

DETECTION OF SHARED COPY NUMBER
VARIATION IN PATIENTS CO-HYBRIDISED TO THE
SAME ARRAY

ALED JONES

90565930

A dissertation submitted to the University of Manchester for the
degree of Master of Science in the Faculty of Medical and Human
Sciences;

The School of Medicine

2016

CONTENTS

1	INTRODUCTION	15
1.1	Array CGH	15
1.2	Copy Number Variation (CNV)	16
1.3	Application of array CGH in the clinical setting	16
1.4	Patient to patient hybridisation	16
1.5	Project rationale	17
1.6	Project aims	19
1.7	Array design	19
1.7.1	Feature extraction	19
1.7.2	Feature extraction file	20
1.8	Use of algorithms to detect copy number variation	21
1.8.1	Z score algorithm	22
1.8.2	Aberration Detection Method (ADM-1/ADM-2)	22
1.8.3	Circular binary segmentation (CBS)	23
1.8.4	Hidden Markov model	24
1.8.5	Algorithm performance	25
1.9	Data storage and manipulation	25
1.9.1	Relational database	25
1.9.2	Non-relational database	26
1.9.3	Python	26
1.9.4	R statistical software package	26
2	METHODS	27
2.1	Shared Imbalance Detection (SID) tool	27
2.2	Requirements gathering	27
2.2.1	Functional requirements	27
2.2.2	Non functional requirements	28
2.2.3	Project management	29
2.3	Database infrastructure	29
2.3.1	SQL server	29

2.3.2	Migration from development server (MYSQL) to production server (MS-SQL)	29
2.3.3	HeidiSQL client	30
2.4	SID - database design	30
2.4.1	Indexes and keys	31
2.5	Importing and processing feature extraction files	31
2.5.1	Python	31
2.5.2	Version control	32
2.6	SID - pprocessing feature extraction files	32
2.7	Detection of CNV	32
2.7.1	Choice of algorithm (Z score)	32
2.7.2	Creating a normal reference range	34
2.7.3	Determining if a probe is normal or ab- normal	34
2.7.4	Duplicated probes	35
2.7.5	Determining if a region is normal or ab- normal	35
2.8	Creation of true positive feature extraction files	37
2.8.1	Training cases	37
2.8.2	Test cases	38
2.8.3	Prospective cases	39
2.9	Definition of true positive and false positive calls	39
2.9.1	Unknown true positive calls	39
3	RESULTS	40
3.1	Training cases	40
3.1.1	Z score threshold	40
3.1.2	Removal of cases with high false positive calls	41
3.1.3	Comparison of true positive and false pos- itive calls	41
3.2	Test Cases	46
3.2.1	A greater sensitivity was seen using more representative feature extraction files	46
3.2.2	Suitability of reference range	46
3.2.3	Proportion of CNV regions covered by calls	46
3.3	Prospective cases	48

3.4	Performance of SID	49
4	DISCUSSION	50
4.1	Suitability of SID	50
4.2	Z score threshold to define a probe as abnormal	50
4.3	The number of probes within true positive calls	51
4.4	Combining overlapping trios of abnormal probes	51
4.5	Benign CNV	51
4.6	Calls smaller than 5 probes	52
4.7	Different thresholds for gain or loss	52
4.8	Alternative approaches to consecutive probes analysis	52
4.8.1	Larger sliding windows	52
4.8.2	Sliding window - Regions of interest	53
4.9	Further work	54
4.9.1	Sex chromosomes	54
4.9.2	Sample type	54
4.9.3	What happens when shared CNV is detected?	54
4.9.4	Exporting results to LIMS	55
4.9.5	Initiation of SID	56
4.9.6	Migration of database to the server	56
4.10	Long term management of SID	56
4.10.1	Unit test	56
4.10.2	Code review	57
4.10.3	Recalculation of the reference range	57
4.10.4	Truncating the database to maintain performance	57
4.10.5	Contingency plan	58
4.10.6	Database backup	58
4.10.7	Validation	59
4.10.8	Documentation	59
4.11	Future considerations	60
4.11.1	Array redesign and multiple array designs	60
5	CONCLUSION	61
A	APPENDIX	62

A.1 Analysis of calls made in samples with the 750 probe deletion on chromosome 22	62
---	----

BIBLIOGRAPHY	65
--------------	----

Word count: 7552

LIST OF FIGURES

Figure 1	A trace showing how a shared imbalance would not be detected	17
Figure 2	A scanned 8 x 60k array	20
Figure 3	Circular binary segmentation algorithm	23
Figure 4	Agile development methodology	29
Figure 5	The SID database schema	30
Figure 6	A visual representation of the SID python script	33
Figure 7	The distribution of signal intensities from 30 arrays	35
Figure 8	The normal distribution	36
Figure 9	Consecutive probe analysis algorithm	36
Figure 10	Trios of abnormal probes are combined into a single call	37
Figure 11	The number of probes classified as normal within known regions of CNV at a range of Z score thresholds	40
Figure 12	Training cases: False positive calls are smaller than true positive calls.	42
Figure 13	Training cases: The lowest confidence Z score within each call at a threshold of 2.374	43
Figure 14	Training cases: The average Z score of probes within calls at a threshold of 2.374	44
Figure 15	Training cases: The lowest confidence probe within calls at a threshold of 3.55	45
Figure 16	Training cases: Average Z scores for calls made with a threshold of 3.55	45
Figure 17	Test cases: the proportion of probes within reported CNV included by a call at a threshold of 3.55	47

Figure 18	Prospective cases: The number of calls per array across a range of threshold in 121 prospective cases	49
Figure 19	Classifying overlapping shared CNV	53
Figure 20	Test cases: The average Z score at a threshold of 2.374	63
Figure 21	Test cases: The lowest confidence probe within calls at a threshold of 2.374	63
Figure 22	Test cases: The number of probes within calls made at a threshold of 2.374	64
Figure 23	Test cases: The highest confidence probe within calls made at a threshold of 2.374	64

LIST OF TABLES

Table 1	Cases with one of the 9 most common CNV were used to create the test set	38
Table 2	Calls made in 86 training cases	41
Table 3	Calls made in 81 training cases after removal of outliers	42
Table 4	Test cases: Calls made at a range of thresholds	46
Table 5	Prospective cases: The prevalence of calls	48
Table 6	Prospective cases: The number and size of calls	49
Table 7	Test cases: The calls made on the arrays made from the cases with the 750 probe deletion on chromosome 22	62

ACRONYMS

ADM Aberration Detection Method

AFU Arbitrary Fluorescence Units

CBS Circular Binary Segmentation

CGH Comparative Genomic Hybridisation

csv comma separated values

CNV Copy Number Variation

Cy3 Cyanine 3

Cy5 Cyanine 5

IDE Integrated Development Environment

GUI graphical user interface

HMM Hidden Markov Model

LIMS Laboratory Information Management System

PGD Preimplantation Genetic Diagnosis

RDBMS Relational Database Management System

SID Shared Imbalance Detection

SOP Standard Operating Protocol

SQL Structured Query Language

Dedicated to Sonia

ABSTRACT

Introduction: Array CGH is a well established technique to detect copy number variation (CNV) in diagnostic genetic laboratories. To reduce the cost of the array service a patient to patient hybridisation approach was adopted, halving the costs of consumables. This however, introduces additional challenges, as arrays are analysed by directly comparing hybridisation partners, if a CNV is present in both hybridisation partners relatively there is no change so the CNV is not detected. This is currently mediated by manually selecting hybridised partners based on phenotype, a time consuming process.

The aim of this project is to investigate a tool which can detect CNV independent of the hybridisation partner, enabling shared CNV to be detected.

Methods: The shared imbalance detection (SID) tool was created, consisting of a database and Python script, which performs a Z score analysis on the signal intensities of each probe. Consecutive probes scored as abnormal in both hybridisation partners are reported as shared CNV. Two true positive datasets were created to train and test the Z score threshold.

Results: A range of Z score thresholds were applied to define a probe as normal or abnormal. A call was made in the expected region in 94% of cases in a training set of 81 arrays, however this dataset was found not to be an accurate representation of shared CNV. A more representative test set of 69 arrays had calls in the expected region in 100% of cases with no false positive calls using a Z score threshold of 3.55. 121 unmodified arrays were analysed in a prospective study.

Discussion: SID can detect shared CNV in line with the departmental analysis protocols and in a timely manner. Long term management and further work, including incorporating analysis of the sex chromosomes and validation and implementation the service is described.

DECLARATION

No portion of the work referred to in the dissertation has been submitted in support of an application for another degree or qualification of this or any other university or other institute of learning;

Aled Jones

INTELLECTUAL PROPERTY

- i The author of this dissertation (including any appendices and/or schedules to this dissertation) owns certain copyright or related rights in it (the "Copyright" and s/he has given The University of Manchester certain rights to use such Copyright, including for administrative purposes.
- ii Copies of this dissertation, either in full or in extracts and whether in hard or electronic copy, may be made only in accordance with the Copyright, Designs and Patents Act 1988 (as amended) and regulations issued under it or, where appropriate, in accordance with licensing agreements which the University has entered into. This page must form part of any such copies made
- iii The ownership of certain Copyright, patents, designs, trade marks and other intellectual property (the "Intellectual Property") and any reproductions of copyright works in the dissertation, for example graphs and tables ("Reproductions"), which may be described in this dissertation, may not be owned by the author and may be owned by third parties. Such Intellectual Property and Reproductions cannot and must not be made available for use without the prior written permission of the owner(s) of the relevant Intellectual Property and/or Reproductions.
- iiii Further information on the conditions under which disclosure, publication and commercialisation of this dissertation, the Copyright and any Intellectual Property and/or Reproductions described in it may take place is available in the University IP Policy (see <http://documents.manchester.ac.uk/display.aspx?DocID=487>), in any relevant Dissertation restriction declarations deposited in the University Library, The University Library's regulations (see <http://www.manchester.ac.uk/library/aboutus/regulations>) and in The University's Guidance for the Presentation of Dissertations.

ACKNOWLEDGMENTS

Firstly I would like to express my extreme thanks to Wook for his vision and constant support, advice and guidance throughout the course of this research project and the whole STP.

Similarly I want to thank Kevin for his expert advice and other colleagues, particularly those within the array team at the Viapath Genetics Laboratories for their help and support which helped make this project so enjoyable.

I'd also like to thank the 2013 STP cohort for being so supportive and forming such a tight and responsive network to help address the many obstacles encountered along the way.

Finally I'd like to thank Sonia for her support through this project and the last three years in general.

INTRODUCTION

1.1 ARRAY CGH

Array Comparative Genomic Hybridisation (CGH) is a commonly used diagnostic test in clinical genetics. The test utilises many probes (from 44,000 to 1,000,000) affixed to a glass slide. Each probe consists of many copies of an oligonucleotide 40-60 base pairs long, designed to target a specific region of the genome. Two samples are hybridised to the array in a competitive reaction, usually a diagnostic sample in competition with a reference sample.

The diagnostic and reference samples are labelled with a fluorescent dye (patient Cyanine 5 (Cy5) and reference Cyanine 3 (Cy3)). Once hybridised the array is scanned, capturing the excitation of each dye, producing a high resolution image of the array [1]. This image undergoes feature extraction to calculate the Cy5 and Cy3 signal intensities at each probe and associated quality scores.

This data can then be analysed for copy number variation by applying one of a number of available algorithms: Probe scores are normalised, split into segments of equal copy number and each segment assigned a copy number status [2].

Dividing the signal intensity of Cy5 by that of Cy3 produces a ratio where equal copy number is 1. This ratio is then logged (base 2) to produce the log2 ratio where equal copy number is 0, a decrease in copy number is indicated by a negative value and increased copy number a positive value.

1.2 COPY NUMBER VARIATION (CNV)

Copy Number Variation (CNV) is a deviation in copy number from a reference genome which typically contains 2 copies of a DNA segment [3]. CNV can occur in recombination and replication events. CNV can be benign polymorphisms or associated with Mendelian, sporadic and complex disease possibly through gene dosage, disruption, fusion or positional effects [4].

1.3 APPLICATION OF ARRAY CGH IN THE CLINICAL SETTING

In 2009, array CGH started replacing karyotyping as the method for detection of CNV in NHS diagnostic genetic services. Array CGH offers a higher resolution, less reliance on analyst interpretation/skill and the advantage of a high throughput practical workflow, however few Trusts actually have this test commissioned as a first line test due to the higher cost. Guy's and St Thomas' NHS Trust adopted a patient to patient hybridisation approach, halving the cost of consumables by replacing the reference sample with another patient sample [1].

The product of the increased resolution of arrays is a higher abnormality pick up rate than karyotyping (25% vs 3.7%) [5, 6].

1.4 PATIENT TO PATIENT HYBRIDISATION

As array CGH compares two samples, the use of a reference sample infers any CNV detected is from the patient. Hybridising two patients removes this assumption, producing two challenges:

1. Is a CNV a duplication in one patient, or a deletion in the second patient?

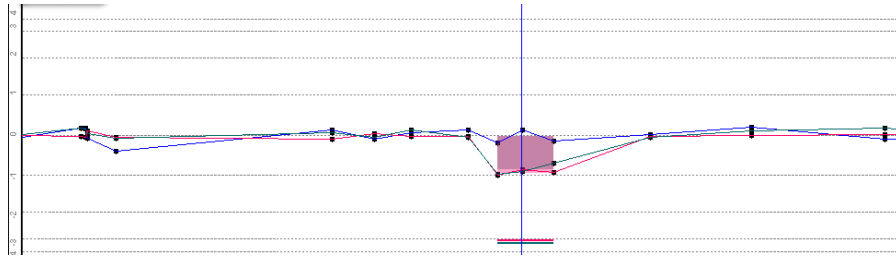


Figure 1: Three array traces are shown. Two arrays (Red and green traces) have the same three probe deletion on chromosome 6. The dark blue trace shows the two arrays hybridised together insilico. The log ratio appears completely normal and the deletion is not detected.

2. If both patients have the same CNV relatively no difference would be seen, resulting in a 'normal' log ratio so the CNV would not be detected.

The first challenge can be overcome by comparing the signal intensities of each dye across the CNV and normal regions. The sample where the signal intensities within the CNV are markedly different from the normal regions contains the CNV.

The second challenge is currently overcome with "a careful consideration of patient referral information" to reduce the risk of this occurring. Hybridisation partners are mismatched on phenotype [1]. This assumes that patients with differing referral reasons eg heart vs renal defects will not have the same underlying CNV. Furthermore, CNV is rare and therefore the risk of two patients with the same CNV being hybridised together is low to start with, even without phenotype mismatching.

1.5 PROJECT RATIONALE

Financial pressures on the NHS are requiring increasingly innovative approaches. Adopting a patient to patient hybridisation approach halves the cost of performing the test [1].

One fear preventing adoption of a patient to patient approach in a clinical setting is the fear of missing a diagnosis due to shared CNV in hybridisation partners [7].

In a clinical setting it is important to have a high test sensitivity and specificity. Any missed diagnoses due to the patient to patient approach would be a false negative, a type two error, reducing the sensitivity of the test. A test must have a high specificity where as a screen can afford to have a higher level of false positives if followed by a test.

A false negative misses a chance of diagnosis, extending the patient's diagnostic odyssey, or may result in a normal report incorrectly being issued for a prenatal test, potentially resulting in an affected child.

The chance of two patients having the same aberration (without phenotype mismatching) has been calculated as 1 in 6000 [8].

Whilst the chance of this occurring is low, and is mediated with strategies such as mismatched phenotypes, with 4000 arrays performed each year this could occur every 1.5 years (without phenotype mismatching).

An array run starts with the creation of a worksheet. This worksheet defines which patients will be hybridisation partners. Hybridisation partners must be phenotype mismatched, a process performed by one clinical scientist and checked by a second. For a run of 96 samples worksheet creation and checking can take three hours. With at least two runs per week nearly one full working day of clinical scientist time is spent creating worksheets. The efficacy can be affected by incorrect, limited or vague phenotype information.

Simplifying the creation of a worksheet will save time, possibly reassigned to a lower grade staff member and may also be automatable, further reducing the cost of processing an array.

More importantly further reducing the chance of missing an CNV improves the service the patient receives, ending, or preventing a diagnostic odyssey.

1.6 PROJECT AIMS

The aim of this project is investigate a tool which is run during data processing and acts as a screen to detect CNV before the array present in both hybridisation partners.

1.7 ARRAY DESIGN

The arrays, reagents, equipment and software used to process and analyse are manufactured by Agilent Technologies [9].

The array is an 8x60K array with a median resolution of 120kb.

The array was custom designed using the online probe catalogue eArray [10]. The array design consists of 46554 probes. 3886 are control probes and 42658 are targeted of which 29 are duplicated leaving 42629 unique probes.

A request for a tool similar to the aims of this project was sent to and declined by Agilent.

1.7.1 *Feature extraction*

After the arrays have been processed in the laboratory and scanned (Figure 2) the array image undergoes feature extraction.

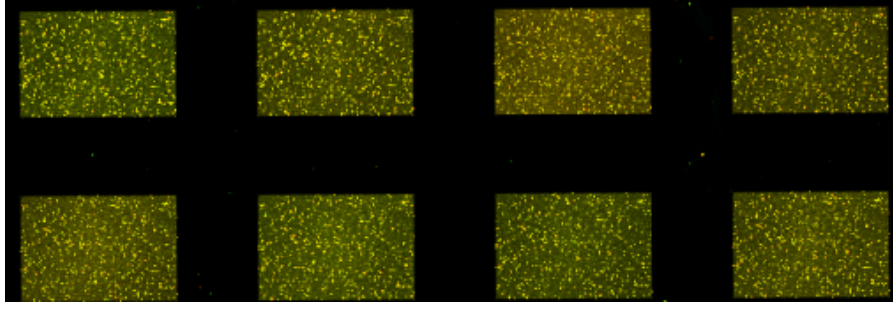


Figure 2: A scanned array slide containing 8 arrays

The feature extraction software converts the raw signal image into measurements for each probe. At each probe the excitation of each dye is captured as the signal intensity. This intensity is processed, including a normalisation step across all probes producing a processed signal intensity normalised to 1000 Arbitrary Fluorescence Units (AFU).

At each probe a log ratio is calculated from the two signal intensities (equation 1)

$$\text{Log}_{10} \frac{\text{Cy5 signal intensity}}{\text{Cy3 signal intensity}} \quad (1)$$

1.7.2 Feature extraction file

The result of feature extraction is one feature extraction file per array, a tab delimited text file of 10-15mb in size.

The feature extraction file can be created in four levels of verbosity: full, compact (default), QC and minimal.

Each feature extraction file consists of three sections, the parameters, stats and features.

1.7.2.1 Parameters

The parameters section contains 38 fields covering information relevant to all 8 arrays on the slide such as the scanner make

and model, time, date and protocols or parameters used to scan the slide.

1.7.2.2 *Stats*

The stats section contains 189 fields covering array specific measurements, such as average intensities and the derivative log ratio spread (DLRS) which is an indicator of quality used during analysis.

1.7.2.3 *Features*

The majority of the file is made up by the features section with one row per probe consisting of 42 fields. 9 fields are identifiers such as the probe catalogue name, the genomic locations and flags to identify probes used as controls. The remaining 33 fields include the raw signal intensity measurement for Cy3 and Cy5 dyes, the log ratio, background signal intensities and flags to mark a probe as having an extreme signal intensity, saturated or non-uniform.

1.8 USE OF ALGORITHMS TO DETECT COPY NUMBER VARIATION

If all goals of this project are met the tool produced will have much the same role as the aberration detection algorithms, taking raw signal intensities and identifying CNV.

There are a number of algorithms which call CNV from CGH data. These algorithms are recursive binary segmentation methods which break chromosomes into segments of equal copy number. Examples include Z scores, Circular Binary Segmentation (CBS), Aberration Detection Method (ADM) 1 and 2, Nexus and Hidden Markov Model (HMM).

1.8.1 *Z score algorithm*

A Z score is a statistical measure of deviation from the mean for a normally distributed population. The Z score algorithm assigns each probe a Z score based on the log ratio. Any probes outside a user defined cut off (R) are classified as an outlier, or significantly away from mean.

The number of probes classified as above (R) and below (R') this threshold and total number of measurements (N) are recorded across small windows of the genome.

These 'windows' can be specified as a number of adjacent measurements or a fixed size eg every 1MB. Within each window the abundance of probes which log ratios which deviate from the mean is measured (r:r').

A Z score is then calculated measuring the significance of the over-abundance of probes with a deviant score in this window [11].

1.8.2 *Aberration Detection Method (ADM-1/ADM-2)*

These algorithms are designed by Agilent[11]. These algorithms do not used fixed windows but segment the genome into intervals of equal copy number using log ratio scores from adjacent probes to best define interval breakpoints.

1.8.2.1 *ADM-1*

Firstly the data is normalised by subtracting the mean log ratio and dividing by the variance, creating a normally distributed population with a mean of 0.

Each chromosome is then broken into intervals of equal copy number and intervals are assigned a score (S(I)) which denotes the difference from the mean. A user defined thresh-

old is set and any intervals which are above this are called as an aberration.

The interval with the highest score is selected and the same process is performed on this segment to further define breakpoints. This is repeated for all intervals with an $S(I)$ above the threshold.

1.8.2.2 ADM-2

The ADM-2 algorithm builds on ADM-1 by including probe quality information to weight probe signal log ratios when assigning $S(I)$.

1.8.3 Circular binary segmentation (CBS)

CBS uses log ratios from adjacent probes to create intervals which, as opposed to ADM-1/2 (which classifies segments as aberrant or not), are grouped into intervals with equal copy number [12].

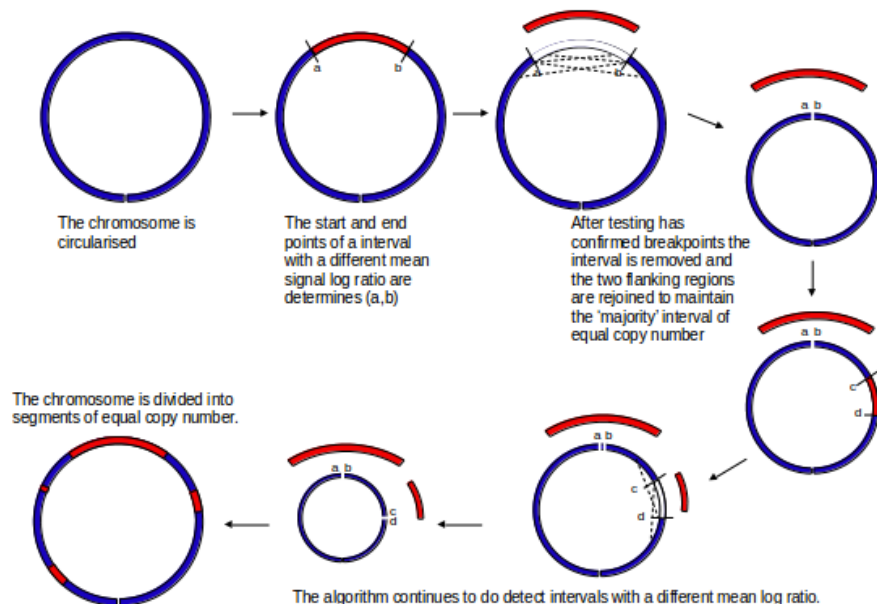


Figure 3: Circular binary segmentation algorithm

Each chromosome is made into a circle (Figure 3). This allows for two break points to be identified, increasing the resolution of detection. When two breakpoints are detected the interval of potential different copy number is removed, and the flanking regions are joined to form a new circle. This allows a t-test to be performed between the mean log ratios on each interval.

Definition of an interval is determined using permutation testing, creating intervals using various breakpoints and looking for the most significant P value. If this P value is above a threshold an interval is created and the process is repeated for all intervals until no more changes are found.

A number of checks are performed to ensure the correct end points have been found including edge effect correction, change point pruning and estimation of the log score distribution.

1.8.4 *Hidden Markov model*

The hidden Markov model (HMM) [11] uses observations (log ratios) to determine the state of a hidden or latent value (copy number). The HMM has a chain property which takes into account the state of the previous probe.

Firstly the data is segmented: the Haar wavelet is used to normalise the data before breakpoints are defined using preset parameters (FDR threshold). This step also calculates probabilities and probability distribution parameters used in later steps.

The HMM is then applied. The forward-backward algorithm is used to calculate the posterior probabilities of the states. Baum-Welch learning uses these to ensure that each state has a well defined value.

Finally the Viterbi algorithm is used to assign a state to each probe.

1.8.5 *Algorithm performance*

Each algorithm performs differently depending on the signal:noise ratio and the size of the aberration[13].

Circular binary segmentation (CBS) is rated as one of the most consistent performers albeit one of the most computationally demanding [14, 13]. Some algorithms' performance varies greatly depending on user defined parameters which may be difficult to establish [14].

The algorithms described above are global segmentation methods which compare segments to the whole genome. Recently local segmentation methods have been described which look for CNV using high resolution data sets [15]. These algorithms have higher power than global algorithms but are less robust so an approach utilising both types of algorithm may produce the best performance [3].

1.9 DATA STORAGE AND MANIPULATION

The tool requires data to be stored, curated and accessed. A database is an efficient tool for this role.

1.9.1 *Relational database*

A relational database is a collection of tables where each row is an instance and each column is an attribute. Tables can be linked to store data in a simple and efficient manner.

Relational databases are stored in a Relational Database Management System (RDBMS) server such as MySQL or MS-SQL which can be installed locally or on a server to allow access

from remote machines. A graphical user interface (GUI) or ‘front end’ enables user interaction with the data, commonly seen in Microsoft Access or an internet browser.

Database design and data manipulation/retrieval utilises the Structured Query Language (SQL).

1.9.2 *Non-relational database*

MongoDB [16] is an alternative database format to relational databases, storing data in BSON format, a binary form of the JSON format. JSON stores data as nested key-value pairs.

MongoDB does not require relationships between tables and fields to be defined, removing the need for complex joins and enables documents with differing formats to be easily stored and compared. MongoDB is faster and easier to scale than relational databases.

1.9.3 *Python*

Python is an object orientated programming language. Python programs can be written to perform a wide range of functions, such as parsing text files, performing complex mathematical and statistical calculations utilising additional modules or packages such as NumPy [17] and SciPy [18].

Connectors are available (including MySQLdb [19], SQLAlchemy [20], PyODBC [21] and PyMongo [22]) enabling Python to interrogate and interact with databases.

1.9.4 *R statistical software package*

The R statistical software package [23] has a repository of packages (Bioconductor [24]) which includes packages which implement CNV algorithms such as CBS within DNACopy [25].

METHODS

2.1 SHARED IMBALANCE DETECTION (SID) TOOL

The Shared Imbalance Detection (SID) tool consists of two components, a database and a python script. The database holds information about the array design, values used during the analysis and tables which are filled during analysis. The python script parses feature extraction files, importing information to the database and performing the Z score analysis with a combination of SQL commands and python functions.

2.2 REQUIREMENTS GATHERING

To ensure the tool is fit for purpose requirements were agreed with key stakeholders in the Array and Bioinformatics teams within the Viapath Genetic Laboratories:

2.2.1 *Functional requirements*

- Use the signal intensities from the feature extraction file to identify CNV independently of the hybridisation partner
- An region of CNV is defined as a minimum of three consecutive abnormal probes which is ended by three consecutive normal probes.
- Report any common or overlapping regions of abnormal copy number present in both hybridisation partners.
- The test is a screen, sensitivity favoured over specificity.

2.2.2 *Non functional requirements*

- In line with the existing departmental resources, experience and skill sets data will be stored in an MS-SQL database hosted by a remote server, ensuring long term support.
- Any scripting should be performed in Python, in line with the existing department skill set to ensure long term support and development and code reviewed as per best practice guidelines.
- Python code must conform with the departmental Python style guide.
- Code must be version controlled using Git.
- Processing should be timely and able to be performed on existing desktop computers.
- The feature extraction files are required for analysis by Genetic Scientists so must not be modified or moved from the current location.
- Compatible with the trust operating system (Windows XP, with an upgrade to Windows 10 planned before 2018).
- Able to be run from multiple computers within the department by multiple users.
- Analysis must be reproducible.
- The database must be in third normalised form to ensure data integrity.
- Standard Operating Protocol (SOP) and validation documentation in line with ISO accreditation requirements.
- Validated in line with departmental and ISO quality procedures.

2.2.3 Project management

- Use an Agile development methodology to regularly test and refine the project (Figure 4).

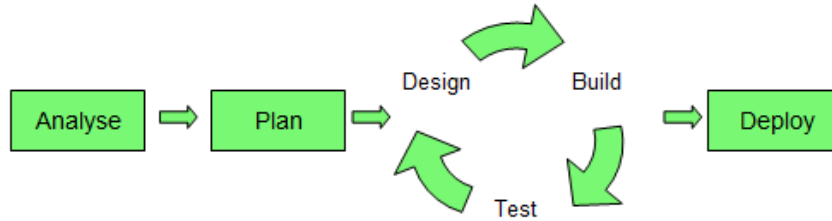


Figure 4: The agile software design methodology promotes regular testing to ensure the design is fit for purpose

- Use smartsheet (<https://www.smartsheet.com>) to set goals and monitor progress with stakeholders
- An electronic lab book is to be kept within the departmental wiki
- Weekly meetings will be held with key stakeholders.

2.3 DATABASE INFRASTRUCTURE

2.3.1 SQL server

Access to the MS-SQL hosted on the department server was not available during development. Installation of an MS-SQL server on the development system was also not possible at the start of the project so development was performed using a MySQL server (v5.6)[26] running locally on a development PC running Windows 10.

2.3.2 Migration from development server (MYSQL) to production server (MS-SQL)

The database was later transferred to a MS-SQL server (SQL-Express 2014 [27]) on the development system using Microsoft SQL Server Migration Assistant (v5.3) for MySQL [28]

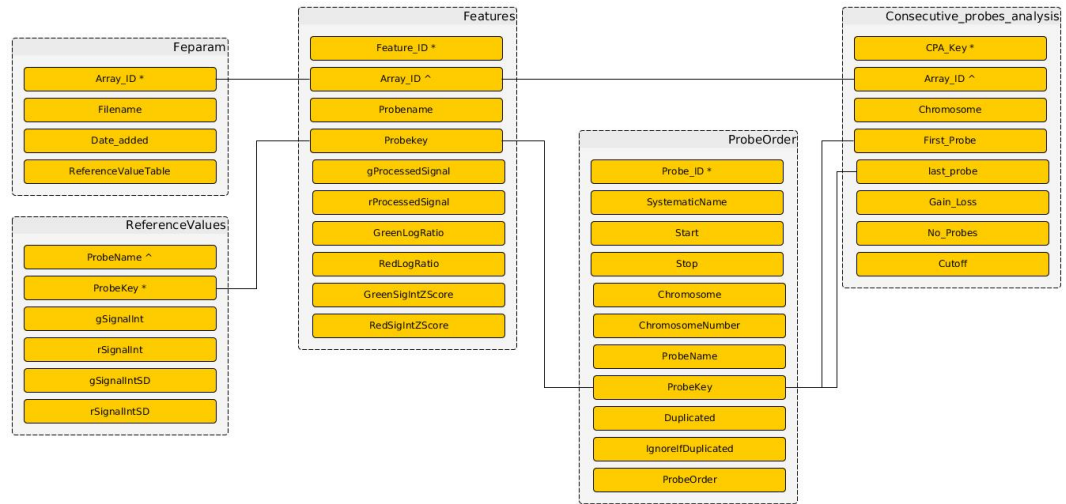


Figure 5: The SID database schema. Five tables are required.^ denotes an indexed field, * denotes a primary key. Lines denote how tables are joined.

The MySQLdb python connector was replaced by PyODBC MSSQL connector [21] and tested to ensure no loss of function.

2.3.3 HeidiSQL client

Database development was performed using the HeidiSQL client (v9.1)[29].

2.4 SID - DATABASE DESIGN

To minimise processing time only essential fields from the feature extraction file is imported to the database. The database is designed in third normalised form, storing numeric keys to prevent anomalies and to maximise the performance of joins. Array wide values are stored a single time in the feparam table, not for each feature.

Five tables are required (Figure 5).The Feparam table holds the filename of the feature extraction file. Upon insertion an auto incrementing number is created as the primary key and used as the unique identifier for that array (Array_ID). The

Feparam table also records the reference table used and an automatic time stamp.

The array design is stored in the Probeorder table. Each probe has a numerical key (probekey) and it's order in genomic position (probeorder). Only non-control probes are stored and replicated probes flagged.

The ReferenceValues table holds the average signal intensity and the standard deviation from the normal control population for each dye at each probe (see 2.7.2).

The Features table holds the analysis for each probe. Of the 42 fields only the probename and processed signals are imported. Remaining fields are populated during the analysis.

Results of analysis are entered into the consecutive_probe_analysis table enabling the aberration to be fully characterised.

2.4.1 *Indexes and keys*

Each table has a primary key which is unique and indexed. Any columns which are joined to the primary key of another table are foreign keys. Any columns used to join table are also indexed to improve performance.

2.5 IMPORTING AND PROCESSING FEATURE EXTRACTION FILES

2.5.1 *Python*

Feature extraction files were parsed using Python (v2.7)[30]. The Anaconda Python distribution (v2.3.0) [31] was selected to make use of the pre-installed packages such as NumPy, SciPy and math packages. The package manager Conda was used to

install the MySQLdb package (mysql-python v1.2.5)[19] to access the MySQL database from python.

The Eclipse Integrated Development Environment (IDE) (Luna v 4.4.1) [32] was used with the PyDev package (v 3.9.0) [33] to develop and run Python scripts.

The Pylint Python package (v 1.4.2) [34] and autopep8 were used within Eclipse to apply the local Python style guide to the code.

2.5.2 *Version control*

Python code was version controlled using Git and backed up using GitHub (<https://github.com/aledj2/FeatureExtraction>).

2.6 SID - PPROCESSING FEATURE EXTRACTION FILES

An overview of the SID python script can be seen in Figure 6 .

2.7 DETECTION OF CNV

2.7.1 *Choice of algorithm (Z score)*

A Z score analysis method was selected for the analysis as the characterisation of breakpoints is not as critical as in routine analysis and this algorithm is simple to implement, maintain and troubleshoot, without the need for any additional software such as R.

A Z score (equation 2) provides a measure of how many standard deviations the signal intensity of a probe is from the mean of a normal population. A threshold can then be applied to score probes as normal or abnormal.

$$Z = \frac{\text{Signal Intensity} - \text{Mean}}{\text{SD}} \quad (2)$$

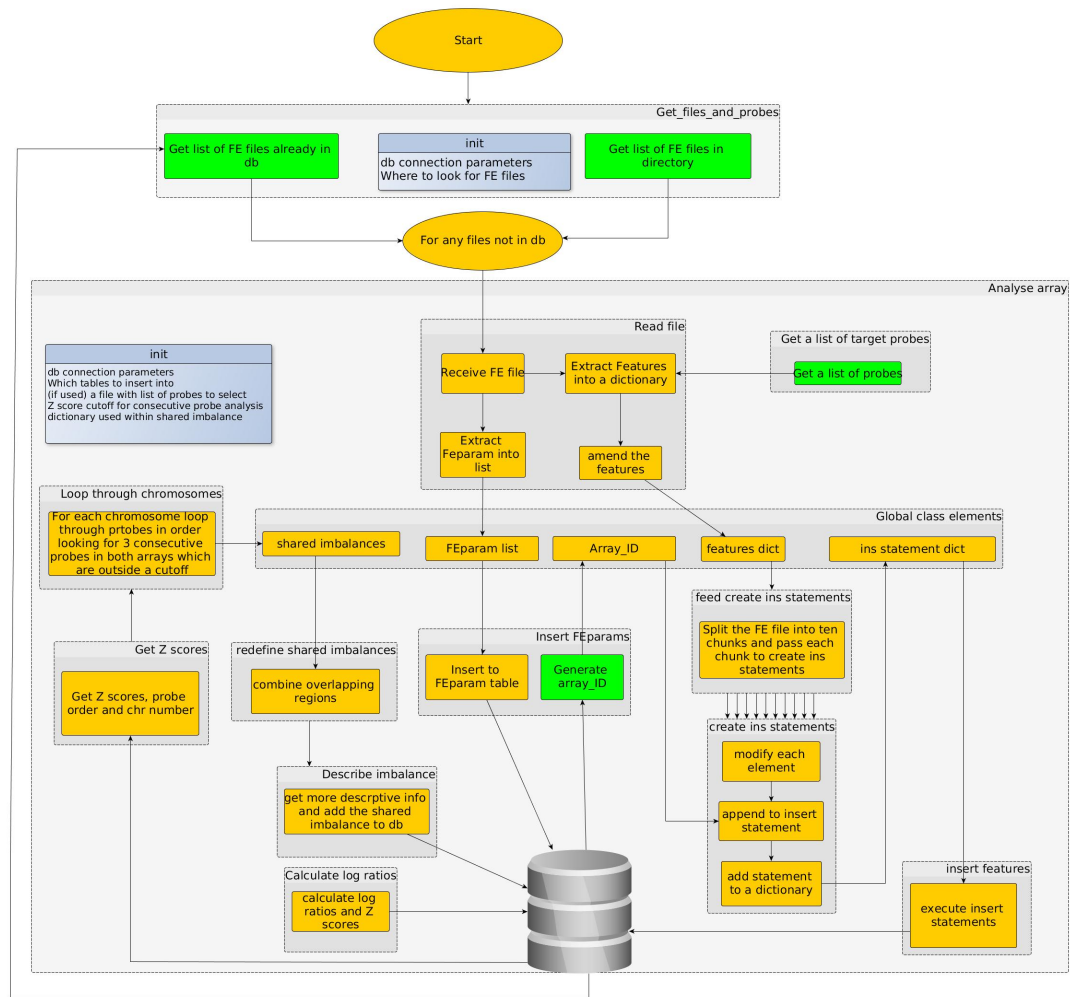


Figure 6: A visual representation of the SID python script. The `Get_files_and_probes` class compares the list of feature extraction files within a directory to the feparam table. Each array yet to be analysed is passed to `Analyse_array` class. The read file module extracts the relevant information into dictionaries. The insert FEparams module enters the filename to the database and captures the primary key assigned (the array_ID). The features section of the feature extraction file is broken into 10 and inserted with the array_ID. The calculate log ratios module applies an update statement to calculate Z scores and log ratios. The get Z scores module takes these Z scores and identifies shared trios of abnormal probes, combining overlapping or neighbouring calls and storing results. Functions coloured green read from the database, yellow process data and blue are `__init__` modules which initialise variables for the class

2.7.2 *Creating a normal reference range*

100 arrays with no reported CNV were selected to create a reference range. Care was taken not to select only high quality arrays as this may result in many false positive calls in poor quality arrays.

The 100 arrays were imported into the database. A Python script calculated the average signal intensity and standard deviation for each dye at each probe. These values are held in the ReferenceValues table which is versioned using the date as the reference may require periodically recalculating.

The ReferenceValues table used is specified within the SID Python script and recorded in the Feparams table, allowing analysis to be repeated after the reference range has been recalculated.

2.7.2.1 *Normal distribution of reference range*

Z scores require normally distributed reference ranges. However some probes will not be normally distributed for example if they are within a common CNV.

A histogram of the signal intensities of all probes from 30 arrays (Figure 7) shows a slightly positively skewed distribution. This it to be expected as the signal intensities cannot be negative and gains in copy number is seen more commonly than losses [4].

2.7.3 *Determining if a probe is normal or abnormal*

The Z score for each dye at each probe is calculated during the import of data and stored within the Features table. The Z score is a measure of how many standard deviations from the mean the signal intensity of the probe is. A threshold is required to define a probe as abnormal.

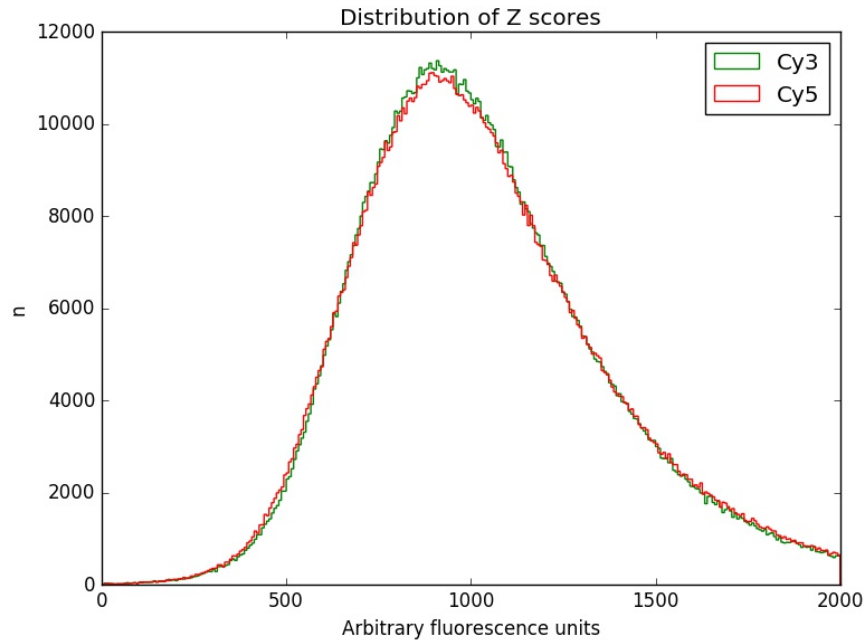


Figure 7: The distribution of signal intensities from all probes from 30 arrays is slightly positively skewed

Z scores are two tailed so probes with a Z score above the threshold are defined as gains and Z scores below the negative threshold are defined as loss (Figure 8).

2.7.4 *Duplicated probes*

For the 29 duplicated probes an average of the signal intensities is taken.

2.7.5 *Determining if a region is normal or abnormal*

To combine abnormal probes into segments a slightly different approach was taken to the Z score algorithm described by Agilent [11].

Aberrations below three probes are not reported so a sliding window of three probes was used, identifying shared CNV of three consecutive probes (Figure 9). The window moved along the entire chromosome creating overlapping tiles (Figure 10).

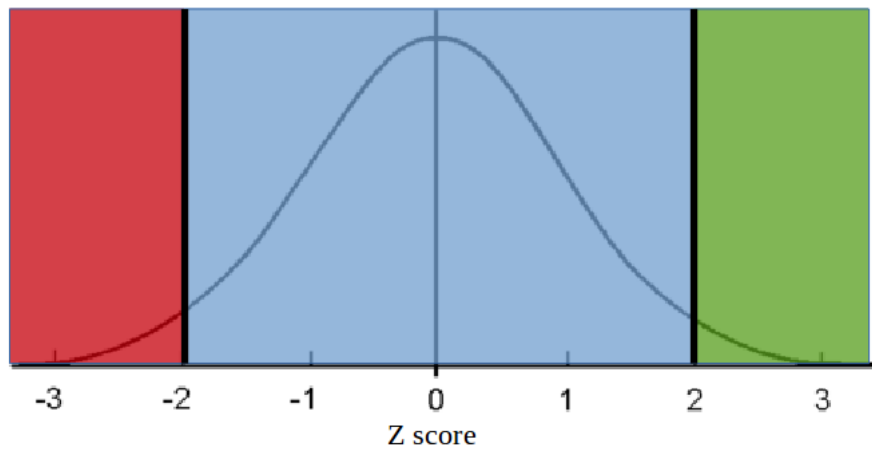


Figure 8: Thresholds define probes as abnormal or normal. A more extreme threshold will define fewer probes as abnormal.

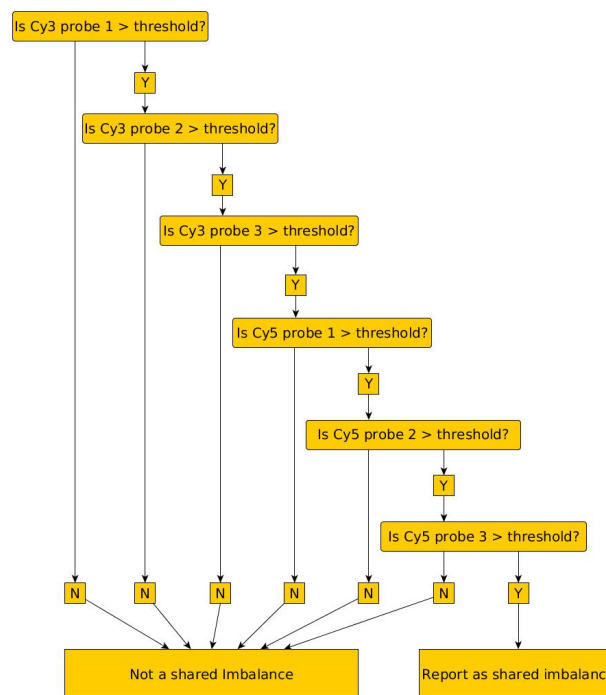


Figure 9: Consecutive probe analysis algorithm. To determine if the sliding window of three probes is abnormal in both samples the Z score of each probe is tested in succession. The threshold is defined in SID python script. This process is repeated assessing if the Z score is less than the negative threshold.

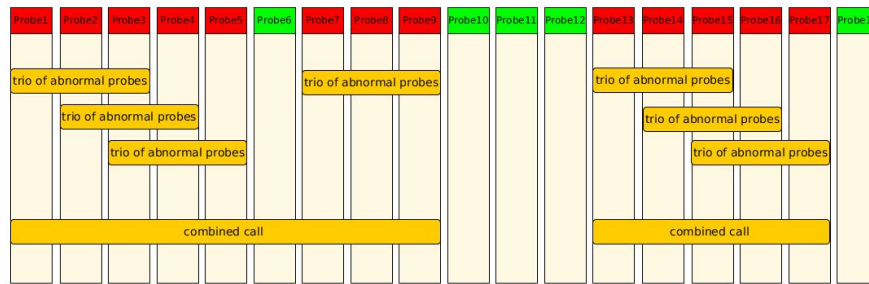


Figure 10: Trios of probes are combined into a single call in line with the analysis strategy. Abnormal probes are represented in red and normal probes represented in green. Trios of shared imbalances are shown. Overlapping and neighbouring calls are combined if they are separated by less than three normal probes (eg probes 10-12).

Trios of shared CNV are then combined in line with routine analysis, combining regions unless separated by 3 or more normal probes (Figure 10).

2.8 CREATION OF TRUE POSITIVE FEATURE EXTRACTION FILES

SID underwent three stages of development.

2.8.1 *Training cases*

A Python script was written to create true positive feature extraction files using feature extraction files with reported aberrations.

Probes within the reported region were identified and a new feature extraction file was created where the Cy5 signal intensity was replaced by the Cy3 signal intensity. This was repeated replacing the Cy3 signal intensity with the Cy5 signal intensity.

The script created one true positive, with signal intensities representing an aberration in both hybridisation partners and one true negative file where both partners have normal signal intensities.

Table 1: Cases with one of the 9 most common CNV were used to create the test set

Aberration	Copies	Number of Probes	Number of arrays
15:22765627-2308509	1	11	9
16:29673953-30198600	1	34	9
16:29673953-30198600	3	34	10
22:18896971-211440514	1	750	8
15:30933080-32514980	1	30	6
16:21806298-22407931	1	15	10
7:72700413-74142327	1	43	6
1:145413387-145747269	3	21	5
6:26440746-26463502	1	3	6

This training set was used to investigate and define the parameters required to call a probe as abnormal.

2.8.2 Test cases

The test set was used to ensure the parameters were not over-fitted to the training set. An updated strategy was applied to create the test set. Two samples which had exactly the same reported CNV were hybridised post-hoc. To ensure the signal intensities would be compared against the correct reference range a sample labelled with Cy3 was paired with one labelled with Cy5.

EvE, a Python script used for the post-hoc hybridisation of arrays in Preimplantation Genetic Diagnosis (PGD) cases, parses two feature extraction files, extracting all measurements for both samples on both arrays, and creates a new feature extraction file where the Cy3 measurements are populated from the Cy3 sample from the first file and the Cy5 measurements from the Cy5 sample from the second file.

The test cases were reported to have one of the most commonly reported CNV, representing a range of sizes and loss and gains (Table 1).

2.8.3 *Prospective cases*

121 unmodified array files used to create the test set were used as a prospective trial to indicate the expected number of calls seen in routine analysis.

2.9 DEFINITION OF TRUE POSITIVE AND FALSE POSITIVE CALLS

A true positive call is defined as any call overlapping with the reported CNV. There may be multiple calls within this region.

A false positive call does not share a single probe with the reported CNV.

2.9.1 *Unknown true positive calls*

In the training and test sets the true positive calls are defined as those which have been reported as clinically significant based on the referral reason however on average an array will have 3-4 regions of CNV [35]. These calls can be false positive calls made by the algorithm or real benign calls in regions of common copy number variation.

Therefore the number of false positive calls will be inflated and the true specificity unknown.

RESULTS

3.1 TRAINING CASES

3.1.1 *Z score threshold*

Investigating the Z score of probes within regions known to be abnormal indicates Z score thresholds that can be applied to mark probes as normal or abnormal (Figure 11).

However, reported regions may contain normal probes and whilst sensitivity is favoured over specificity a balance must be found to prevent many false positive calls.

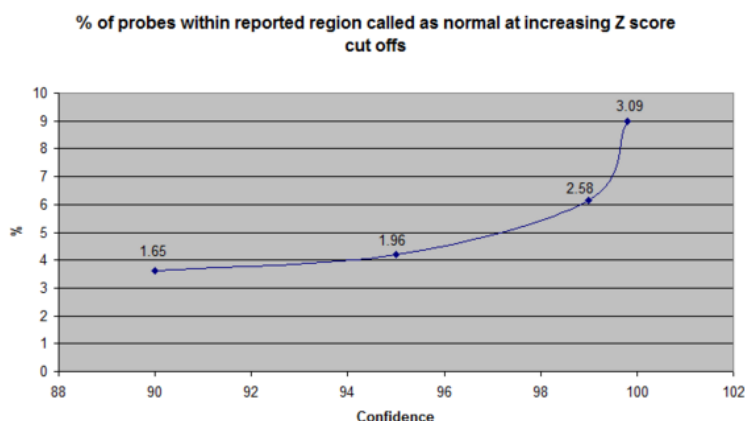


Figure 11: The percentage of probes called as normal within known abnormal regions increases with the Z score threshold. Abnormal regions can include normal probes so even at a Z score threshold of 1.65 over 3% of probes are called as normal. The Z score is shown above each point

3.1.2 Removal of cases with high false positive calls

A range of Z score thresholds were applied, starting at ± 2.374 , defining abnormal probes as the upper 1% and lower 1% of the reference range, and increasing up to ± 4.25 (Table 2).

Table 2: A very high number of false positive calls was seen in the 86 training cases using a range of Z score thresholds.

	Z Score						
	2.374	3	3.5	3.55	3.75	4	4.25
% True Positives (n)	92 (79)	92 (79)	92 (79)	92 (79)	92 (79)	92 (79)	88 (76)
% Arrays with false positives (n)	67 (58)	45 (39)	37 (32)	36 (31)	33 (28)	30 (26)	30 (26)
Number of false positive calls	10228	7681	6177	6044	5538	4944	4394
Average number of false positive calls per array (max)	192.98 (2128)	225.91 (1733)	228.78 (1468)	232.46 (1443)	240.78 (1351)	235.43 (1230)	209.24 (1108)

At the lowest Z score threshold (2.374) the false negative rate was 8% with 67% of cases containing a false positive call.

The same Z score threshold resulted in 10228 false positive calls, of which 94% came from just five arrays (range 1612-2128). The average number of false positive calls in the other 81 arrays was 10.8 (max 104) (Table 3).

Visual inspection of the array traces and the raw data from these 5 arrays gave no indication why so many false positives calls were made or how this could be predicted. Therefore these five arrays were considered a different population and excluded from further analysis.

3.1.3 Comparison of true positive and false positive calls

The remaining 81 arrays show that at least one call was detected within the abnormal region in 94% of cases with a Z score up to 4 (Table 3).

Using a Z score of 2.374 nearly 2 in 3 arrays had a false positive call. A Z score threshold of 4 reduced the frequency of false positive calls to 1 in 4 arrays, with an average 1.66 and maximum 6 false positive calls per array.

Table 3: The calls made in 81 training cases using a range of Z score thresholds. The true positive call was missed in at least 5 cases across all thresholds, with a false positive call made in at least 1 in 4 cases. The number of false positive calls, and the average number of calls made decreased as the Z score threshold increased.

	Z Score						
	2.374	3	3.5	3.55	3.75	4	4.25
% True Positives (n)	94 (76)	94 (76)	94 (76)	94 (76)	94 (76)	94 (76)	90 (73)
% Arrays with false positives (n)	65 (53)	42 (34)	33 (27)	32 (26)	28 (23)	26 (21)	26 (21)
Number of false positive calls	571	114	50	48	41	35	31
Average number of false positive calls per array (max)	10.8 (104)	3.35 (20)	1.85 (7)	1.84 (7)	1.78 (7)	1.66 (6)	1.48 (5)

3.1.3.1 Number of probes within a call

The length of false positive calls are shorter than true positive calls. Increasing the minimum number of probes within a call would remove these false positive calls but would also remove true positive calls and fail to meet the functional requirement that a call is a minimum of 3 consecutive probes (Figure 12).

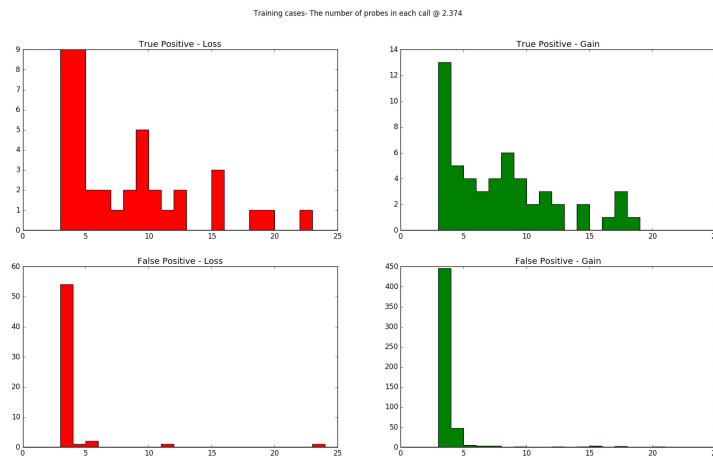


Figure 12: Training cases: False positive calls are smaller than true positive calls.

3.1.3.2 *Difference between true and false positive Z scores using a threshold of 2.374*

Alternatively increasing the Z score threshold to above the lowest scoring probe in abnormal regions may prevent the region being called (Figure 9).

Examining the minimum (Figure 13) and average (Figure 14) Z scores for each true and false positive calls show that the two populations are not disparate therefore increasing the Z score threshold may also remove true positives.

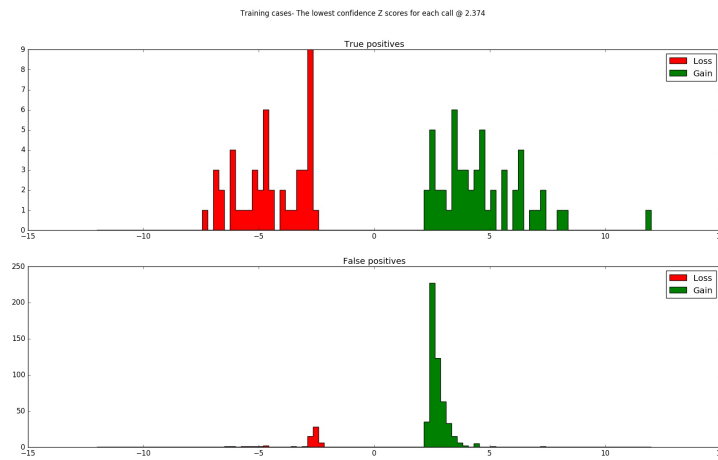


Figure 13: Training cases: The lowest confidence Z score within each call at a threshold of 2.374. There is an overlap between the confidence of the lowest probe within true and false positive calls

However the increased number of probes within a true positive region (Figure 12) and the higher average Z score (Figure 14) may mean the true positive regions remain.

3.1.3.3 *Difference between true and false positive Z scores using a threshold of 3.55*

Increasing the Z score threshold to 3.55 showed that the average and lowest confidence Z scores in remaining false positive calls were similar to that of true positives (Figures 15 and 16).

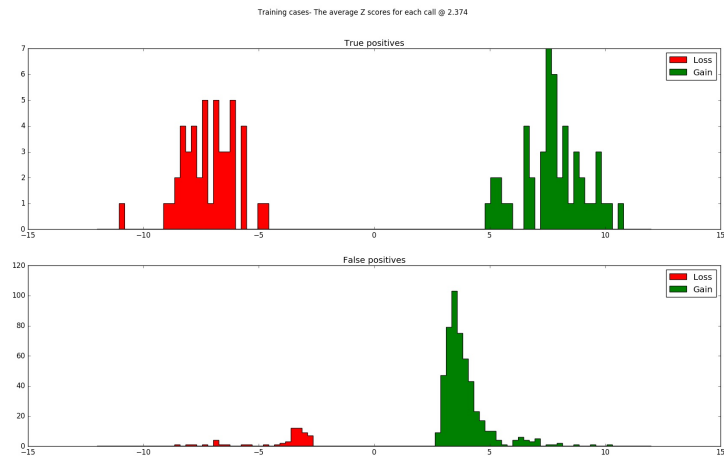


Figure 14: Training cases: The average Z score of probes within calls at a threshold of 2.374. The average Z score of true positive calls are further from the mean than the average Z score of false positive calls.

3.1.3.4 Suitability of reference range

There were no recurring false positives calls to suggest the reference range is unsuitable.

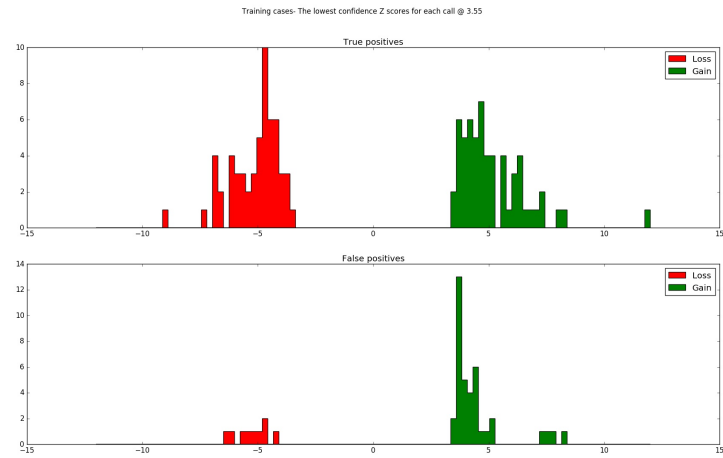


Figure 15: Training cases: At a threshold of 3.55 the lowest confidence probes within true and false calls are very similar

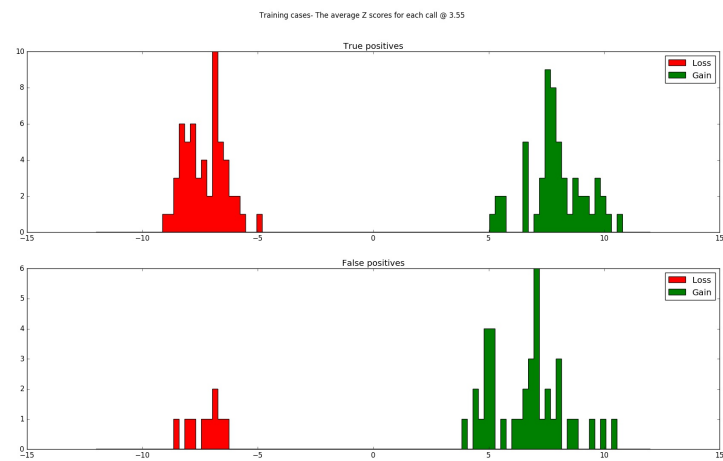


Figure 16: Training cases: The average Z score for probes within true and false positive calls with a threshold of 3.55 are very similar.

3.2 TEST CASES

The training cases are not truly representative of an array where both hybridisation partners have the same CNV. Cy3 and Cy5 dyes are incorporated with varying efficacy so the reference value will vary. When calculating the Z score for the training set comparing the signal intensity of one dye to the reference value of the other dye may introduce error.

The same thresholds were applied to 69 test cases (Table 4).

3.2.1 *A greater sensitivity was seen using more representative feature extraction files*

There was 100% sensitivity using a Z score below 4 and there were no false positive calls using a Z score above 3.55 (Table 4).

Table 4: Test cases: A call was made in the expected CNV region in all cases with a threshold below 4. There were no false positive calls made with a threshold above 3.55

	Z score						
	2.375	3	3.5	3.55	3.75	4	4.25
% True Positive arrays (n)	100 (69)	100 (69)	100 (69)	100 (69)	100 (69)	100 (69)	97 (67)
% Arrays with false positives (n)	13 (9)	8.7 (6)	1.4 (1)	0	0	0	0
Number of false positive calls (max)	48 (28)	9 (3)	1 (1)	0	0	0	0

3.2.2 *Suitability of reference range*

There were no recurring false positive calls.

3.2.3 *Proportion of CNV regions covered by calls*

As shown in Figure 17 with the exception of 8 arrays with the 750 probe deletion on chromosome 22 (Table 1) at least 25% and on average 66% of the probes within the known CNV were within a call. Cases with the chromosome 22 deletion is described in section A.1

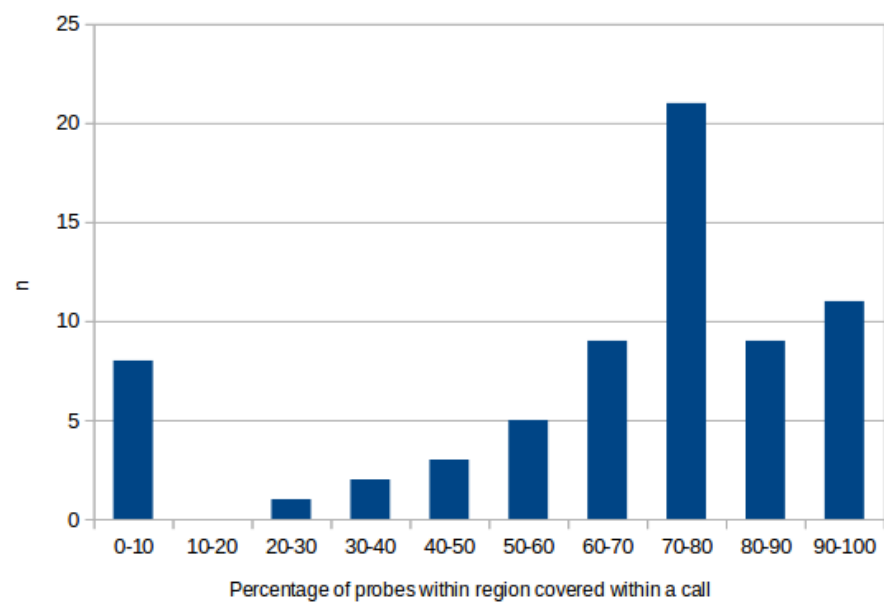


Figure 17: The proportion of probes within reported CNV included by a call at a threshold of 3.55.

3.3 PROSPECTIVE CASES

121 routine arrays were processed. The expected results are unknown but the number of calls is indicative of the number of calls required to be assessed during analysis.

Using a Z score threshold of 2.374 there were a significant number of calls, with more than 1 in 3 cases having calls to investigate. A threshold of 3.55 resulted in a call in at least 1 in 9 arrays (Table 5).

The maximum number of calls in a single array is quite high, with two arrays with more than 50 calls at a threshold of 3.55 (Figure 18), reminiscent of the 5 arrays excluded from the training set

It is unknown if any of these calls are actually true positives.

The average size of each call is just over three probes (Table 6)

Table 5: Prospective Cases: The proportion of the 121 arrays where a call has been made, and the total number of calls made across a range of thresholds. At a threshold of 3.55 12% of arrays contain a call. Across these 14 cases there are 246 calls, of which 119 is made in a single case.

Z score	% of arrays with a call (n)	Total number of calls (max calls in a single array)
2.374	37 (45)	1760 (555)
3	16 (19)	591 (219)
3.55	12 (14)	246 (119)
3.75	12 (14)	186 (98)
4	10 (12)	125 (66)
4.25	7 (9)	90 (50)

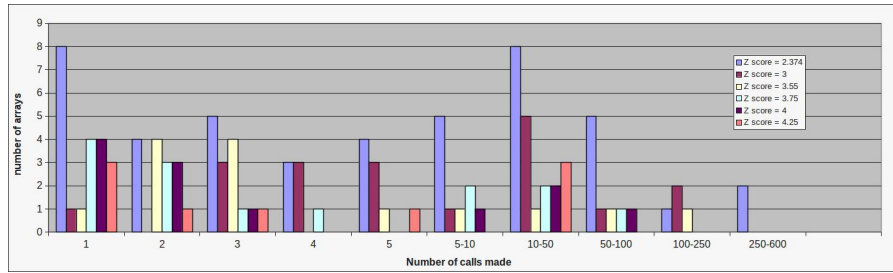


Figure 18: Prospective cases: The number of calls per array across a range of threshold in 121 prospective cases

Table 6: Prospective cases: The average and maximum number of probes within a call during analysis of 121 prospective arrays.

Z score	Average probes in call	Maximum probes in call	Number of calls
2.374	3.24	8	1760
3	3.19	7	591
3.55	3.14	6	246
3.75	3.13	6	180
4	3.1	5	125
4.25	3.09	5	90

3.4 PERFORMANCE OF SID

When uploading a run of 96 arrays, initially each array is processed in approximately two minutes. However, due to the increasingly large features table, the time taken to process an array can increase to around 15 minutes.

DISCUSSION

4.1 SUITABILITY OF SID

The Shared Imbalance Detection (SID) tool, written in Python and utilising a MS-SQL database can scan a directory, parsing feature extraction files and by applying a Z score algorithm can detect regions of shared CNV from manually created true positive cases with a sensitivity of 100% in a timely fashion.

4.2 Z SCORE THRESHOLD TO DEFINE A PROBE AS ABNORMAL

Analysis of the training cases identified suitable thresholds to be applied however the training cases were not truly representative of routine feature extraction files.

Analysis of the test cases suggest that using a threshold of between 3.55 and 4 to define a probe as abnormal would detect all shared copy number variation without any false positive calls.

Applying the lower threshold to prospective arrays resulted in calls in 1 in 9 cases. However some cases contained a large number of calls, which would increase the analysis burden. Further investigation into why some arrays contain a much higher number of calls and a strategy to predict and manage these arrays may be required.

4.3 THE NUMBER OF PROBES WITHIN TRUE POSITIVE CALLS

The proportion of probes called as abnormal within the true positive calls was not assessed in detail as each call serves as a flag to prompt further investigation by the analyst who would further describe the exact shared CNV.

Whilst the entire CNV may not be called, the calls may be spread out across the CNV, giving an indication of the shared CNV.

Why so little of the large chromosome 22 deletion was not called should be investigated further.

4.4 COMBINING OVERLAPPING TRIOS OF ABNORMAL PROBES

A strategy has been created to combine trios of shared abnormal probes in line with existing analysis protocols.

4.5 BENIGN CNV

It is expected that hybridisation partners would share one or more of the many known, common benign CNV regions. SID would correctly identify these regions as abnormal, however the above analysis would classify these calls as false positives. It may be worth excluding these regions from future specificity and sensitivity calculations

Should SID be adopted by the clinical service stakeholders should be consulted to determine if calls within known benign CNV regions should be filtered out, or if the experience of clinical scientists would make analysis of common calls trivial.

4.6 CALLS SMALLER THAN 5 PROBES

The vast majority of false positive calls in the test set and prospective arrays contain less than 5 probes (Table 6 and Figure 22).

The average signal intensity of probes within a true positive call is higher than that of a false positive call (Figure 20) so a secondary analysis could be applied to further reduce false positive calls.

4.7 DIFFERENT THRESHOLDS FOR GAIN OR LOSS

The proportional signal intensity change between diploid and triploid and diploid and monosomy is the same ($\pm 50\%$). Therefore the Z score would be the same for monosomy and trisomy.

Similar numbers of false positive gains and losses suggest using the same threshold for losses and gains does not introduce bias however should further work identify a lack of sensitivity for either gain or loss separate thresholds could be applied.

4.8 ALTERNATIVE APPROACHES TO CONSECUTIVE PROBES ANALYSIS

4.8.1 *Larger sliding windows*

The Z score algorithm implemented differs to that suggested by Agilent who suggest a sliding window and the proportion of probes outside a threshold used to define a region as abnormal [11].

This approach may result in a call which does not meet the specification (Figure 19a).

The consecutive probe approach applied within SID is designed to match the routine analysis approach however a sliding window approach would prevent aberrations such as those in Figure 19b being missed.

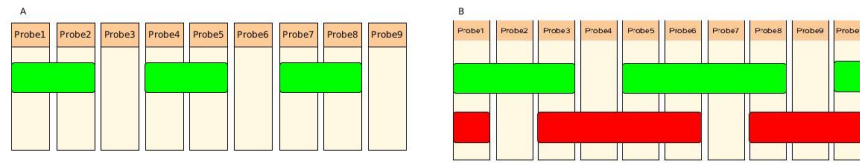


Figure 19: Classifying overlapping shared CNV. a. Using a sliding window approach and measuring the proportion of probes outside a threshold may result in calls which do not meet the functional specification of three consecutive probes. b. The existing algorithm may miss clearly overlapping CNV where three consecutive abnormal probes are not shared

4.8.2 Sliding window - Regions of interest

A variation of the sliding window approach was investigated, using defined regions of interest (every CNV which has been reported more than twice) as a window. The same Z score analysis was performed and the proportion of abnormal probes within the region used to determine if a region is abnormal or not.

Despite this approach being targeted to known clinically relevant CNV and the potential to expand the list of regions this approach is more limited than the consecutive probe approach so was not pursued.

4.9 FURTHER WORK

4.9.1 *Sex chromosomes*

Further work is required to include sex chromosomes in the analysis. The gender of the sample affects the signal intensities of probes on the sex chromosomes.

To identify CNV on the sex chromosomes a gender specific reference range is required. The average signal intensity across non pseudoautosomal regions of the X and Y chromosomes could determine gender and therefore which reference range should be applied.

4.9.2 *Sample type*

This project has only assessed postnatal blood samples. The array service also processes prenatal, solid tissue and PGD samples, which have variable resolution and performance due to the quality of the sample.

Differing analysis strategies and Z score thresholds would be required for each sample type in line with existing analysis approaches in addition to a method for specifying the sample type of each array.

4.9.3 *What happens when shared CNV is detected?*

Once shared CNV is detected further investigation is required. The level of investigation required would dictate the acceptable levels of specificity of SID. If a thorough investigation is required the time and cost incurred would negate any savings made by replacing phenotype mismatching with SID.

The signal intensities can be examined within the analysis software which would confirm or deny the presence of CNV in

line with the current approach to identify if an call is a duplication in one hybridisation partner or a deletion in the other. However having more than 50 calls on a single array would make this approach unfeasible.

Another potential approach is EvE, a tool which allows post hoc hybridisation of samples. Each sample could be independently analysed against another sample, such as a reference control, to see if CNV is present.

EvE takes about one minute to run and the analysis could be performed in less than 10 minutes. However should this be required routinely any time savings would be diminished. EvE has been validated for use with the PGD service but has not yet been validated for use with postnatal blood samples.

An extreme approach is to repeat the array, pairing each sample with a different hybridisation partner. If this approach is taken the specificity must be well understood to avoid a large number of samples being repeated, removing any cost savings.

The approach taken would require a thorough consultation with clinical scientists and bioinformaticians,

4.9.4 *Exporting results to LIMS*

The departmental Laboratory Information Management System (LIMS) system is very widely used, managing all patient information, analysis and results. The SID database will be a separate database hosted on the same MS-SQL server.

A function could therefore be added to the end of the SID python script to export the results directly to the LIMS system. The filename can be used to link to the two patient records using an existing query.

The LIMS system would require an extra table to hold the shared CNV. Clinical scientists should be consulted to determine where and how this information should be displayed.

Confirmation that analysis has completed successfully should be included for samples with no shared CNV.

4.9.5 *Initiation of SID*

Array runs occur at scheduled times during the week. SID scans a directory and compares the files in the directory to those already processed (stored in the feparam table) which would enable the script to be scheduled to run at specific times.

Alternatively the script could be set off manually by the array technician once the feature extraction process has completed. An SOP and training would be required.

4.9.6 *Migration of database to the server*

The database has not yet been migrated to the departmental MS-SQL server. Interacting across a network into a remote server may increase the time taken to process each sample and the stability of this connection tested.

An increase in processing time may be tolerable however a process to identify and remedy interrupted or failed processing will be required.

4.10 LONG TERM MANAGEMENT OF SID

4.10.1 *Unit test*

After any changes are made to either the python script or database a unit test, utilising a truth set from a subset of the test set, should be performed to ensure no loss of function.

4.10.2 *Code review*

The code and database schema should be version controlled using git and backed up remotely using github.

Any changes made to the SID python script should take place on the development branch and code reviewed in line with the departmental standard operating protocol (and best practice guidelines (unpublished)) before changes are adopted into the live code.

4.10.3 *Recalculation of the reference range*

The measured signal intensity can be affected following any maintenance or change to the array scanner. Following any maintenance a new reference range should be calculated using the first array run post maintenance.

The reference range should be created in a new table and the SID python script updated to apply and record this in the feparam table.

4.10.4 *Truncating the database to maintain performance*

Each array run consists of 48 or 96 arrays. As the run is processed and the Features table grows in size the time taken to process each array increases up to 7 fold.

Approaches to control SQL transactions failed to improve this however should this processing time become prohibitive alternative insert approaches can be considered, such as creating and inserting the features from a comma separated values (csv) file. This approach proved slower when uploading single arrays.

Truncation of the Features table, following analysis of each array or array run would maintain performance.

Alternatively a features table could be created for each run, however a table holding 96 samples would be around 800MB which, with at least two runs per week, would quickly become prohibitively large. Given the relative short processing time, high throughput of arrays and long term availability of feature extraction files if a re-analysis is required it would be sensible to re-run SID than to return to stored values.

The probeorder, feparam and ReferenceValues tables should not be truncated. The feparam table holds the filename used by the python script to prevent feature extraction files being reprocessed. Archiving ReferenceValues table ensures reproducible analysis after a new reference range has been created while the probeorder table preserves the array design.

The consecutive_probe_analysis table could be truncated once the shared CNV has been exported to the LIMS.

4.10.5 *Contingency plan*

Should SID be implemented and then become unavailable contingency plans include temporarily suspending the array service or reverting to the current strategy of manually phenotype mismatching hybridisation partners.

4.10.6 *Database backup*

SID will be incorporated to Trust IT's existing database backup strategy, utilising a third party software to perform daily backups, keeping the last three backups.

4.10.7 *Validation*

ISO laboratory quality accreditation (ISO15189) requires thorough validation of all tests and equipment. To validate SID 100 arrays can be analysed and 3 common CNV regions assessed. The predictions made by SID can be compared to the signal intensities inspected manually on the array trace. The workload of manually assessing 300 regions is in line with validation of other tests within the department.

4.10.8 *Documentation*

Three documents are required to accompany SID:

ISO accreditation requires a validation document which details the steps taken to confirm SID performs as expected.

A SOP or user manual describes how SID is run and explaining some common error messages.

A development manual would be aimed at members of the bioinformatics department who will be responsible for developing and maintaining the script and database. This will describe in depth how SID works, dependencies and a change log.

The development manual would also contain instructions to perform tasks which the user cannot perform, for example completely removing arrays to enable the array to be reprocessed and creation and adoption of new reference ranges.

These documents would be stored in the controlled document management system (Q-Pulse) and on the laboratory wiki which hosts the current standard operating procedures for the array service.

4.11 FUTURE CONSIDERATIONS

4.11.1 *Array redesign and multiple array designs*

The department currently has three different array designs: the routine design used in this project, a high resolution array used only in specific cases and a legacy design. To completely avoid phenotype mismatching SID must be able to analyse all array designs.

Each array design would require a probeorder table and a reference range. The array design is contained within the filename. SID would need updating to extract this to apply the correct tables.

Similarly any future changes to the routine array design would require the same measures.

CONCLUSION

In line with the requirements set out in chapter 2.2 SID is capable of processing feature extraction files, and by analysing each hybridisation partner independently identify CNV present in both hybridisation partners.

Using a threshold of between 3.55 and 4 a call was present within the known CNV in 100% of the 69 manually created true positive test cases.

A number of calls outside these regions were also identified which may be false positive calls or true shared CNV within a common CNV region. An indication of the expected prevalence of calls was shown in 121 prospective arrays.

However before SID can be implemented within the clinical service further work is required. These steps have been discussed, including a strategy for analysing the sex chromosomes and migrating SID to the departmental server. Both SID and the departmental LIMS system must be modified to enable shared CNV to be reported to the clinical scientist.

The long term management of SID was described.

APPENDIX

A.1 ANALYSIS OF CALLS MADE IN SAMPLES WITH THE 750 PROBE DELETION ON CHROMOSOME 22

At a threshold of 3.55 all calls from samples with the chromosome 22 deletion were made within the 73 probe region (Table 7).

Table 7: Test cases: At a threshold of 3.55 all calls made on the arrays made from the cases with the 750 probe deletion on chromosome 22 were within a 73 probe section (probes 39733 to 39805). Of these 73 probes between 17% and 74% were covered by a call.

Array_ID	first called probe	last called probe	% of 73 probes covered (n)
63	39735	39805	72.6 (53)
64	39733	39797	64.4 (47)
65	39735	39803	52.1 (38)
66	39735	39803	74 (54)
67	39745	39784	23.3 (17)
68	39745	39784	23.3 (17)
69	39745	39784	23.3 (17)
70	39743	39803	17.8 (13)

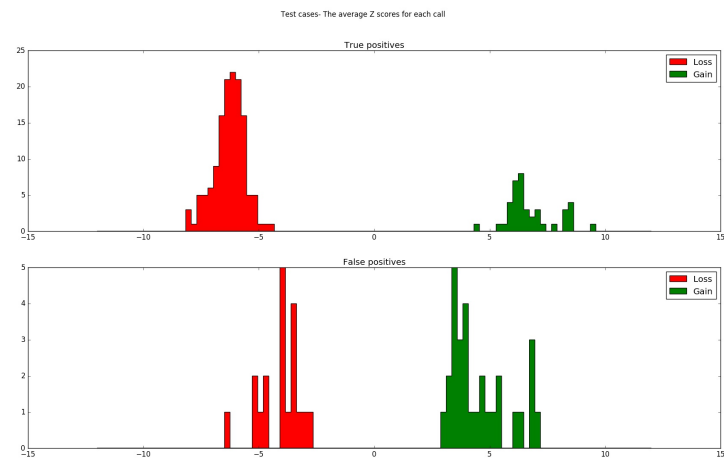


Figure 20: Test cases: The average Z score in true positive calls made at a threshold of 2.374 are more significant in true positive calls than false positive calls however there is overlap between the two groups

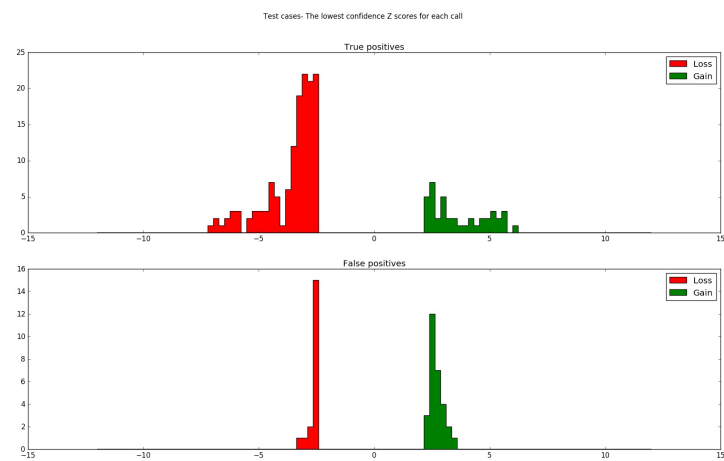


Figure 21: Test cases: There is overlap between the lowest confidence probe within true and false positive calls made at a threshold of 2.374

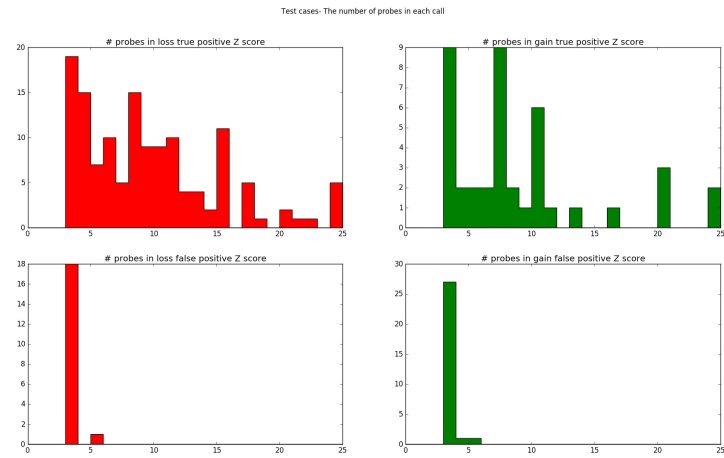


Figure 22: Test cases: At a threshold 2.374 true positive calls contain more probes than false positive calls.

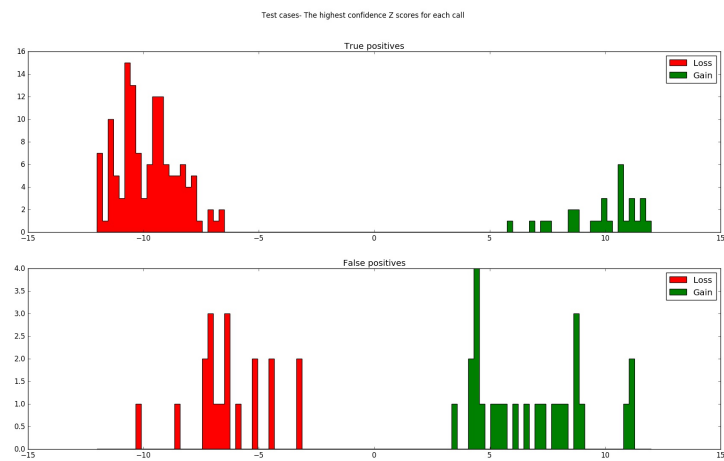


Figure 23: Test cases: At a Z score threshold of 2.374 there is an overlap between the highest confidence probe within true and false positive calls.

BIBLIOGRAPHY

- [1] Joo Wook Ahn, Kathy Mann, Sally Walsh, Marwa Shehab, Sarah Hoang, Zoe Docherty, Shehla Mohammed, and Caroline Mackie Ogilvie. Validation and implementation of array comparative genomic hybridisation as a first line test in place of postnatal karyotyping for genome imbalance. *Molecular Cytogenetics*, 3:9, April 2010. PMID: 20398301 PMCID: PMC2885406.
- [2] Phillipe Hupe. GLAD-Introduction. <http://bioinfo-out.curie.fr/projects/glad/node2.html>, November 2004.
- [3] Siddharth Roy and Alison Motsinger Reif. Evaluation of calling algorithms for array-CGH. *Frontiers in Genetics*, 4, October 2013. PMID: 24298279 PMCID: PMC3829466.
- [4] Feng Zhang, Wenli Gu, Matthew E. Hurles, and James R. Lupski. Copy number variation in human health, disease, and evolution. *Annual Review of Genomics and Human Genetics*, 10(1):451–481, 2009. PMID: 19715442.
- [5] Joo Wook Ahn, Susan Bint, Anne Bergbaum, Kathy Mann, Richard P Hall, and Caroline Mackie Ogilvie. Array CGH as a first line diagnostic test in place of karyotyping for postnatal referrals - results from four years' clinical application for over 8,700 patients. *Molecular Cytogenetics*, 6:16, April 2013. PMID: 23560982 PMCID: PMC3632487.
- [6] Bertrand Jordan,. *Microarrays in diagnostics and biomarker development: current and future applications*. Springer, Heidelberg ; New York, 2012.
- [7] John Dunlop and Jake Miller. Why haven't you adopted a patient to patient approach? Personal communication, July 2015.

- [8] JooWook Ahn. The frequency of abnormalities suggest a 1 in 6000 chance of two patients randomly having the same imbalance, Personal Communication, June 2015.
- [9] Agilent Technologies. CGH & CGH+SNP microarrays. <http://www.genomics.agilent.com/en/CGH-CGH-SNP-Microarrays/>, accessed 2015-01-07.
- [10] Agilent Technologies. eArray. <https://earray.chem.agilent.com/earray/>, accessed 2015-01-07.
- [11] Agilent Technologies. Agilent CytoGenomics 2.0 Reference Guide, 2011.
- [12] E. S. Venkatraman and Adam B. Olshen. A faster circular binary segmentation algorithm for the analysis of array CGH data. *Bioinformatics (Oxford, England)*, 23(6):657–663, March 2007. PMID: 17234643.
- [13] Hanni Willenbrock and Jane Fridlyand. A comparison study: applying segmentation to array CGH data for downstream analyses. *Bioinformatics*, 21(22):4084–4091, November 2005. PMID: 16159913.
- [14] Weil R. Lai, Mark D. Johnson, Raju Kucherlapati, and Peter J. Park. Comparative analysis of algorithms for identifying amplifications and deletions in array CGH data. *Bioinformatics (Oxford, England)*, 21(19):3763, October 2005. PMID: 16081473 PMCID: PMC2819184.
- [15] Yue S. Niu and Heping Zhang. THE SCREENING AND RANKING ALGORITHM TO DETECT DNACOPY NUMBER VARIATIONS. *The annals of applied statistics*, 6(3):1306–1326, September 2012.
- [16] MongoDB, Inc. MongoDB. <https://www.mongodb.org/>, accessed 24/02/2016.
- [17] S. van der Walt, S.C. Colbert, and G. Varoquaux. The NumPy array: A structure for efficient numerical computation. *Computing in Science Engineering*, 13(2):22–30, March 2011.

- [18] SciPy. www.Scipy.org, accessed 15/01/2015
- [19] Andy Dustman. MySQLdb. <https://github.com/farcepest/MySQLdb1>, January 2014, accessed 15/1/2015.
- [20] Michael Bayer. SQLAlchemy - The Database Toolkit for Python, <http://www.sqlalchemy.org/>, accessed 15/01/2015.
- [21] pyodbc. <https://github.com/mkleehammer/pyodbc>, accessed 25/03/2016.
- [22] PyMongo. Python Driver (PyMongo) Getting Started With MongoDB 3.0.4, <https://docs.mongodb.org/getting-started/python/client/>, accessed 24/02/2016.
- [23] R: A language and environment for statistical computing, version 3.1.2, R Foundation for Statistical Computing, <http://www.R-project.org/>, accessed 15/01/2015.
- [24] Wolfgang Huber, Vincent J Carey, Robert Gentleman, Simon Anders, Marc Carlson, Benilton S Carvalho, Hector Corrada Bravo, Sean Davis, Laurent Gatto, Thomas Girke, Raphael Gottardo, Florian Hahne, Kasper D Hansen, Rafael A Irizarry, Michael Lawrence, Michael I Love, James MacDonald, Valerie Obenchain, Andrzej K Oles, Herve Pages, Alejandro Reyes, Paul Shannon, Gordon K Smyth, Dan Tenenbaum, Levi Waldron, and Martin Morgan. Orchestrating high-throughput genomic analysis with Bioconductor. *Nature Methods*, 12(2):115–121, January 2015.
- [25] Venkatraman E. Sechan and Adam Olshen. Dnacopy : Dna copy number data analysis, <http://master.bioconductor.org/packages/release/bioc/html/DNACopy.html>, accessed 10/01/2015.
- [26] MySQL. MySQL Community Server version 5.6, <https://dev.mysql.com/downloads/mysql/5.6.html#downloads>, accessed 11/04/2016.

- [27] SQL Server Express Edition, Microsoft, <https://www.microsoft.com/en-gb/server-cloud/products/sql-server-editions/sql-server-express.aspx>, accessed 22/03/2016.
- [28] Microsoft SQL Server Migration Assistant v5.3 for MySQL, Microsoft, <https://www.microsoft.com/en-us/download/details.aspx?id=42657>, accessed 22/03/2016.
- [29] HeidiSQL. HeidiSQL - MySQL, MSSQL and PostgreSQL made easy, <http://www.heidisql.com/>, accessed 09/04/2016.
- [30] Python Software Foundation. Python Language Reference, version 2.7, July 2010 <http://www.python.org>, accessed 19/01/2015.
- [31] Continuum. Why Anaconda?, <https://www.continuum.io/why-anaconda>, accessed 10/04/2016.
- [32] The Eclipse Foundation. Eclipse IDE <https://www.eclipse.org>, accessed 11/04/2016.
- [33] PyDev. Python IDE for Eclipse, <http://www.pydev.org/>, accessed 11/04/2016.
- [34] Pylint. Code analysis for Python, <https://www.pylint.org/>, accessed 10/04/2016.
- [35] JooWook Ahn. The average number of calls for each patient during analysis of arrayCGH cases., Personal communication Feb 2016.

COLOPHON

This document was typeset using the typographical look-and-feel classicthesis developed by André Miede. The style was inspired by Robert Bringhurst's seminal book on typography "*The Elements of Typographic Style*". classicthesis is available for both L^AT_EX and L^yX:

<http://code.google.com/p/classicthesis/>

Happy users of classicthesis usually send a real postcard to the author, a collection of postcards received so far is featured here:

<http://postcards.miede.de/>