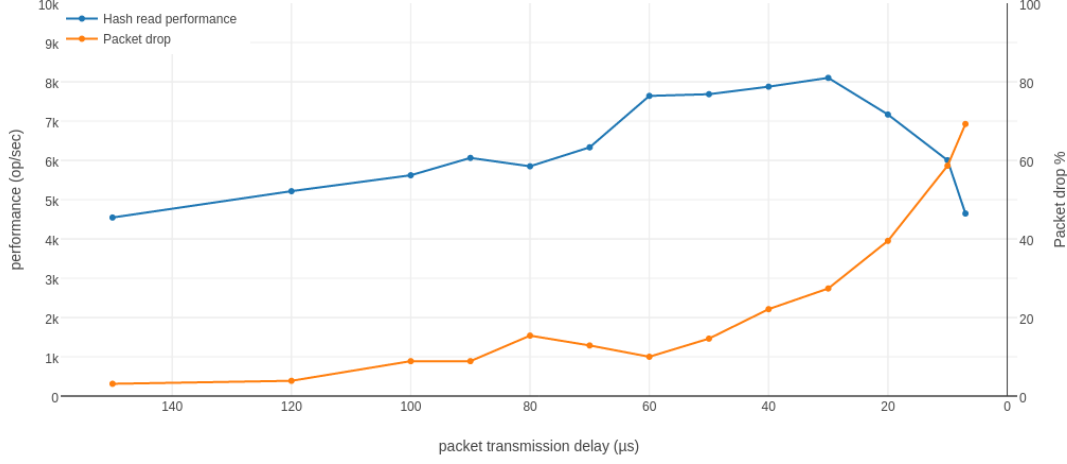


Fig. 4.2: HASH PULL performance with variable delay

operations with the decrease of transmission delay. As can be seen, large packet transmission delays of 50-100μs are redundant, as they do not reduce packet drops, while being a high cost on performance. Naturally the more packets we can send in one second influences the speed of replies, thus improving performance. This hits a maximum at 40μs, with almost 10,000 op/sec, in which transmission is at a high rate with no loss of reliability. Transmission delays between 40-10μs result on a decrease of performance, because although packets are sent more frequently, a lot of them are dropped (up to about 35%), being also extremely unreliable. Naturally performance of is inversely proportional to packet drop rate, as unsuccessful queries do not contribute to the final result. We reach the worst case at 10-5μs delay, when the destination core cannot cope with the speed of incoming packets, simply crashing and ceasing to respond. A very similar behaviour is visible for *pull* operations, as plotted on figure 4.2.

These SpiDB experiments were performed with 100,000 *put* operations with keys and values of size 4-bytes each. More information on the data gathered for these experiments can be found on the appendix under section B.2.

These results are largely influenced by the latency of the SDP datagram, driving the need of such delays. A simulation framework described by Navaridas et al. [8] has shown how latency scales with system size by simulating a range of system configurations from 16×16 (256 nodes) to 256×256 (65,536 nodes) on SpiNNaker. The average and maximum latencies obtained after a simulation of 100,000 cycles are plotted in figure 4.3 [3]. End-to-end latency scales linearly with the number of nodes, i.e. $O(N)$. The maximum latency is that needed to reach the most distant nodes. This means when designing a SpiNNaker application, a developer needs to take into account the locality of communications and the expected latency, which can be used to determine transmission rates.

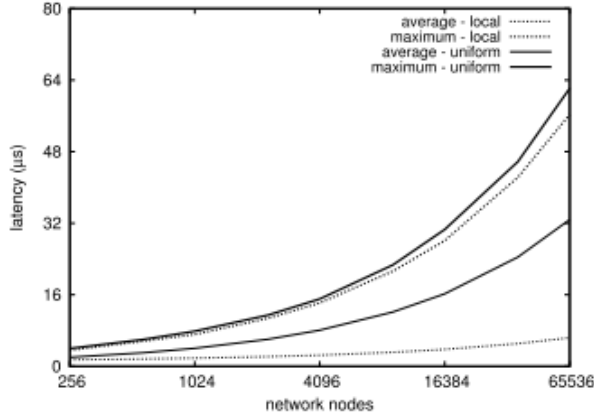


Fig. 4.3: Latency scalability under uniform and local traffic distributions [3]

Local distribution emulates close proximity destinations. Traffic following a uniform distribution represents a very pessimistic scenario in which locality is not exploited.

4.2 Performance benchmark

In order to evaluate the performance of my Key-Value database management system, I compared it with two of the most popular RAM-based Key-Value pair databases: **Memcached** [17] and **Redis** [18]. NoSQL databases have recently become particularly popular [19] and these two systems have been on the market for multiple years, thus serving as a good comparison to my project, SpiDB.

Figure 4.4 shows a performance plot of the SpiNNaker Database (SpiDB) against these competitors for 100,000 consecutive *put* operations. It can be seen that SpiDB runs at about 7,000-8,000 operations per second, which is slower than the others by a factor of 8. Performance for all systems was mostly constant given input size, up to the maximum SDP data size of 256-bytes. More specification and data gathered can be found at appendix under section B.1.

The current SpiDB implementation is new and can be highly enhanced with modifications to the network communication and the addition of caching, which increases the overall performance, enhancing the application in comparison to the other systems. Nevertheless there is a limit to how much improvement can be brought using the SDP protocol, as it has a network and memory bottleneck, earlier discussed in sections 4.1 and 4.3. This means improvements to the code would never allow it to reach performances above around 10,000 operations per second, which is still very slow in comparison with the others.

Another weakness to the approach is the fact that the *root* core is a central point of failure, penalising performance. As of now the SpiNNaker hardware can only host one IP server, at a single chip, to communicate over Ethernet to *host*, meaning incoming and outgoing queries must pass through it. One solution to this is using another external connection, present on SpiNN-4 machines, named *spinnlink*, which allows distributed data transfers.

Although from an initial evaluation it may seem that SpiNNaker is not capable of hosting a database efficiently, that can be proven wrong. It must be outlined that this performance cap is the limitation of a single, independent SpiNNaker board. The SpiNNaker hardware was designed to be inter-connected,

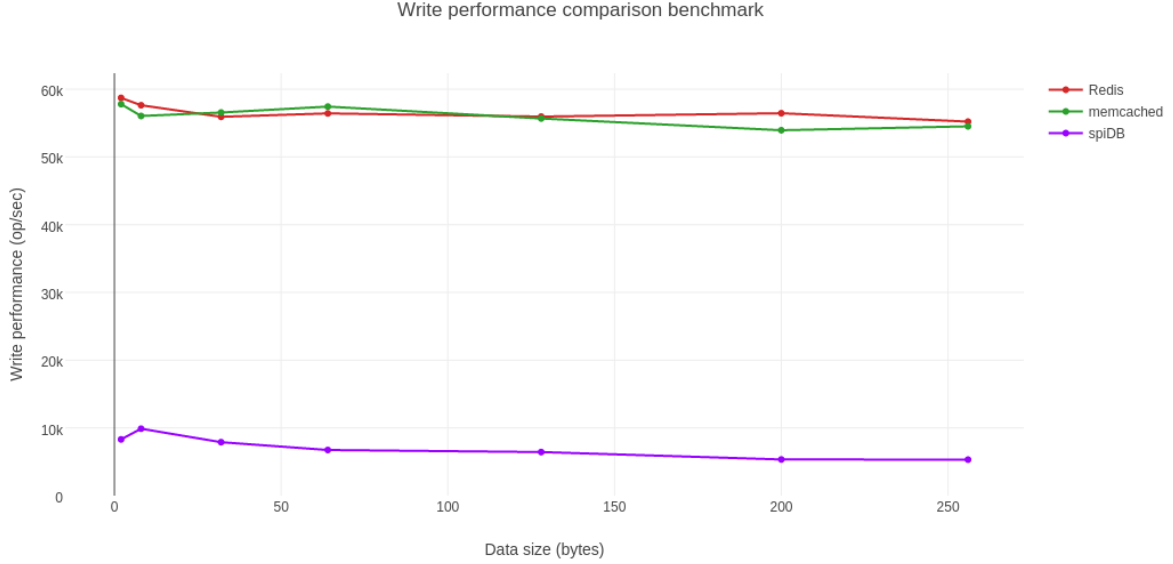


Fig. 4.4: PUT operation performance benchmark

allowing transfers of data between boards, which can be used to add another degree of scalability. This means multiple boards can be connected to the host/user and process a set of separate queries in a fully parallel way. Ideally this would result in a linear performance increase proportional to the number of boards connected, meaning that two boards will run the database operations twice faster than one board alone and so forth. From this analysis I can state that an instance of SpiDB running on 8-10 different 4-chip boards will perform better than *Redis* or *memcached* on commercial hardware. The SpiNNaker hardware is expensive to produce for now, but with a decrease on its prices in the future, it can be a very strong candidate for hosting a DBMS.

This scalability work has not yet been fully tested, thus it may serve as constructive speculation for now. It should also be noted that realistically a perfectly linear performance increase with the increase of horsepower is not common, as other factors need to be taken into account, such as synchronisation, fault tolerance, etc.

4.3 SDRAM bandwidth utilisation

As discussed in section 2.1, a single SpiNNaker chip contains 18 cores sharing 128-MB of dynamic memory (SDRAM), used to store/retrieve data and transfer certain packets between cores. An important stage of the design of a SpiNNaker application is managing the share of available resources, such as SDRAM, to maximise collective utilisation and performance.

In order to find an efficient way to collectively access SDRAM, its bandwidth has to be taken into account. Figure 4.5, plotted by E. Painkras et al. [4], shows the aggregate SDRAM bandwidth utilisation of a number of cores on a single SpiNNaker chip, increased progressively from 1 to 14. The plot outlines the DMA transfer rate for each core as they simultaneously access the shared resource [5]. It is clear

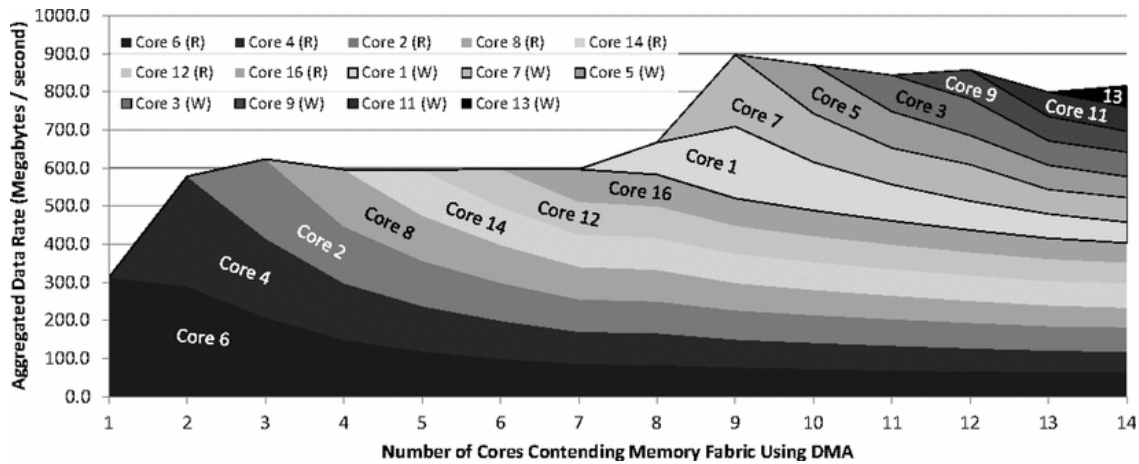


Fig. 4.5: SDRAM bandwidth utilisation [4]

from the figure that the read channel saturates just over 600 MB/s and the write channel adds around 300 MB/s on top of that, for a total aggregate bandwidth of 900 MB/s. Because memory reads and writes share the same network channel, an increase in the number of writes decreases the read bandwidth. Memory bandwidth, although saturated, is distributed fairly among all cores [4].

The progression of Core 6 on the plot shows that if the DMA controller is free, a single core can read from shared memory at 300 MB/s, but if such memory is being constantly accessed by other processors, each will only retrieve data at a rate of around 50 MB/s (6 times slower than doing so exclusively). This evaluation is important because it shows the limitation of all cores frequently accessing shared memory simultaneously, which is exactly the case of the SpiDB *pull* and *SELECT* queries.

It is worth noting that the core's private memory (DTCM), although 2000 times smaller than SDRAM (64-KB vs 128-MB), is 50 times faster to access and does not suffer from concurrent access limitations, which makes it a strong candidate for the implementation of a *cache* memory. This has not been yet implemented, but with appropriate synchronisation it would highly reduce SDRAM access and reduce response times. On SpiDB another enhancement to minimise unnecessary SDRAM access is the use of a hash-table structure, in which the location of values inserted on the database is always known, avoiding a linear scan through the system.

Chapter 5

Conclusion

SpiNNaker is a highly distributed hardware architecture originally designed to simulate the interaction between billions of neurons in the brain. As described in chapter 2, the SpiNNaker hardware is composed of a grid of hundreds of processors with the ability to communicate over different protocols and collectively reach the solution to a problem.

This project involved making use of such hardware to build a Database Management System, namely SpiDB, divided into a NoSQL Key-Value store and an SQL Relational Database. The implementation was tested in various ways in order to find the optimal solution for processing database queries on SpiNNaker, such as the use of both a *naive* and *hash* storage strategies, covered in chapter 3.

This project aims to explore SpiNNaker from a perspective it was not built for, which in this case is a database system. This allows benchmarking and evaluation of results which outline the architecture's strengths and weaknesses under a non-neuromorphic application. Such evaluation may be significant feedback for the current SpiNNaker research.

In the context of a general purpose application, as discussed on chapter 4, the main SpiNNaker limitations demonstrated by this project were:

- SDRAM bandwidth saturation with concurrent memory access
- Central point of failure and slow-down with Ethernet communication
- Frequent packet drops with bursts of SDP and P2P packets
- Limited core private memory DTCM to store incoming packets and data

The design decisions and approaches proposed by me have shown themselves to not be able to fully overcome the communications unreliability issue and the central point of failure, but insights and analysis are provided for further improvements. The evaluation presented can be used as guidelines for the development of any general purpose application on SpiNNaker and even for neuromorphic applications themselves.

SpiNNaker has also demonstrated strengths which a general purpose application can profit from:

- High scalability on number of cores and boards when arranged as a tree structure
- Low network latency for a large number of cores

The project aim has been successfully reached and important achievements were made on the way.

5.1 Achievements

This year long project allowed me to not only build an application which may positively impact the SpiNNaker research, by tackling it with challenges it was not designed for, but also granted me with valuable achievements, both personal and professional:

- Collaboration on an IEEE article publication, alongside Alan Stokes, Thomas Heinis and Yordan Chaparov, currently under pending approval, namely *Exploring the Energy Efficiency of Neuromorphic Hardware for Data Management*.
- Contribution to a large, international, research project on emerging hardware.
- Enhanced knowledge on low-level C and Python, languages which I previously had very little experience with.

5.2 Future work

This project has been the tip of the iceberg of what a fully-functional SpiNNaker database management system can be. All code written by me is now part of the official open-source SpiNNaker Software Stack. This means SpiDB is likely to expand in the future or serve as an example application running on SpiNNaker, accessible by researchers and developers around the globe. A large amount of my code is being currently used by the Imperial College London professor Thomas Heinis and the student Yordan Chaparov, who plan on improving it while decreasing power consumption. The topic of *Databases on SpiNNaker* has recently had great attention within the SpiNNaker team, which may result on it becoming a full Ph.D research.

The most significant features for future research I evaluate to be:

- Improving **scalability** and testing on the large scale *million core machine*.
- Enhancing **reliability** and increasing **query sizes**, perhaps by implementing a protocol on top of the SDP and MC layer. As of now packets are limited to 256-bytes, with unreliability during busy times.

These would certainly take full advantage of the SpiNNaker hardware for hosting a Database Management System, bringing useful feedback and value.

SpiNNaker databases should also have the following features, sorted in decreasing order of importance given the project's aims:

- **Caching** and **indexing**: frequently accessed areas of shared SDRAM memory can be cached at the smaller but much faster private DTCM.
- **Self balancing** during idle times. While no queries are being executed, cores could distribute their contents in a balanced way for faster retrieval. Indexing or other pre-processing could also be executed on the meantime.
- **Additional operations** supporting table merges, triggers and aggregations.
- **Security** and **multi-user access**: different sections of the database can have restricted access through credentials checking.
- An application **server** allowing queries to be requested over the internet from different locations.
- A user-friendly, powerful **Graphical User Interface**

Bibliography

- [1] “The free on-line dictionary of computing.” <http://www.dictionary.com/browse/database-management-system>. Accessed: 28-03-2016.
- [2] S. Furber, D. Lester, L. Plana, J. Garside, E. Painkras, S. Temple, and A. Brown, “Overview of the spinnaker system architecture,” *IEEE Trans. Comput.* vol. PP, no. 99, pp. 1,7–9, 2012.
- [3] C. Patterson, J. Garside, E. Painkras, S. Temple, L. Plana, and S. F. J. Navaridas, T. Sharp, “Scalable communications for a million-core neural processing architecture,” *J. Parallel Distrib. Computing*, vol 72, pp. 6–10, January 2012.
- [4] E. Painkras, L. Plana, J. Garside, S. Temple, F. Galluppi, C. Patterson, D. Lester, A. Brown, and S. Furber, “Spinnaker: A 1-w 18-core system-on-chip for massively-parallel neural network simulation,” *IEEE Journal of solid-state circuits*, vol. 48, pp. 1–4,8,9, August 2013.
- [5] E. Painkras, L. Plana, J. Garside, S. Temple, S. Davidson, J. Pepper, D. Clark, C. Patterson, and S. Furber, “Spinnaker: A multi-core system-on-chip for massively-parallel neural net simulation,” pp. 1–4, 2012.
- [6] “Spinnaker project.” <http://apt.cs.manchester.ac.uk/projects/SpiNNaker/project/>. Accessed: 01-04-2016.
- [7] A. Brown, J. Reeve, K. Dugan, and S. Furber, “Discrete event-based neural simulation using the spinnaker system,” pp. 1–3, 2015.
- [8] J. Navaridas, L. A. Plana, J. Miguel-Alonso, M. Lujan, and S. Furber, “Spinnaker: impact of traffic locality, causality and burstiness on the performance of the interconnection network,” *Proceedings of the 7th ACM International Conference on Computing Frontiers, CF10*, pp. 1–2, 2010.
- [9] L. Plana, S. Furber, S. Temple, M. Khan, Y. Shi, J. Wu, and S. Yang, “A gals infrastructure for a massively parallel multiprocessor,” *Design & Test of Computers, IEEE*, vol. 24, no. 5, pp. 454–456, 2007.
- [10] “The spinnaker chip.” <http://apt.cs.manchester.ac.uk/projects/SpiNNaker/SpiNNchip/>. Accessed: 01-04-2016.
- [11] S. Furber, *SpiNNaker Datasheet - Version 2.02*. University of Manchester, January 2011.
- [12] ARM DDI 0311C, *ARM968E-S Tech. Ref. Manual*, January 2004.
- [13] “Arm968 processor.” <https://www.arm.com/products/processors/classic/arm9/arm968.php>. Accessed: 01-04-2016.
- [14] J. Wu and S. Furber, “A multicast routing scheme for a universal spiking neural network architecture,” *The Computer Journal*, vol. 53, no. 3, pp. 280–288, 2010.

- [15] S. Temple, *AppNote 4 - SpiNNaker Datagram Protocol (SDP) Specification - Version 1.01*. SpiNNaker Group, School of Computer Science, University of Manchester, November 2011.
- [16] S. Temple, *ybug - System Control Tool for SpiNNaker - Version 1.30*. SpiNNaker Group, School of Computer Science, University of Manchester, April 2014.
- [17] “Memcached.” <https://memcached.org>. Accessed: 21-04-2016.
- [18] “Redis.” <http://redis.io/>. Accessed: 21-04-2016.
- [19] J. Han, E. Haihong, G. Le, and J. Du, “Survey on nosql database,” pp. 363–366, 2011.

Appendix A

Queries

This chapter contains code snippets of the main SpiDB operations.

A.1 Query formats

This section contains the query structs and enums used for internal and external communication of the database. The two supported variable types are: *int* and *char[]* (string).

```
1 typedef enum var_type { UINT32=0, STRING } var_type;
2 typedef uint32_t id_t;
3 typedef uint32_t info_t;
4
5 typedef enum spiDBcommand {
6     PUT = 0,
7     PULL,
8
9     CLEAR,
10
11     PUT_REPLY,
12     PULL_REPLY,
13
14     CREATE_TABLE,
15     INSERT_INT0,
16     SELECT,
17     SELECT_RESPONSE,
18
19     PING
20 } spiDBcommand;
21
22 typedef struct Response{
23     id_t      id;
24     spiDBcommand cmd;
25
26     bool      success;
27     uint8_t   x;
28     uint8_t   y;
29     uint8_t   p;
```

```

30
31     uchar          pad[3];
32
33     uchar          data[256];
34 } Response;

```

A.2 Key-Value Database

A.2.1 Definitions

```

1  typedef struct putQuery{
2      spiDBcommand    cmd;
3      id_t            id;
4
5      info_t          info;
6      uchar          k_v[256];
7  } putQuery;
8
9  typedef struct pullQuery{
10     spiDBcommand    cmd;
11     id_t            id;
12
13     info_t          info;
14     uchar          k[256];
15 } pullQuery;
16
17 typedef struct pullValue{
18     var_type    type;
19     size_t      size;
20
21     uchar      pad[3]; //padding
22
23     uchar      data[256];
24 } pullValue;
25
26 typedef struct pullValueResponse{
27     spiDBcommand    cmd;
28     id_t            id;
29
30     uchar          pad[3]; //padding
31
32     pullValue      v;
33 } pullValueResponse;

```

A.2.2 PUT

```
1 #define MULTIPLE_OF_4(n) ((n+3)&(!3))
2 #define try(cond) do { if (!cond) return false; } while (0)
3
4 address_t append(address_t* address, void* data, size_t size_bytes){
5     try(address && data && size_bytes);
6
7     address_t old_addr = *address;
8
9     sark_mem_cpy(*address, data, size_bytes);
10    *address += (size_bytes+3) >> 2;
11
12    return old_addr;
13 }
14
15 size_t put(address_t* addr, info_t info, void* k, void* v){
16
17     size_t k_size = k_size_from_info(info);
18     size_t v_size = v_size_from_info(info);
19
20     try(k_size && v_size && k && v);
21
22     append(addr, &info, sizeof(info_t));
23     append(addr, k, k_size);
24     append(addr, v, v_size);
25
26     return sizeof(info_t) +
27           MULTIPLE_OF_4(k_size) +
28           MULTIPLE_OF_4(v_size);
29 }
```

A.2.3 PULL

```
1 var_type k_type_from_info(info_t info){
2     return (info & 0xF0000000) >> 28;
3 }
4
5 size_t k_size_from_info(info_t info){
6     return (info & 0x0FFF0000) >> 16;
7 }
8
9 var_type v_type_from_info(info_t info){
10    return (info & 0x0000F000) >> 12;
11 }
12
13 size_t v_size_from_info(info_t info){
14    return (info & 0x00000FFF);
15 }
```

```

15 }
16
17 pullValue* pull(address_t addr, info_t info, uchar* k){
18     try(info);
19     try(k);
20
21     var_type k_type = k_type_from_info(info);
22     size_t k_size = k_size_from_info(info);
23
24     info_t curr_info = 0;
25
26     while((curr_info = *addr) != 0){
27
28         addr++;
29
30         var_type read_k_type = k_type_from_info(curr_info);
31         size_t read_k_size = k_size_from_info(curr_info);
32
33         var_type v_type = v_type_from_info(curr_info);
34         size_t v_size = v_size_from_info(curr_info);
35
36         if(read_k_size == 0 || v_size == 0){
37             continue;
38         }
39
40         size_t k_size_words = (read_k_size+3) >> 2;
41         size_t v_size_words = (v_size+3) >> 2;
42
43         if(read_k_type != k_type || read_k_size != k_size){
44             addr += k_size_words + v_size_words;
45             continue;
46         }
47
48         uchar* k_found = (uchar*)addr;
49         addr += k_size_words;
50
51         uchar* v_found = (uchar*)addr;
52         addr += v_size_words;
53
54         if(arr_equals(k, k_found, k_size)){
55             pullValue* v = (pullValue*)
56                 sark_alloc(1, sizeof(pullValue));
57             v->size = v_size;
58             v->type = v_type;
59             sark_mem_cpy(v->data, v_found, v->size);
60
61             return v;
62         }
63         else{
64             continue;

```

```

65     }
66 }
67
68 return NULL;
69 }

```

A.3 Relational Database

A.3.1 Definitions

```

1  typedef struct Column {
2      uchar      name[16];
3      var_type   type;
4      size_t     size;
5  } Column;
6
7  typedef struct Table {
8      uchar      name[16];
9      size_t     n_cols;
10     size_t     row_size;
11     size_t     current_n_rows;
12     Column     cols[4];
13 } Table;
14
15 typedef struct createTableQuery { //CREATE TABLE
16     spiDBcommand cmd;
17     id_t         id;
18
19     Table        table;
20 } createTableQuery;
21
22 typedef struct Entry{
23     id_t         row_id;
24     uchar      col_name[16];
25     size_t     size;
26     var_type   type;
27
28     uchar      pad[3]; //padding for the value
29
30     uchar      value[256];
31 } Entry;
32
33 typedef struct insertEntryQuery { //INSERT INTO
34     spiDBcommand cmd;
35     id_t         id;
36
37     uchar      table_name[16];

```



```

38     Entry          e;
39 } insertEntryQuery;
40
41 typedef enum {
42     EQ = 0,        // =
43     NE,            // != or <>
44     GT,            // >
45     GE,            // >=
46     LT,            // <
47     LE,            // <=
48     BETWEEN,
49     LIKE,
50     IN
51 } Operator;
52
53 typedef enum {
54     LITERAL_UINT32 = 0,
55     LITERAL_STRING,
56     COLUMN
57 } OperandType;
58
59 typedef struct Operand {
60     OperandType type;
61     uchar        value[64];
62 } Operand;
63
64 typedef struct Condition {
65     Operand      left;
66     Operator      op;
67     Operand      right;
68 } Condition;
69
70 typedef struct selectResponse {
71     spiDBcommand cmd;
72     id_t          id;
73
74     Table*        table;
75     address_t     addr;
76
77     uchar         n_cols;
78     uchar         col_indices[MAX_NUMBER_OF_COLS];
79 } selectResponse;
80
81 typedef struct selectQuery { //SELECT FROM
82     spiDBcommand cmd;
83     id_t          id;
84
85     uchar         table_name[16];
86     uchar         col_names[MAX_NUMBER_OF_COLS][16];
87

```

```
88 |      Condition      condition ;  
89 | } selectQuery ;
```

A.3.2 Queries

Code containing the SpiDB SQL relational database queries can be found at
<https://github.com/SpiNNakerManchester/SpiNNakerGraphFrontEnd/tree/master/examples/spiDB>.

Appendix B

Experiments

This chapter describes the machine specifications and code instances used to run the experiments and evaluations of this report.

B.1 Specification

All experiments and instances of SpiDB reported on this dissertation were run on an ASUS X555L quad core Intel(R) Core(TM) i5-5200U CPU @ 2.20GHz, 8gb DDR3 RAM @ 1600 MHz, 64-bit Ubuntu 15.04 machine connected over Ethernet to a 4-chip SpiNN-3 board, produced in September 2011.

The *Memcached* and *Redis* benchmarks were run under such specifications, as well as the transmission rate experiments.

B.2 Communication Reliability

B.2.1 Sending & receiving SDP packets

This section contains code snippets to send SDP packets from single source to single destination with variable delays/waits between them. These were used in the experiments described on section 4.1, using core 2 and 1 of chip 0,0.

C code to send SDP packets with variable delay to a single core:

Listing B.1: Source

```
1 void send_sdp_packets(uint32_t number_of_packets, uint32_t delay,
2   uint32_t chip, uint32_t core){
3     sdp_msg_t msg;
4     msg.flags      = 0x87;
5     msg.tag        = 0;
6
7     msg.src_addr   = spin1_get_chip_id();
8     msg.src_port   = (SDP_PORT << PORT_SHIFT) | spin1_get_core_id();
9
10    msg.dest_addr   = chip;
11    msg.dest_port   = (SDP_PORT << PORT_SHIFT) | core;
12    msg.length      = sizeof(sdp_hdr_t);
```

```

13     for(uint i = 0; i < number_of_packets; i++){
14         if(!spin1_send_sdp_msg(&msg, SDP_TIMEOUT)){
15             log_error("Failed to send");
16         }
17         sark_delay_us(delay);
18     }
19 }

```

C code to log how many SDP packets were successfully received at destination:

Listing B.2: Destination

```

1  uint rcv = 0;
2
3  void sdp_packet_callback(register uint mailbox, uint port) {
4      rcv++;
5      spin1_msg_free((sdp_msg_t*)mailbox);
6      return;
7  }
8
9  ...
10
11 spin1_callback_on(SDP_PACKET_RX, sdp_packet_callback, 0);

```

B.2.2 Storing packets

The following code snippet (listing B.3) describes an efficient way of storing incoming SDP packets, to reduce packet drop, and process them under a lower priority, as discussed in section 4.1.

Listing B.3: Storing incoming packets into a queue

```

1  //buffer to store incoming packages
2  sdp_msg_t** msg_cpies = (sdp_msg_t**)
3      sark_alloc(Queue_SIZE, sizeof(sdp_msg_t));
4  uint i = 0;
5
6  void sdp_packet_callback(register uint mailbox, uint port) {
7      //immediately store incoming packet contents into a queue
8
9      i = (i+1)%Queue_SIZE;
10     register sdp_msg_t* m = msg_cpies[i];
11     sark_word_cpy(m, (sdp_msg_t*)mailbox,
12         sizeof(sdp_hdr_t) + SDP_DATA_SIZE);
13
14     spin1_msg_free((sdp_msg_t*)mailbox);
15
16     // If there was space, add packet to the ring buffer
17     if (circular_buffer_add(sdp_buffer, (uint32_t)m)) {

```

```

18         if (!processing_events) {
19             processing_events = true;
20
21             //trigger lower priority request processing
22             if (!spin1_trigger_user_event(0, 0)){
23                 log_error("Unable to trigger user event.");
24             }
25         }
26     }
27     else {
28         log_error("Unable to add SDP packet to circular buffer.");
29     }
30 }
31
32 void process_requests(uint arg0, uint arg1){
33
34     uint32_t mailbox;
35     do {
36         if (circular_buffer_get_next(sdp_buffer, &mailbox)) {
37             sdp_msg_t* msg = (sdp_msg_t*)mailbox;
38
39             ...
40
41         }
42         else {
43             processing_events = false;
44         }
45     } while (processing_events);
46
47     ...
48
49     //priority assignment
50     spin1_callback_on(SDP_PACKET_RX, sdp_packet_callback, 0);
51     spin1_callback_on(USER_EVENT, process_requests, 2);

```

B.2.3 Plot data

The following table contains data from multiple runs of SpiDB instances with 100,000 *put* operations under different transmission delays, used to produce figure 4.1.

interval (μ s)	packet drop (%)	performance (ops/sec)
100	0.076	6175
90	0.670	6608
80	0.475	7042
70	0.445	7408
60	0.150	8394
50	0.265	9057
40	1.010	9882
30	11.510	9918
20	22.910	9506
10	34.140	8960
7	33.965	9778
4	98.432	243