

Breast Cancer Intrinsic Subtypes: A Critical Conception in Bioinformatics

Heloisa Helena Zaccaron Milioli

B.Sc. in Biological Sciences

M.Sc. in Genetics

*Thesis submitted in fulfilment of the requirements for the degree of
Doctor of Philosophy*



The University of Newcastle
Faculty of Science and Information Technology
School of Environmental and Life Sciences

Callaghan, NSW
Australia

September, 2016

Statement of Originality

The thesis contains no material which has been accepted for the award of any other degree or diploma in any university or other tertiary institution and, to the best of my knowledge and belief, contains no material previously published or written by another person, except where due reference has been made in the text. I give consent to the final version of my thesis being made available worldwide when deposited in the University's Digital Repository¹, subject to the provisions of the Copyright Act 1968.

September, 2016

Heloisa Helena Zaccaron Milioli

Prof. Pablo Moscato

¹ Unless an Embargo has been approved for a determined period

Statement of Authorship

I hereby certify that the work embodied in this thesis contains a published paper/s/scholarly work of which I am a joint author. I have included as part of the thesis a written statement, endorsed by my supervisor, attesting to my contribution to the joint publication/s/scholarly work.

September, 2016

Heloisa Helena Zaccaron Milioli

Prof. Pablo Moscato

Acknowledgements

I would like to express my deep gratitude to Prof Pablo Moscato. I appreciated the guidance and encouragement he has provided throughout my time at Centre for Bioinformatics, Biomarker Discovery and Information-Based Medicine (CIBM). I have been extremely lucky to have a supervisor who cared so much about my work, and who responded to my questions and queries so promptly. I would like to thank my co-supervisors, A/Prof Regina Berretta and Dr Jannette Sakoff, for their professional advices and useful critiques.

Special thanks should be given to Dr Carlos Riveros for his patient assistance and unconditional support. For all the extra hours he spent working with me, the constructive criticism and friendly advice during my PhD. I am sincerely grateful for sharing his truthful and illuminating views on a number of issues related to the breast cancer project. I am also grateful to Dr Renato Vimieiro for the valuable help in data management, extensive collaboration and magic proofreading. In particular, I thank Inna Tishchenko for her precious effort and intelligence in the analysis of the data.

I would like to extend my thanks to all CIBM students and collaborators who contributed with valuable discussions and enthusiastic encouragements: Ademir Cristiano Gabardo, Ahmed Shamsul Arefin, Amer Abu Zaher, Amir Salehipour, Chloe Warren, Claudio Sanhueza, Francia Jimenez, Leila Moslemi Naeni, Luke Mathieson, Mohammad Nazmul Haque, Nasimul Noman, Natalie de Vries, Nisha Puthiyedth, Shannon Fenn. Thanks for sharing the experience, either positive or negative. I also acknowledge the Hunter Medical Research Institute (HMRI) and University of Newcastle (UoN) staffs for sharing an amazing productive environment.

I express my warm thanks to Jennie Thomas for her enthusiastic support of students and researchers through a number of grants and scholarships. Being part of your *family* is a great honour and an enormous pleasure. Thanks for believing in my research and for funding my dreams. I am also grateful to A/Prof David Wild for guiding me throughout my visit to Bloomington, US.

Special thanks to my beloved family for their unconditional support!

Words cannot express how grateful I am to my mother, father, aunt, brother and nephew for all of the sacrifices that you've made on my behalf. Your positive energies have sustained me thus far. I would also like to thank my in-laws for striving towards my goal. Finally, I would like express appreciation to my beloved husband, Jorge André Martins. I would not be here without him, his patience and love.

*To the best grandmother,
Helena Mafioletti Zaccaron
(Wherever you are)*

Table of Contents

Acknowledgements	V
Table of Contents	IX
List of Figures	XIII
List of Tables	XV
List of Equations	XVII
Abbreviations	XIX
Achievements	XXIII
CHAPTER 1	1
1. INTRODUCTION AND OVERVIEW	1
1.1 Breast Cancer: an Overview	2
1.2 Bioinformatics Resources and Tools	4
1.3 Research Motivation	6
1.3.1 Research Questions	7
1.4 Research Aims and Thesis Structure	7
1.5 References	11
CHAPTER 2	16
2. BREAST CANCER: CURRENT STATUS AND PERSPECTIVES	16
2.1 Breast Carcinogenesis	17
2.2 The Breast Tumour Classification	19
2.3 Intrinsic Subtypes	24
2.3.1 Luminal A and B	25
2.3.2 HER2-enriched	26
2.3.3 Basal-like	26
2.3.4 Normal-like	28
2.3.5 Other groups	28
2.4 Novel Integrative Clusters	29
2.5 Predicting Molecular Subtypes	30
2.6 References	32

CHAPTER 3 **40**

3.	MICROARRAY TECHNOLOGIES AND ‘OMICS’ DATA SETS	40
	3.1 Microarray technologies	41
	3.1.1 Illumina Approach	44
	3.1.2 Affymetrix Platforms	45
	3.2 The METABRIC Breast Cancer Data Set	46
	3.2.1 Biospecimen Collection and Ethics Approval	46
	3.2.2 Gene Expression Data Description	49
	3.2.3 Genotype Calling	49
	3.2.4 The Breast Cancer Cohort	50
	3.3 ROCK: Integrative Breast Cancer Data	50
	3.4 References	52

CHAPTER 4 **56**

4.	IDENTIFICATION OF NOVEL BIOMARKERS FOR BREAST CANCER SUBTYPING	56
	4.1 Introduction	57
	4.2 Methods	58
	4.2.1 Study Design and Computing Resources	58
	4.2.2 Selection of Biomarkers Using the CM1 Score	60
	4.2.3 The Quality of CM1 List Based on Ensemble Learning	61
	4.2.4 Statistical Analysis	61
	4.2.5 Survival Analysis	64
	4.3 Results	64
	4.3.1 Section Description and Resources	64
	4.3.2 Using the CM1 List to Differentiate Breast Cancer Subtypes	65
	4.3.3 The High Levels of Agreement Between CM1 and PAM50 Lists	71
	4.3.4 The Use of an Ensemble Learning with the CM1 List Improves the Subtype Distribution in the METABRIC and ROCK Data Sets	76
	4.3.5 Breast Cancer Intrinsic Subtypes Defined by Clinical Markers and Survival Curves	79
	4.4 Discussion	85
	4.5 Conclusion	86
	4.6 References	87
	4.7 Supporting Information	91

CHAPTER 5 **125**

5.	ITERATIVELY REFINING THE METABRIC SUBTYPE LABELS	125
5.1	Introduction	126
5.2	Methods	127
5.2.1	Transcriptomic Data Set	127
5.2.2	The Refinement Method	127
5.2.3	The CM1 Score	129
5.2.4	Statistical Analysis	129
5.2.5	Clinical Data and Survival Curves	129
5.3	Results and Discussion	130
5.3.1	Discriminative Probes Used to Assign Intrinsic Subtype Labels in the Refinement Process	130
5.3.2	New Subtype Labels Reveal More Reliable Distribution of Clinical Markers and Survival Outcomes	131
5.4	Conclusion	134
5.5	References	135
5.6	Supporting Information	137

CHAPTER 6 **155**

6.	META-FEATURES FOR PREDICTING BREAST CANCER INTRINSIC SUBTYPES	155
6.1	Introduction	156
6.2	Methods	157
6.2.1	Ethics Statement and Data Description	157
6.2.2	Study Design and Computing Resources	158
6.2.3	Statistical Analysis	161
6.3	Results and Discussion	162
6.3.1	Thirteen Meta-features Define Breast Cancer Intrinsic Subtypes	162
6.3.2	An Ensemble Learning Approach Validates the Quality of Meta-features for Predicting Subtypes	168
6.3.3	Expanding Prediction Models Based on Microarray Data	171
6.4	References	172
6.5	Supporting Information	177

CHAPTER 7 **181**

7. BASAL-LIKE BREAST CANCER SUBTYPE	181
7.1 Introduction	184
7.2 Methods	186
7.2.1 Breast Cancer Data Sets	186
7.2.2 Probe Selection Approach	187
7.2.3 Clustering Basal-like Breast Cancer Samples	188
7.2.4 Validation across Data Sets	189
7.2.5 Network Analysis	189
7.2.6 MicroRNA Differential Expression	190
7.2.7 Copy Number Aberration Profiles	190
7.3 Results	191
7.3.1 Survival-related Probes Defining Basal-like Subgroups	191
7.3.2 Basal I and Basal II Validated across Independent Data Sets and Microarray Platforms	200
7.3.3 Clinical Features and Survival Outcomes Supporting the Basal-like Subgroups	200
7.3.4 MicroRNAs Differentially Expressed between Basal I and Basal II	203
7.3.5 Copy Number Aberration Profiles Further Differentiating Basal-like Subgroups	206
7.4 Discussion	209
7.4.1 Survival-related Probes Defining the Molecular Signature of Basal-like Breast Cancer Subgroups	209
7.4.2 MicroRNA Expression Levels Differentiating Basal I from Basal II	210
7.4.3 Genomic Aberrations Further Characterise Basal II and Basal I Subgroups	212
7.4.4 Consensus on the Analysis of Basal-like Breast Cancer Subtypes: a Literature Overview	213
7.5 Conclusion	215
7.6 References	216
7.7 Supporting Information	225
CHAPTER 8	241
8. CONCLUDING REMARKS	241
8.1 Final Statements	242
8.2 Future Directions	246
8.3 Closing Note	248

List of Figures

Figure 3.1 Conceptual view of a cRNA microarray processing.	42
Figure 4.1 The step-by-step process	59
Figure 4.2 The gene expression profile of the balanced top ten probes selected for each of the five breast cancer intrinsic subtypes across 997 samples from the discovery set.	69
Figure 4.3 Gene expression patterns of the 42 probes selected using the CM1 score	70
Figure 4.4 The mRNA log ₂ normalised expression values of 7 novel highly discriminative biomarkers across the five intrinsic subtypes	71
Figure 4.5 Class distribution in the METABRIC discovery and validation, and ROCK set	77
Figure 4.6 Similarity between subtypes distribution in the METABRIC discovery and validation sets, and in the ROCK set	79
Figure 4.7 ER marker distribution across subtypes in the METABRIC data sets	81
Figure 4.8 PR marker distribution across subtypes in the METABRIC data sets	82
Figure 4.9 HER2 distribution across subtypes in the METABRIC data sets.....	83
Figure 4.10 The survival curves for METABRIC discovery and validation sets	84
Figure 4.11 The mRNA log ₂ normalised expression values of 42 probes (A and B) in the CM1 list across the five intrinsic subtypes in the METABRIC discovery and validation, and ROCK 97	
Figure 5.1 Refinement Method	128
Figure 5.2 The heat map of refined intrinsic features selected using CM1 score	131
Figure 5.3 The survival curves for original and refined labels in the METABRIC discovery and validation sets.....	133
Figure 5.4 Mean Final Classifier Performance, as measured by Fleiss' κ against the final ensemble learning labels of all samples, across the 10 different refinement runs	141
Figure 5.5 Evolution of performance of classifiers along iterations in a typical refinement run. The κ values are measured against final ensemble learning labels	142
Figure 5.6 MST- <i>k</i> NN clustering, coloured according to the original METABRIC labels defined by the PAM50 method.....	145
Figure 5.7 MST- <i>k</i> NN clustering, coloured according to the refined labels using an iterative process	146
Figure 5.8 MST- <i>k</i> NN clustering, coloured according to the IntClust classification proposed by Curtis et al. (2012)	147
Figure 6.1 Summary systematic approach	159
Figure 6.2 Meta-features selected with the CM1 score in the METABRIC discovery set	164

Figure 6.3 Gene expression patterns of the 13 meta-features selected using the CM1 score and (α, β) -k-Feature set 165

Figure 6.4 Pairwise expression patterns across intrinsic subtypes in the METABRIC discovery and validation sets 166

Figure 6.5 Individual expression patterns across intrinsic subtypes in the METABRIC discovery and validation sets 167

Figure 6.6 Graph representing an instance of the (α, β) -k-Feature Set; as per the data defined in Table 6.5. 178

Figure 6.7 Graph containing a feasible solution for the (α, β) -k-Feature Set problem; as per the data defined in Table 6.5. 179

Figure 7.1 Heat map of the 80-genes signature in METABRIC training set..... 196

Figure 7.2 Minimum Spanning Tree of the 80-probe signature 197

Figure 7.3 Survival curves in the METABRIC and ROCK data sets..... 201

Figure 7.4 The box plot of miRNAs differentiating Basal I and Basal II subgroups 205

Figure 7.5 Copy number aberration of basal subgroups in METABRIC data set 207

Figure 7.6 The heat map of 400 probes in METABRIC training set 233

Figure 7.7 Network analysis of multiple drug targets for breast cancer therapy 238

Figure 8.1 t-SNE graph of METABRIC samples coloured according to PAM50 244

Figure 8.2 t-SNE graph of METABRIC samples coloured using the refined labels..... 244

List of Tables

Table 2.1 Primary Tumour (T).....	21
Table 2.2 Regional Lymph Nodes (N).....	22
Table 2.3 Distant Metastasis (M).....	22
Table 2.4 Anatomic stage/prognostic groups.....	23
Table 3.1 METABRIC microarray data description.....	47
Table 3.2 Data accession – gene expression and genotyping information.....	48
Table 3.3 Data accession – microRNA expression information.....	49
Table 3.4 Overview of the ten data sets in the ROCK online portal.....	51
Table 4.1 CM1 List.....	66
Table 4.2 Scores and ranks for the CM1 list.....	67
Table 4.3 The ensemble learning overall performance on assigning labels to samples in the METABRIC discovery and validation sets, and ROCK test set.....	73
Table 4.4 Contingency tables for predicted labels using classifiers trained with the CM1 list ..	73
Table 4.5 Contingency tables for predicted labels using classifiers trained with the PAM50 list ..	73
Table 4.6 Contingency tables for predicted labels using classifiers trained with CM1 and PAM50 lists ..	74
Table 4.7 Agreement of the 24 classifiers on assigning labels using Fleiss' kappa statistic	75
Table 4.8 Agreement measured by the Adjusted Rand Index between different labelling.....	76
Table 4.9 The CM1 score calculated for each breast cancer subtype	91
Table 4.10 Summary performance of the classifiers using the CM1 list	92
Table 4.11 Summary performance of the classifiers using the PAM50 list.....	94
Table 4.12 The agreement between sample labelling with Fleiss' Kappa measure and the Jensen-Shannon divergence of two probability distributions	95
Table 4.13 The Jensen-Shannon divergence of two probability distributions	96
Table 5.1 Contingency table for predicted labels vs. initial subtypes (rows and columns, respectively).....	130
Table 5.2 Number of samples for each clinical marker in the METABRIC data set according to the PAM50 method and refinement process	132
Table 5.3 Refined subtype labels in the METABRIC data set	137
Table 5.4 List of the 24 classifiers used in the ensemble learning.....	137
Table 5.5 Average agreement of classifiers per subtype.....	138

Table 5.6 Probe appearance after ten iterative processes and the respective annotation based on Dunning et al. (2010) and Illumina array data.....	139
Table 5.7 The percentage of PAM50 labels matching integrative clusters (IntClust 1-10) in the METABRIC study.....	148
Table 5.8 The percentage of Refined labels matching integrative clusters (IntClust 1-10) in the METABRIC study.....	149
Table 6.1 List of meta-features selected with CM1 score and (α, β) -k Feature set.....	163
Table 6.2 Contingency tables for predicted labels using ensemble learning trained with 13 meta-features Discovery set Validation set	168
Table 6.3 Performance of 22 Weka classifiers on predicting labels in the METABRIC discovery and validation sets	169
Table 6.4 Fleiss' kappa values and Adjusted Rand Index for the discovery and validation sets.....	170
Table 6.5 An example of numerical matrix with five features and six samples belonging to class <i>F</i> or <i>G</i>	177
Table 7.1 The 80-genes signature related to survival	198
Table 7.2 Clinical information of patients and tumour samples in the METABRIC data set ...	202
Table 7.3 MicroRNAs differentiating basal-like breast cancer subgroups.....	203
Table 7.4 MicroRNAs and corresponding target genes.....	204
Table 7.5 Cytobands associated with significant CNA acquisitions	208
Table 7.6 Basal-like samples classification for the validation set	225
Table 7.7 Basal-like samples classification for the validation set	225
Table 7.8 The centroids computed for differentiating Basal I and Basal II.....	225
Table 7.9 The functional annotation of G1 probes according to DAVID	225
Table 7.10 The functional annotation of G2 probes according to DAVID	225
Table 7.11 The functional annotation of G3 probes according to DAVID	225
Table 7.12 MicroRNAs differentiating Basal I and Basal II	226
Table 7.13 MicroRNAs and gene targets in Basal I.....	227
Table 7.14 MicroRNAs and gene targets in Basal II.....	230
Table 7.15 Summary gene targets and corresponding drugs	237

List of Equations

Equation 4.1 CM1 score	60
Equation 4.2 Cramer's V	62
Equation 4.3 Average sensitivity	62
Equation 4.4 Fleiss' kappa.....	63
Equation 4.5 Adjusted Rand Index	63
Equation 7.1 Normalisation	189

Abbreviations

AACR	Australasian Association of Cancer Registries
ACS	American Cancer Society
AIHW	Australian Institute of Health and Welfare
AJCC	American Joint Committee on Cancer
AR	Androgen receptor
ARI	Adjusted Rand Index
BL1	Basal-like 1
BL2	Basal-like 2
BLBC	Basal-like breast cancer
BLIA	Basal-like immune-activated
BLIS	Basal-like immune-suppressed
ChIP-chip	Chromatin immunoprecipitation on chip
CIBEX	Center for information biology gene expression database
CIBM	Centre for Bioinformatics, Biomarker Discovery and Information-Based Medicine
CGH	Comparative genomic hybridization
CNA	Copy number aberration
CNV	Copy number variation
CTD	Comparative Toxicogenomic Database
DamID	DNA adenine methyltransferase identification
DAVID	Database for Annotation, Visualization and Integrated Discovery
DDBJ	DNA Data Bank of Japan
DNA	Deoxyribonucleic acid
EBI	European Bioinformatics Institute
EGA	European Genome-Phenome Archive
EpCAM	Epithelial cell adhesion molecule
ER	Oestrogen receptor
FGED	Functional Genomics Data Society
FOIPPA	Freedom of Information and Protection of Privacy Act
FS	Feature Selection
GEO	Gene Expression Omnibus
HER2	Human epidermal growth factor receptor 2
HREC	Human Research Ethics Committee

HTC	High content screening
HTS	High-throughput screening
ICGC	International Cancer Genomics Consortium
IDC	Invasive ductal carcinoma
IHC	Immunohistochemical
IHGSC	International Human Genome Sequencing Consortium
ILC	Invasive lobular carcinoma
IM	Immunomodulatory
JS	Jensen Shannon
Ki-67	Antigen identified by monoclonal antibody Ki-67
kNN	k nearest neighbours
LAR	Luminal androgen receptor
lincRNA	long intergenic non-coding RNA
MA	Memetic algorithm
MCC	Matthews' Correlation Coefficient
MDL	Minimum Description Length Principle
METABRIC	Molecular Taxonomy of Breast Cancer International Consortium
MIAME	Minimum Information About a Microarray Experiment
microRNA	miRNA
MGED	Microarray Gene Expression Data Society
MS	Menopausal status
MST	Minimum Spanning Tree
NCBI	National Center for Biotechnology Information
NOS	Not otherwise specified
NPI	Nottingham prognostic score
NSC	Nearest Shrunken Centroids
NST	No special type
ORF	Open reading frame
PIPA	Personal Information Protection Act
PIPEDA	Personal Information Protection and Electronic Documents Act
PR	Progesterone receptor
PRC	Priority Research Centres
RHD	Research Higher Degree
RNA	Ribonucleic acid
ROCK	Research Online Cancer Knowledgebase
RT-PCR	Reverse-transcriptase Polymerase chain reaction

SAM	Sentrix® Array Matrix
SCM	Subtype Classification Model
SNP	Single nucleotide polymorphism
SSP	Single Sample Predictor
TCGA	The Cancer Genome Atlas
TEND	Trends in the Exploration of Novel Drug targets
TNBC	Triple-negative breast cancer
TNM	Tumour size, nodes, metastasis
TTD	Therapeutic Target Database
UCSC	University of California Santa Cruz
WEKA	Waikato Environment for Knowledge Analysis

Achievements

During my PhD, I applied for grants; submitted manuscripts for publication; and attended workshops, conferences and seminars. The relevant achievements are listed as follows:

Grants Awarded

- Hunter Medical Research Institute, 2014.

JENNIE THOMAS MEDICAL RESEARCH TRAVEL GRANT (AUD \$10,000)

- Hunter Cancer Research Alliance, 2015.

HCRA TRAVEL GRANT (AUD \$1,000)

- Hunter Cancer Research Alliance, 2016.

HCRA PhD Research Award 2016 (AUD \$5,000).

- EMBL Australia PhD, 2016.

Travel Grant to attend the 18th EMBL PhD Symposium (AUD \$3,000).

- XII ELAG Course Fellowship (USD \$700)

Instituto Genética Para Todos – Brazil (unable to attend)

Papers Published in Journals

MILIOLI, H.H.; VIMIEIRO, R.; RIVEROS, C.; TISHCHENKO, I.; BERRETTA, R.; MOSCATO, P. (2015) The discovery of novel biomarkers improves breast cancer intrinsic subtype prediction and reconciles the original PAM50 labels in the METABRIC data set. *PLoS One*; 10(7): 0129711. doi: 10.1371/journal.pone.0129711

MILIOLI, H.H. (2015). The IMPAKT of breast cancer research: fundamental science and clinical medicine. *Future Science OA*; (0). doi: 10.4155/fso.15.69

MILIOLI, H.H.; VIMIEIRO, R.; TISHCHENKO, I.; RIVEROS, C.; BERRETTA, R.; MOSCATO, P. (2016) Iteratively refining breast cancer intrinsic subtypes in the METABRIC dataset *BioData Mining*; 9:2. doi: 10.1186/s13040-015-0078-9

TISHCHENKO, I.; **MILIOLI, H.H.**; RIVEROS, C.; MOSCATO, P. (2016) Extensive Transcriptomic and Genomic Analysis Provides New Insights about Luminal Breast Cancers. *PLoS One*; 11(6): e0158259. doi: 10.1371/journal.pone.0158259

MILIOLI, H.H. Life as an early career researcher: interview with Heloisa Helena Milioli. *Future Science OA*; 1(4) (2016). doi: 10.4155/fsoa-2016-0033.

MILIOLI, H.H.; TISHCHENKO, I.; RIVEROS, C.; BERRETTA, R.; MOSCATO, P. Basal-like breast cancer: molecular profiles, clinical features and survival outcomes. *BMC Med Genomics*; 10(1):19 (2017). doi: 10.1186/s12920-017-0250-9.

MILIOLI, H.H.; RIVEROS, C.; VIMIEIRO, R.; BERRETTA, R.; MOSCATO, P. Meta-features modelling gene expression imbalances: an innovative strategy for breast cancer subtype prediction. Manuscript in preparation to be submitted for publication at Genomics, Proteomics and Bioinformatics (GPB).

Abstracts Published

MILIOLI, H.H.; TISHCHENKO, I.; RIVEROS, C.; SAKOFF, J.; BERRETTA, R.; MOSCATO, P. Consensus on breast cancer cell lines classification for an effective and efficient clinical decision-making. *IMPAKT 2015 Breast Cancer Conference. Annals of Oncology* 26 (suppl 3):iii32-iii33 (2015). doi: 10.1093/annonc/mdv121.08

MILIOLI, H.H.; TISHCHENKO, I.; RIVEROS, C.; BERRETTA, R.; MOSCATO, P. Molecular classification of basal-like breast cancer subtypes based on predictive survival markers. *IMPAKT 2015 Breast Cancer Conference. Annals of Oncology*. 26 (suppl 3):iii17-iii18 (2015). doi: 10.1093/annonc/mdv117.11

MILIOLI, H.H., TISHCHENKO, I., RIVEROS, C., BERRETTA, R. & MOSCATO, P. (2015) Basal-like breast cancer subgroups uncovered by genomic and transcriptomic profiles and overall survival outcomes. *Hunter Cancer Research Alliance Annual Symposium. Asia-Pacific Journal of Clinical Oncology* 11(Suppl. 5):6-19 (2015). doi: 10.1111/ajco.12444

TISHCHENKO, I., **MILIOLI, H.H.**, RIVEROS, C. & MOSCATO, P. How intrinsic are luminal breast cancer subtypes? *Hunter Cancer Research Alliance Annual Symposium. Asia-Pacific Journal of Clinical Oncology* 11(Suppl. 5):6-19 (2015). doi: 10.1111/ajco.12444

MILIOLI, H.H., SANHUEZA, C., BERRETTA, R. & MOSCATO, P. (2015) ABSTRACT P40 Breast Cancer Molecular Portraits of Intrinsic Subtypes and Integrative Clusters in the METABRIC Data Set. *Hunter Cancer Research Alliance Annual Symposium. Asia-Pacific Journal of Clinical Oncology* 12(Suppl. 6):13-34 (2016). doi: 10.1111/ajco.12618

Oral Presentations

MILIOLI, H.H.; VIMIEIRO, R.; TISHCHENKO, I.; RIVEROS, C.; BERRETTA, R.; MOSCATO, P. Refining the breast cancer molecular subtypes in the METABRIC data set. *World Congress on Controversies in Breast Cancer (CoBRA), 2015. Melbourne, AU.*

MILIOLI, H.H.; SANHUEZA, C.; RIVEROS, C.; BERRETTA, R.; MOSCATO, P. Breast cancer molecular portraits of intrinsic subtypes and integrative clusters in the METABRIC data set. **Young Scientist Award** 2nd *World Congress on Controversies in Breast Cancer (CoBrCa) 2016. Barcelona, Spain.*

MILIOLI, H.H.; TISHCHENKO, I.; RIVEROS, C.; BERRETTA, R.; MOSCATO, P. Basal-like breast cancers uncovered by genomic and transcriptomic profiles and patients' overall survival. *Sydney Cancer Conference (SCC) 2016. Sydney, AU.*

Poster Sessions

MILIOLI, H.H.; VIMIEIRO, R.; RIVEROS, C.; SAKOFF, J.; BERRETTA, R.; MOSCATO, P. Breast Cancer Subtypes Individuation Driving Novel Drug Targets for Tailored Therapies. *Translational Cancer Research Conference, 2013. Newcastle, AU.*

MILIOLI, H.H.; VIMIEIRO, R.; RIVEROS, C.; BERRETTA, R.; MOSCATO, P. Identification of novel biomarkers for predicting breast cancer intrinsic subtypes. *ASMR Satellite Scientific Meeting, 2014. Newcastle, AU.*

MILIOLI, H.H.; RIVEROS, C.; VIMIEIRO, R.; MOSCATO, P. Meta-features as predictors of breast cancer intrinsic subtype in the METABRIC gene expression data set. **Best Poster Award (Bronze Medal)** *International Conference on Bioinformatics, 2014. Sydney, AU.*

RIVEROS, C.; **MILIOLI, H.H.**; VIMIEIRO, R.; BERRETTA, R.; MOSCATO, P. Discovery of gene interactions by GPU-enabled computation of pairwise expression level metafeatures. *International Conference on Bioinformatics, 2014. Sydney, AU.*

MILIOLI, H.H.; RIVEROS, C.; VIMIEIRO, R.; TISHCHENKO, I.; BERRETTA, R.; MOSCATO, P. Using an iterative approach to reclassify sample subtypes in the METABRIC breast cancer data set. **Best Poster Award (Third place)** *BioInfoSummer, 2014. Melbourne, AU.*

MILIOLI, H.H.; TISHCHENKO, I.; RIVEROS, C.; BERRETTA, R.; MOSCATO, P. Basal-like breast cancer subsets revealed by survival predictor genes. *ASMR Satellite Scientific Meeting, 2015. Newcastle, AU.*

MILIOLI, H.H.; TISHCHENKO, I.; RIVEROS, C.; BERRETTA, R.; MOSCATO, P. Molecular classification of basal-like breast cancer subtypes based on predictive survival markers. *IMPAKT 2015 Breast Cancer Conference. Brussels, BE.*

MILIOLI, H.H.; TISHCHENKO, I.; RIVEROS, C.; SAKOFF, J.; BERRETTA, R.; MOSCATO, P. Consensus on breast cancer cell lines classification for an effective and efficient clinical decision-making. *IMPAKT 2015 Breast Cancer Conference. Brussels, BE.*

MILIOLI, H.H.; RIVEROS, C.; VIMIEIRO, R.; MOSCATO, P. Meta-features predicting gene expression imbalances across breast cancer intrinsic subtypes. *EMBL Australia PhD Symposium, 2015. Melbourne, AU.*

TISHCHENKO, I., **MILIOLI, H.H.**, RIVEROS, C. & MOSCATO, P. How intrinsic are luminal breast cancer subtypes? *Hunter Cancer Research Alliance Annual Symposium, 2015. Newcastle, AU.*

MILIOLI, H.H., TISHCHENKO, I., RIVEROS, C., BERRETTA, R. & MOSCATO, P. Basal-like breast cancer subgroups uncovered by genomic and transcriptomic profiles and overall survival outcomes. *Hunter Cancer Research Alliance Annual Symposium, 2015. Newcastle, AU.*

NAENI, L., MILIOLI, H.H., TISHCHENKO, BERRETTA, R. & MOSCATO, P. (2015) A New Clustering Approach Identifies Candidate Biomarkers for Breast Cancer Subtyping. *BioInfoSummer, 2015. Sydney, AU.*

MILIOLI, H.H.; RIVEROS, C.; VIMIEIRO, R.; MOSCATO, P. Meta-features predicting gene expression imbalances across breast cancer intrinsic subtypes. **Best Poster Presentation** *BioInfoSummer, 2015. Sydney, AU.*

Other Presentations

Confirmation Year Presentation

Faculty of Science and IT. The University of Newcastle, 2013.

RHD candidates are required to submit the 'Confirmation Year Report' and present the research overview. In August 2013, I presented the preliminary results in the Faculty of Science and IT as an open seminar.

HCRA, ECR and PhD Student (HEAPS) Seminar Series

Hunter Medical Research Institute. The University of Newcastle, 2014 and 2015.

The HEAPS seminar series are organised by the Hunter Cancer Research Alliance (HCRA) for RHD students and supervisors. It is an opportunity for researchers to practice presenting (and critiquing) work in a local and highly supportive environment. In 2014 and 2015, I presented and discussed the results of my research as well as supported other researchers' work.

HUBS3302 Bioinformatics Mini-Conference

Faculty of Health and Medicine. The University of Newcastle, 2014 and 2015.

The purpose of this event is to inspire students in the field and, specially, in their final project for the discipline. In the 2014 and 2015 Bioinformatics Mini-Conference, organised by Belinda Goldie, I presented my research on breast cancer.

Science and Engineering Challenge

Faculty of Engineering and Built Environment. The University of Newcastle, 2014, 2015 and 2016.

The 'Science and Engineering Challenge' organise a number of events aimed at challenging students of all different ages in Science and Engineering. As part of the team, I coordinated activities in Tamworth (2014), Muswellbrook (2014), Dubbo (2015), Newcastle (2015), Central Coast (2016) and Narrabri (2016), and presented my research to the Rotary International (Australian Rotary Districts) in Tamworth and Dubbo.

Faculty Progress Seminar

Faculty of Science and IT. The University of Newcastle, 2015.

Students in the Faculty of Science and IT are required to present a Progress Seminar after completing 2 to 3 years of a PhD. In June 2015, I discussed the overall aims and results of my research and outlined my thesis to fellow RHD candidates and academics in the school.

Google Computer Science for High Schools

Faculty of Engineering and Built Environment. The University of Newcastle, 2015 and 2016.

The University of Newcastle's Computer Science 4 High Schools (CS4HS) is an introductory workshop for in-service and pre-service teachers (both at primary and secondary level), and career advisors focused on developing competencies included in the recently approved Digital Technologies curriculum and is accredited by BOSTES. In three events, I had the opportunity to explain the relevance of computer science to analyse biological/medical data.

Relevant Activities

Course: **Winter School in Mathematical and Computational Biology**

University of Queensland (UQ), Brisbane, 2013.

The winter school introduced mathematical and computational biology and bioinformatics to advanced undergraduate and postgraduate students, postdoctoral researchers and others working in the field. Important topics, such as mathematics, statistics, computer science, information technology, biology, chemistry and medical sciences and engineering, were selected for each day. Lectures and interactive discussions were ministered by national and international authorities.

Course: **European Molecular Biology Laboratory (EMBL) Australia PhD Course**

Australian National University (ANU), Canberra, 2014.

EMBL Australia offered to sixty students a unique introduction to research with the annual EMBL Australia PhD Course. The two-week program shows students how their research fits into the bigger picture of science, and introduces a range of fields including: bioinformatics, developmental biology, genomics, systems biology and regenerative medicine.

Course: **European Molecular Biology Laboratory (EMBL) Australia PhD Course**

Welcome Genome Campus, Hinxton, UK, 2016.

This course introduced a wide range of post-genome techniques including practical experience in performing (1) high-throughput RNAi screening, (2) microarray gene expression analysis and interpretation, using a range of commercial and academic software tools, (3) next-generation sequencing and alignment; (4) protein-protein interaction networks and integration with other data sources, and (5) pathway analysis. Laboratory work was based on the training of state-of-the-art methods and complementary approaches to address biological and medical questions.

Training: **Collaborative Research Training in Human Genetics and Bioinformatics**

Centre for Bioinformatics, Biomarker Discovery and Information-Based Medicine (CIBM). The University of Newcastle, 2014.

The CIBM established a research-training program in 2014 that contributed to improve the capacity of young investigators to conduct human genetics and bioinformatics research. The training promoted scientific collaborations between the University of Newcastle and international (undergraduate) students. The proposed program provided opportunities to generate expertise that could contribute to the

long-term goal of harnessing genetic knowledge and bioinformatics skills to diagnose, prevent, or treat diseases. Training activities were coordinated, facilitated and monitored by Prof. Pablo Moscato, A/Prof Regina Berretta and PhD student Heloisa Helena Milioli.

Short-term Exchange Program: **Cheminformatics and Chemogenomics Research Group (CCRG)**

Indiana University (IU), Bloomington USA, 2015.

Further investigation on cheminformatics and toxicogenomics has been developed in collaboration with A/Prof. David J. Wild (May/June 2015), at the School of Informatics and Computing in Bloomington (USA). These approaches were used to delineate drug-targets for basal-like breast cancer, one of the most aggressive subtypes with limited therapy response. Further research, however, is required to design and perform in vitro tests.

Organising Committee: **Australian Society for Medical Research (ASMR) Satellite Scientific Meeting**

Hunter Medical Research Institute (HMRI), Newcastle, 2015.

This event showcases the recent research achievements of Hunter scientists, encourages postgraduate and student interactions and fosters collaboration between researchers within the Faculty of Health and Medicine, HMRI and the international community. In the 2015 edition, I was member of the committee.

Abstract

Breast cancers have been uncovered by high-throughput technologies that allow the investigation at the genomic, transcriptomic and proteomic levels. In the early 2000s, the gene expression profiling has led to the classification of five intrinsic subtypes: luminal A, luminal B, HER2-enriched, normal-like and basal-like. A decade later, the spectrum of copy number aberrations has further expanded the heterogeneous architecture of this disease with the identification of 10 integrative clusters (IntClusts). The referred classifications aim at explaining the diverse phenotypes and independent outcomes that impact clinical decision-making. However, intrinsic subtypes and IntClusts show limited overlap. In this context, novel methodologies in bioinformatics to analyse large-scale microarray data will contribute to further understanding the molecular subtypes. In this study, we focus on developing new approaches to cover multi-perspective, highly dimensional, and highly complex data analysis in breast cancer. Our goal is to review and reconcile the disease classification, underlying the differences across clinicopathological features and survival outcomes. For this purpose, we have explored the information processed by the Molecular Taxonomy of Breast Cancer International Consortium (METABRIC); one of the largest of its type and depth, with over 2000 samples. A series of distinct approaches combining computer science, statistics, mathematics, and engineering have been applied in order to bring new insights to cancer biology. The translational strategy will facilitate a more efficient and effective incorporation of bioinformatics research into laboratory assays. Further applications of this knowledge are, therefore, critical in order to support novel implementations in the clinical setting; paving the way for future progress in medicine.

Keywords

Breast cancer, Intrinsic subtypes, Integrative clusters, IntClusts, Microarray, Gene expression, Copy number aberration, MicroRNA, METABRIC, Feature selection, Data mining, Ensemble learning, Prediction models, Classification

