MOLECULAR AND REGULATORY PROFILING OF MURINE TRANSCRIPTION FACTOR CTCFL IN NON-GERM CELLS AND MALE GERMLINE STEM CELLS

MAISARAH BINTI AB SAMAD

UNIVERSITI SAINS MALAYSIA

2022

MOLECULAR AND REGULATORY PROFILING OF MURINE TRANSCRIPTION FACTOR CTCFL IN NON-GERM CELLS AND MALE GERMLINE STEM CELLS

by

MAISARAH BINTI AB SAMAD

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

April 2022

ACKNOWLEDGEMENT

All praises and gratitude are to Allah SWT, the Lord to whom every single creature in heaven and the earth belongs. Alhamdulillah for giving us the opportunity and help us endlessly in finishing the thesis. I would like to express my gratitude to Professor Dr Shaharum Shamsuddin, my principal supervisor, for the chance he gave me to participate in the USM-RIKEN Interdisciplinary Collaboration for Advanced Sciences (URICAS). I am also very grateful for his advice, support, and constructive comment throughout the journey. My deepest gratitude goes to Professor Kuniya Abe, my field supervisor from RIKEN BioResource Research Center (RIKEN BRC), Tsukuba, Japan, for the opportunity to work towards this thesis and his support over the last four years of my PhD. I would also like to thank all of my colleagues in the Technology and Development Team for Mammalian Genome Dynamics in RIKEN BRC (Dr Suzuki Shinnosuke, Dr Ura Hiroki, Dr Sugimoto Michihiko, Dr Tada Yuhki, Mrs Koga Yumiko, Mr Cho Dooseon and Mrs Kusayama Miwako) for their advice and generous assistance throughout my stay in RIKEN BRC. I would like to thank RIKEN for awarding me with the International Program Associate fellowship, which has allowed me to carry out this research in RIKEN BRC, Tsukuba. Japan. I would also like to thank the Public Service Department of Malaysia (JPA) for awarding me the Program Pelajar Cemerlang fellowship and School of Health Sciences, USM, to award me the Graduate Research Assistance scheme. Sincere gratitude also to my co-supervisor, Dr Daruliza Kernain Mohd Azman, for the continuous support received. Finally, I would like to thank my mother, family and friends for their patience, understanding, encouragement and support throughout everything.

TABLE OF CONTENTS

| ACK | NOWLED | GEMENT | ii | |
|------|-----------|---|-------|--|
| TABL | E OF CO | NTENTS | iii | |
| LIST | OF TABI | LES | xii | |
| LIST | OF FIGU | RES | . xvi | |
| LIST | OF SYMI | BOLSx | xviii | |
| LIST | OF ABBI | REVIATIONS | xxix | |
| LIST | OF APPE | ENDICESx | xxii | |
| ABST | RAK | X | xxiii | |
| ABST | RACT | | XXXV | |
| СНА | PTER 1 | INTRODUCTION | 1 | |
| 1.1 | Backgrou | and of Study | 1 | |
| 1.2 | Problem | Problem Statement | | |
| 1.3 | Rationale | Rationale of Study | | |
| 1.4 | Research | Research Objectives | | |
| 1.5 | Thesis O | rganisation | 9 | |
| СНА | PTER 2 | LITERATURE REVIEW | 11 | |
| 2.1 | An Overv | view of Transcription Factors and Histone Modifications | 11 | |
| | 2.1.1 | Chromatin organisation | 11 | |
| | 2.1.2 | Transcription factors | 13 | |
| | 2.1.3 | Histone modifications | 15 | |
| 2.2 | Transcrip | otion Factor: CCCTC-binding Factor (CTCF) | 20 | |
| | 2.2.1 | Background | 20 | |
| | 2.2.2 | CTCF target sites (CTSes) | 20 | |
| | 2.2.3 | CTCF in the development and cancer | 24 | |
| 23 | Histone N | Modifying Compley: Polycomb Repressor Compley 2 (PRC2) | 27 | |

| | 2.3.1 | Background | 27 |
|-----|-----------|--|----|
| | 2.3.2 | PRC2 recruitment to chromatin | 31 |
| | 2.3.3 | The role of PRC2 in development and cancer | 32 |
| | 2.3.4 | Relationship between PRC2 and CTCF in gene expression regulation | 35 |
| 2.4 | Testicula | ar Transcription Factor: CCCTC-binding factor-like (CTCFL) | 39 |
| | 2.4.1 | Background | 39 |
| | 2.4.2 | Characterisation of CTCFL binding to target sites | 42 |
| | 2.4.3 | The role of CTCFL in male germline development | 44 |
| | 2.4.4 | The role of CTCFL in cancer | 48 |
| 2.5 | Genome- | wide Analysis Methods for Molecular and Regulatory Profiling | 51 |
| | 2.5.1 | Background | 51 |
| | 2.5.2 | Transcriptomic Profiling Methods | 51 |
| | | 2.5.2(a) Microarray | 52 |
| | | 2.5.2(b) RNA sequencing (RNA-seq) | 55 |
| | | 2.5.2(c) Comparing microarray and RNA-seq | 60 |
| | 2.5.3 | Chromatin Immunoprecipitation and Sequencing (ChIP-seq) | 62 |
| СНА | PTER 3 | GENERAL METHODS AND MATERIALS | 69 |
| 3.1 | Materials | S | 69 |
| 3.2 | Mammal | ian Cell Culture | 69 |
| | 3.2.1 | Mammalian cell lines | 69 |
| | | 3.2.1(a) Mouse embryonic stem cells (ESCs) | 69 |
| | | 3.2.1(b) Immortalised testicular stromal cells (JK1 cells) | 70 |
| | | 3.2.1(c) Mouse male germline stem cells (GSCs) | 70 |
| | 3.2.2 | Preparation of media for cell culture | 72 |
| | | 3.2.2(a) Culture medium for ESCs | 72 |
| | | 3.2.2(b) Culture medium for JK1 cell line and inactivated MFFs | 72 |

| | | 3.2.2(c) Growth medium for primary MEFs72 |
|-----|------------|--|
| | | 3.2.2(d) Basal medium and culture medium for GSCs72 |
| | 3.2.3 | Cryopreservation of cells |
| | 3.2.4 | Revival of cells |
| | 3.2.5 | Preparation of inactivated MEFs as feeder cells76 |
| | | 3.2.5(a) Extraction and subculturing of MEFs76 |
| | | 3.2.5(b) Treatment of MEFs with mitomycin C77 |
| | 3.2.6 | Subculturing of Mammalian Cell Lines |
| | | 3.2.6(a) Subculturing ESCs78 |
| | | 3.2.6(b) Subculturing of JK1 cells79 |
| | | 3.2.6(c) Subculturing of GSCs80 |
| 3.3 | Extraction | on Nucleic Acids and Assessment for Quality and Quantity |
| | 3.3.1 | Plasmid extraction |
| | 3.3.2 | Total RNA extraction from mammalian cell lines |
| | 3.3.3 | Quantity and quality assessment of nucleic acids |
| 3.4 | Reverse | Transcription for cDNA Synthesis |
| 3.5 | Quantita | tive Polymerase Chain Reaction (qPCR) |
| | 3.5.1 | TaqMan® probe qPCR assay87 |
| | 3.5.2 | SYBR® Green qPCR assay90 |
| | 3.5.3 | Relative quantification of qPCR data93 |
| 3.6 | Protein I | Expression Analysis by Immunofluorescence Staining95 |
| 3.7 | Protein I | Expression Analysis by Western blot |
| | 3.7.1 | Lysate extraction from mammalian cells |
| | 3.7.2 | Protein quantification by Bradford protein assay99 |
| | 3.7.3 | Polyacrylamide gel electrophoresis (PAGE) |
| | 3.7.4 | Protein transfer by tank blotting (Western blot)h100 |
| | 3.7.5 | Immunodetection |

| | 3.7.6 | Mild antibody stripping | 102 |
|------|-----------|---|-----|
| 3.8 | Microard | ray Hybridisation | 103 |
| | 3.8.1 | Preparation of RNA Spike Mix for sample preparation | 103 |
| | 3.8.2 | Preparation of Cyanine 3 (Cy3) labelling reaction | 106 |
| | | 3.8.2(a) Step A: Preparation of T7 primer annealing reaction mixture | 106 |
| | | 3.8.2(b) Step B: Preparation of cDNA synthesis | 106 |
| | | 3.8.2(c) Step C: Preparation of transcription reaction | 107 |
| | 3.8.3 | Purification of Cy3-labelled complementary RNA (cRNA) | 110 |
| | 3.8.4 | Quantification and quality assessment of cRNA | 111 |
| | 3.8.5 | Hybridisation | 113 |
| | 3.8.6 | Microarray washing | 115 |
| | 3.8.7 | Microarray scanning and feature extraction | 116 |
| 3.9 | Smart-se | eq2 Library Preparation for RNA sequencing | 116 |
| | 3.9.1 | Isolation of mRNA from 10 ng total RNA samples | 120 |
| | 3.9.2 | Reverse transcription (RT) and first-strand cDNA synthesis | 121 |
| | 3.9.3 | Preamplification of cDNA | 124 |
| | 3.9.4 | Purification of PCR products | 127 |
| | 3.9.5 | Library preparation | 128 |
| 3.10 | Chromat | tin Immunoprecipitation (ChIP) | 132 |
| | 3.10.1 | Chromatin crosslinking and shearing by sonication | 134 |
| | 3.10.2 | Chromatin immunoprecipitation | 135 |
| | 3.10.3 | Capturing of antibody/chromatin complex, elution and reverse crosslinks | 135 |
| | 3.10.4 | Purification of Input and ChIP-ed DNA | 136 |
| | 3.10.5 | Analysis of ChIP-ed DNA samples by qPCR | 137 |
| 3.11 | Library 1 | Preparation for ChIP-seq | 139 |
| | 3.11.1 | End repair, A-tailing and adapter ligation | 139 |

| | 3.11.2 | PCR amp | olification | 142 |
|------|-----------|---------------------|--|-----|
| 3.12 | Library (| Quantificat | ion for RNA-seq and ChIP-seq | 148 |
| | 3.12.1 | Assessme | ent of library fragment size | 148 |
| | 3.12.2 | Library d | ilution and quantification by qPCR | 148 |
| | | | ANT CTCFL EXPRESSION IN PLURIPOTENT FERENTIATED SOMATIC CELLS | 154 |
| 4.1 | Introduct | ion | | 154 |
| 4.2 | Materials | s and Meth | ods | 156 |
| | 4.2.1 | Establish | ment of ESC-CTCFL and JK1-CTCFL cell lines | 156 |
| | | 4.2.1(a) | The pPB-3×FLAG-CTCFL plasmid | 156 |
| | | 4.2.1(b) | Verification of 3×Flag-Ctcfl construct in pPB-3×FLAG-CTCFL plasmid | 158 |
| | | 4.2.1(c) | Transfection of pPB-3×FLAG-CTCFL plasmid in ESCs and JK1 cell lines | 160 |
| | | 4.2.1(d) | Evaluation of 3×FLAG-CTCFL protein expression | 164 |
| | | 4.2.1(e) | Evaluation of 3×FLAG-CTCFL transcript expression by RT-qPCR | 169 |
| | 4.2.2 | | gene expression changes by microarray | 170 |
| | | | Microarray hybridisation and gene expression analysis | 170 |
| | | 4.2.2(b) | Validation of differentially expressed (DE) genes | 172 |
| | 4.2.3 | | al annotation analysis and expression profiling of ally expressed (DE) genes | 173 |
| | | 4.2.3(a) | Functional annotation analysis | 173 |
| | | 4.2.3(b) | Gene expression profile in adult tissues | 173 |
| | | 4.2.3(c) | Gene expression clustering | 174 |
| 4.3 | Results | | | 174 |
| | 4.3.1 | Expression CTCFL of | on of 3×FLAG-CTCFL in ESC-CTCFL and JK1- | 174 |

| | 4.3.2 | expressio | ially expressed (DE) genes upon 3×FLAG-CTCFL on in ESC-CTCFL(DOX) and JK1-CTCFL(DOX) | . 183 |
|-----|------------|-------------|--|-------|
| | 4.3.3 | | ons of aberrant CTCFL expression to the gene | . 188 |
| | | 4.3.3(a) | Aberrant CTCFL expression alters cell developmental genes in pluripotent ESCs | .188 |
| | | 4.3.3(b) | Aberrant CTCFL expression in differentiated JK1 cells | .195 |
| | | 4.3.3(c) | Several CTA genes were upregulated by aberrant CTCFL expression in non-germ cells | .202 |
| 4.4 | Discussion | on | | . 204 |
| 4.5 | Limitatio | on of Study | <i>T</i> | . 209 |
| 4.6 | Summar | у | | . 210 |
| | | | GENE EXPRESSION KNOCKDOWN IN MALE | . 212 |
| 5.1 | Introduct | ion | | . 212 |
| 5.2 | Materials | s and Meth | ods | . 215 |
| | 5.2.1 | Ctcfl kno | ckdown in GSCs using DsiRNA | . 215 |
| | | 5.2.1(a) | Culture and maintenance of GSCs | .215 |
| | | | Dicer-substrate RNA interference (DsiRNA) for <i>Ctcfl</i> knockdown in GSCs | |
| | | 5.2.1(c) | Transfection of DsiRNAs into GSCs | .217 |
| | | 5.2.1(d) | Gene expression analysis by RT-qPCR | .219 |
| | 5.2.2 | Preparati | on of Smart-seq2 libraries for RNA-seq | . 220 |
| | | 5.2.2(a) | Preparation of low quantity RNA and pre- amplification of cDNA samples | .220 |
| | | 5.2.2(b) | Library preparation and sequencing | .221 |
| | | 5.2.2(c) | Sequence alignment and gene counts | .223 |
| | 5.2.3 | Gene exp | pression analysis for Smart-seq2 libraries | . 225 |
| | | 5.2.3(a) | Differential gene expression analysis | .225 |

| | | 5.2.3(b) | Downstream analysis |
|-----|-------------------|------------|---|
| 5.3 | Results | | |
| | 5.3.1 | Knockdov | wn of <i>Ctcfl</i> expression in GSCs |
| | | 5.3.1(a) | Preliminary assessment of DsiRNAs226 |
| | | 5.3.1(b) | Ctcfl knockdown in GSCs230 |
| | 5.3.2 | Smart-sec | q2 libraries for GSCs(neg) and GSC-Ctcfl(kd)234 |
| | 5.3.3 | • | of differentially expressed (DE) genes in GSC236 |
| | | 5.3.3(a) | Differential gene expression analysis236 |
| | | 5.3.3(b) | The functional annotation analysis of the DE genes $\dots 240$ |
| | | 5.3.3(c) | Ctcfl knockdown in GSCs represses the genes involved with apoptosis and TNF signalling pathways |
| | | 5.3.3(d) | The DE genes in GSC-Ctcfl(kd) were almost exclusive |
| | | 5.3.3(e) | Expression profile of DE genes based on the landmark genes in adult and neonatal male germ cells |
| 5.4 | Discussio | on | |
| 5.5 | Limitatio | n of Study | 259 |
| 5.6 | Summary | , | |
| | PTER 6 ING PRO | | CTERISATION OF CTCFL GENOME-WIDE263 |
| 6.1 | Introducti | ion | |
| 6.2 | Materials | and Meth | ods |
| | 6.2.1 | Preparatio | on and sequencing of ChIP-seq samples266 |
| | | 6.2.1(a) | In silico analysis of transcriptional regulators to differentially expressed (DE) genes using LISA266 |
| | | 6.2.1(b) | ChIP preparation of ESC-CTCFL(DMSO) and ESC-CTCFL(DOX) |
| | | 6.2.1(c) | Validation of protein enrichment to the target sites269 |

| | | 6.2.1(d) | Preparation of ChIP-seq libraries and sequencing269 |
|-----|---------|-----------|--|
| | 6.2.2 | - | rocessing and enrichment analysis for ChIP-seq |
| | | 6.2.2(a) | Reads alignment to reference genome272 |
| | | 6.2.2(b) | Enrichment analysis of CTCFL, CTCF and SUZ12274 |
| | | 6.2.2(c) | Quality assessment of enrichment |
| | 6.2.3 | - | downstream analysis for CTCFL, CTCF and |
| | | 6.2.3(a) | The intersection of CTCF and CTCFL peak sets275 |
| | | 6.2.3(b) | Read density score and clustering of peak sets275 |
| | | 6.2.3(c) | SUZ12 binding sites for downstream analysis276 |
| | | 6.2.3(d) | Binding peak sets annotation and differential enrichment analysis |
| | | 6.2.3(e) | Functional annotation analysis for peak sets277 |
| | 6.2.4 | | the regulatory mode CTCFL toward the ally expressed (DE) genes |
| | | 6.2.4(a) | CTCFL, CTCF and SUZ12 binding sites at the promoters of DE genes in ESC-CTCFL(DOX)277 |
| | | 6.2.4(b) | Histone enrichment analysis |
| | | 6.2.4(c) | CTCF and/or CTCFL binding sites at the promoters of DE genes in JK1-CTCFL(DOX) and GSC-Ctcfl(kd) |
| 6.3 | Results | ••••• | |
| | 6.3.1 | | libraries for 3×FLAG-CTCFL (CTCFL), CTCF and ESC-CTCFL |
| | 6.3.2 | Prelimina | ary ChIP-seq analysis of CTCFL, CTCF and SUZ12 283 |
| | 6.3.3 | Analysis | of CTCFL, CTCF and SUZ12 enrichment289 |
| | | 6.3.3(a) | CTCF and CTCFL co-occupy the gene promoters289 |
| | | 6.3.3(b) | SUZ12 enrichment in ESC-CTCFL(DOX)298 |
| | | 6.3.3(c) | SUZ12 enrichment intersects with CTCF and CTCFL binding sites |

| | 6.3.4 | | the regulatory mode CTCFL towards the ally expressed (DE) genes | 306 |
|-------|-----------|------------|--|-----|
| | | 6.3.4(a) | CTCFL, CTCF and SUZ12 binding sites dominate the promoter of upregulated genes in ESC-CTCFL(DOX) | 306 |
| | | 6.3.4(b) | Bivalent chromatin enriched the intersected CTCFL, CTCF and SUZ12 binding sites in ESCs | 314 |
| | | 6.3.4(c) | CTCF and/or CTCFL binding sites at the promoter of DE genes in JK1-CTCFL(DOX) and GSC-Ctcfl(kd) | 320 |
| 6.4 | Discussio | on | | 324 |
| 6.5 | Limitatio | n of Study | · | 330 |
| 6.6 | Summary | / | | 332 |
| СНАН | PTER 7 | GENERA | AL DISCUSSION AND CONCLUSION | 334 |
| 7.1 | General I | Discussion | | 334 |
| 7.2 | Conclusio | on | | 341 |
| 7.3 | Direction | s for Futu | re Research | 342 |
| REFE | RENCES | | | 345 |
| A PPF | NDICES | | | |

LIST OF TABLES

| | Page |
|------------|---|
| Table 2.1 | Transcriptional role and localisation of histone modifications |
| | (Gates et al., 2017; Zhao and Garcia, 2015) |
| Table 3.1 | The composition of 200 ml basal medium for the culture of |
| | GSCs74 |
| Table 3.2 | Preparation of DNase digestion reaction mixture84 |
| Table 3.3 | Preparation of reverse transcription reaction set up86 |
| Table 3.4 | Reaction mixture setup for TaqMan® probe qPCR assay88 |
| Table 3.5 | The two-step qPCR cycle conditions for TaqMan® probe qPCR |
| | assay89 |
| Table 3.6 | Reaction mixture setup for SYBR® Green qPCR assay91 |
| Table 3.7 | The three-step qPCR cycle conditions for SYBR® Green qPCR |
| | assay92 |
| Table 3.8 | Preparation of solutions for immunofluorescence staining96 |
| Table 3.9 | List of solutions for lysate extraction, PAGE, Western blotting |
| | and antibody stripping98 |
| Table 3.10 | Dilutions of Spike Mix for Cy3-labelling using 50 ng total RNA105 |
| Table 3.11 | Preparation of master mix and reaction mixture for cDNA |
| | synthesis108 |
| Table 3.12 | Preparation of master mix and reaction mixture for transcription109 |
| Table 3.13 | Preparation of fragmentation and hybridisation reaction mixtures114 |
| Table 3.14 | Essential solutions for Smart-seq2 library preparation119 |
| Table 3.15 | Preparation of RT mix |
| Table 3.16 | Cycling conditions for reverse transcription123 |
| Table 3.17 | Components and preparation of preamplification PCR mix125 |

| Table 3.18 | Cycling conditions for preamplification PCR126 |
|------------|---|
| Table 3.19 | The tagmentation reaction mix |
| Table 3.20 | The preparation of Index adapter-ligated fragments130 |
| Table 3.21 | Amplification of Index adapter-ligated fragments by PCR131 |
| Table 3.22 | List of solutions prepared for ChIP133 |
| Table 3.23 | The relative quantification was calculated using the two-step Percent Input Method, which normalised both background levels and input chromatin used for the ChIP |
| Table 3.24 | Reaction setups for end-repair, A-tailing and adapter ligation were adapted from the manufacturer's protocol |
| Table 3.25 | Reaction mixture setup for cycle determination |
| Table 3.26 | The protocol for quantifying enriched DNA fragments using SYBR® Green qPCR assay |
| Table 3.27 | Library amplification reaction setup |
| Table 3.28 | Amplification of enriched ChIP DNA fragments by PCR147 |
| Table 3.29 | The setup of the reaction mixture for library quantification150 |
| Table 3.30 | The protocol for library quantification |
| Table 3.31 | The calculation of library concentration using the quantification data and the average fragment length |
| Table 4.1 | Restriction enzyme digestion mixture |
| Table 4.2 | Plasmids and reagents used for transfection mixture preparation161 |
| Table 4.3 | List of samples prepared for microarray hybridisation171 |
| Table 4.4 | The top ten highly significant GO terms (p-value < 0.01) associated with upregulated genes in ESC-CTCFL(DOX)190 |
| Table 4.5 | The highly significant GO terms (p-value < 0.01) associated with downregulated genes in ESC-CTCFL(DOX) |
| Table 4.6 | The highly significant GO terms (p-value < 0.01) associated with upregulated genes in IK1-CTCFL (DOX) |

| Table 4.7 | The significant GO terms (p-value < 0.05) associated with downregulated genes in JK1-CTCFL(DOX) |
|-----------|--|
| Table 5.1 | The preparation of transfection mixture to deliver the DsiRNA constructs to the GSCs |
| Table 5.2 | The list of samples and combinations of indexes used in tagging for library amplification |
| Table 5.3 | Mapping summary for Smart-seq2 libraries using HISAT2 displays the high alignment rate (> 90%) from each library to the reference genome. |
| Table 5.4 | Functional annotation clusters of biological processes identified by DAVID for the upregulated genes |
| Table 5.5 | The top five functional annotation clusters of biological processes identified by DAVID for the downregulated genes |
| Table 5.6 | List of significant KEGG pathways associated with the downregulated genes |
| Table 5.7 | Clusters of biological processes associated with downregulated genes in Cluster 1 (differentiation) and Cluster 4 (SSC) |
| Table 6.1 | The preparation of immunocomplexes for a single condition268 |
| Table 6.2 | List of samples and their parameters for library preparation271 |
| Table 6.3 | Mapping summary and the alignment rate for ChIP-seq libraries284 |
| Table 6.4 | Number of genes with the three classes of CTCF and CTCFL binding sites at the respective genomic features |
| Table 6.5 | The top five biological processes associated with the three classes of CTCF and/or CTCFL binding sites |
| Table 6.6 | The top ten significant biological processes associated with SUZ12(H3K27me3+) loci |
| Table 6.7 | The top five biological processes associated with the three classes of CTCF and CTCFL binding sites that intersected with SUZ12 enriched sites |

| Table 6.8 | The top five biological processes involved with the encoded |
|-----------|--|
| | CTCF&CTCFL binding sites, grouped in Cluster 1 and Cluster 2 |
| | loci sets 318 |

LIST OF FIGURES

| | Page |
|------------|--|
| Figure 1.1 | The overall research design conducted in this study was based on three specific objectives |
| Figure 2.1 | Levels of chromatin organisation |
| Figure 2.2 | The regulators of histone modifications |
| Figure 2.3 | The formation of chromatin loops and CTCF binding motifs. A. CTCF and cohesin generate chromatin loops through loop extrusion and form TAD. The image was adapted from Pongubala and Murre (2021). B. The combination of CTCF ZFs binds to the core and secondary binding motifs. The image was adapted from Wu <i>et al.</i> (2020). |
| Figure 2.4 | The compositions of PRC2, including the core and accessory proteins. PRC2 exists in two distinct multimeric forms, PRC2.1 and PRC2.2. The figure was adapted from Mortimer (2019)30 |
| Figure 2.5 | The interactions of CTCF and PRC2 in normal development. A. CTCF is required for Polycomb-mediated imprinting of the Igf2/H19 locus in murine cells and directly interacts with SUZ12 (Li <i>et al.</i> , 2008). B. CTCF form boundaries to establish the polycomb domain. The domains house the polycomb-bound developmental genes (Dowen <i>et al.</i> , 2014). C. CTCF forms a transcriptional co-repressor complex with SUZ12-PRC2 to block gene expression that can inhibit the Schwann cell differentiation (Wang <i>et al.</i> , 2020a). |
| Figure 2.6 | The molecular architecture of CTCFL protein. The three domains of CTCF and CTCFL are N-terminal, ZF DNA binding domain and C-terminal domain. The illustrations were adapted from Klenova <i>et al.</i> (2002) and Hore <i>et al.</i> (2008) |

| Figure 2.7 | and single 1×CTS. The three classes of CTCF and/or CTCFL binding sites are CTCF&CTCFL, CTCF-only and CTCFL-only sites. The figure was adapted from Pugacheva <i>et al.</i> (2015)43 |
|-------------|--|
| Figure 2.8 | CTCFL highest expression in Type B spermatogonia and primary spermatocytes (red box). The schematic overview of the main stages in mouse spermatogenesis was adapted from Ernst <i>et al.</i> (2017) |
| Figure 2.9 | The expression of <i>Ctcfl</i> in the adult mice testis using single-cell expression profiling. A. The t-distributed stochastic neighbour embedding (t-SNE) plot displays the highest level of <i>Ctcfl</i> expression in spermatogonia and leptotene spermatocytes. The plot was adapted from Jung <i>et al.</i> (2019a) using Mouse Testis Single-cell RNA-seq Atlas Online Tool. B. Single-cell RNA-seq analysis of <i>Ctcf</i> and <i>Ctcfl</i> expression at the significant stages of mouse spermatogenesis, SG: Spermatogonia, SPC: Spermatocytes, RS: Round Spermatids. The data was acquired from Green <i>et al.</i> (2018) and reanalysed and plotted by Rivero-Hinojosa <i>et al.</i> (2021) |
| Figure 2.10 | Workflow of sample preparation for Agilent array processing, adapted from Agilent product package insert |
| Figure 2.11 | An overview of the experimental steps in the RNA-seq protocol. The figure was adapted from Van den Berge <i>et al.</i> (2019)56 |
| Figure 2.12 | The ChIP assay for preparation of immune-precipitate material and various analysis methods. The figure was adapted from Collas (2010) |
| Figure 2.13 | The ChIP-seq analysis workflow includes the sample preparation, sequencing, and computational analysis. The figure was adapted from Nakato and Shirahige (2016) and Nakato and Sakata (2020)65 |
| Figure 3.1 | The GSCs were derived from self-renewal, undifferentiated spermatogonia or SSCs from neonatal mouse testis. The image was adapted from Kanatsu-Shinohara and Shinohara (2013)71 |

| Figure 3.2 | The formula A. the double delta Ct ($\Delta\Delta$ Ct) method and B. Pffafl |
|------------|--|
| | method. Either method was used for relative quantification, |
| | depending on the E value94 |
| Figure 3.3 | The formula for calculating the cRNA yield and reaction activity112 |
| Figure 3.4 | The flowchart for Smart-seq2 library preparation, adapted from Picelli <i>et al.</i> (2014b) |
| | Ficeiii et al. (20146)117 |
| Figure 3.5 | Dilution of the library for quantification. The figure was adapted |
| | from the KAPA Library Quantification protocol guide149 |
| Figure 4.1 | The flow chart illustrates the experimental design conducted in this chapter |
| Figure 4.2 | The pPB-3×FLAG-CTCFL expression plasmid contained the |
| C | 3×FLAG-CTCFL recombinant construct and puromycin resistant |
| | gene (puro ^{R)} . Transactivation of TRE promoter by doxycycline |
| | occurred for conditional expression of 3×FLAG-CTCFL in the |
| | transfected cell line |
| Figure 4.3 | The experimental design for transfection of the pPB-3×FLAG- |
| | CTCFL plasmid into the ESCs and JK1 cells |
| Figure 4.4 | The setup for immunofluorescence staining. Four-chambered 35 |
| | mm dishes were used for ESC-CTCFL and JK1-CTCFL165 |
| Figure 4.5 | Preparation for immunoblotting. A. The setup for cell cultures in |
| | (DMSO) and (DO)X conditions using the 100 mm dishes. B. The |
| | setup for sample loading in the pre-cast gel. M, protein ladder |
| | marker, R1, replicate one, R2, replicate two |
| Figure 4.6 | Ectopic 3×FLAG-CTCFL expression plasmid. The gel |
| | electrophoresis result shows the size of the $3 \times Flag$ -Ctcfl |
| | construct near 2 kb (2,007 bp) |
| Figure 4.7 | Immunofluorescence staining for analysis of 3×FLAG-CTCFL |
| | expression probed by anti-FLAG in ESC-CTCFL(DMSO) and |
| | ESC-CTCFL(DOX) cells. The images were visualised via $200\times$ |
| | magnification using a fluorescence microscope |

| Figure 4.8 | Immunofluorescence staining for analysis of 3×FLAG-CTCFL |
|-------------|---|
| | expression probed by anti-FLAG in JK1-CTCFL(DMSO) and |
| | JK1-CTCFL(DOX) cells. The images were visualised via 200× |
| | magnification using a fluorescence microscope |
| Figure 4.9 | The total cell fluorescence for anti-FLAG against 3×FLAG- |
| | CTCFL expression in (DMSO) and (DOX) conditions for ESC- |
| | CTCFL $(n = 6)$ and JK1-CTCFL $(n = 5)$. The differences of |
| | fluorescence signal between (DMSO) and (DOX) conditions |
| | were analysed by Student's T-test (** p-value < 0.01)179 |
| Figure 4.10 | Immunoblotting of ESC-CTCFL and JK1-CTCFL cells. The |
| 8 | expression of the 3×FLAG-CTCFL protein was assessed in |
| | (DMSO) and (DOX) conditions. The anti-FLAG antibody probed |
| | the 3×FLAG-CTCFL protein. The α-TUBULIN was used as the |
| | loading control |
| Figure 4.11 | The relative <i>Ctcfl</i> expression by RT-qPCR. <i>Ctcfl</i> was |
| rigule 4.11 | significantly expressed in ESC-CTCFL(DOX) and JK1- |
| | CTCFL(DOX) cells. The comparison between (DMSO) and |
| | (DOX) conditions was analysed using T-test (** p-value < 0.01 , n |
| | = 3). |
| E' 4.10 | |
| Figure 4.12 | The Venn diagrams display the intersection of upregulated and |
| | downregulated genes between ESC-CTCFL(DOX) and JK1- |
| | CTCFL(DOX). Common DE genes in both cell lines were listed |
| | below the Venn diagrams, and <i>Ctcfl</i> was listed as an upregulated |
| | gene for ESC-CTCFL(DOX) but not JK1-CTCFL(DOX)184 |
| Figure 4.13 | Significant induction of Ctcfl 3' UTR in ESC-CTCFL(DOX) |
| | analysed by Student's T-test (*p-value < 0.05, n = 5) while Ctcfl |
| | expression in JK1-CTCFL(DOX) was negligible (n = 3)186 |
| Figure 4.14 | Validation of several upregulated genes upon 3×FLAG-CTCFL |
| | expression by RT-qPCR analysis. A. The $Prss50$ (n = 4), $Stra8$ (n |
| | = 3), $Adgrg1$ (n $=$ 3), $Cited1$ (n $=$ 4) and $Plekhg4$ (n $=$ 4) were |
| | significantly upregulated in ESC-CTCFL(DOX). B. In JK1- |
| | CTCFL(DOX), the expression of $Adgrg1$ (n = 3), $Cited1$ (n = 3) |

| | and $Rimbp3$ (n = 4) were significantly induced. The differences between (DMSO) and (DOX) conditions were analysed by Student's T-test, (*p < 0.05, **p < 0.01) |
|-------------|---|
| Figure 4.15 | The scatterplot by REVIGO visualised the clusters of summarised GO terms for the significant biological process associated with upregulated genes in ESC-CTCFL(DOX). The arrow indicates the cluster of biological processes for development |
| Figure 4.16 | The expression profile of DE genes in ESC-CTCFL(DOX) in adult tissues. Hierarchical gene clustering heatmaps for A. upregulated, and B. downregulated genes displayed the gene clusters associated with tissue-selective expression. The similarity matrix plot analysed the expression profile similarity for the C. upregulated and D. downregulated genes in the adult tissues. |
| Figure 4.17 | Hierarchical clustering of upregulated genes in JK1-CTCFL(DOX) compared to other cell types. The black line indicates the upregulated genes (n = 34) overlapped in JK1-CTCFL(DOX) and undifferentiated ESC-CTCFL(DMSO) |
| Figure 4.18 | The expression profile of DE genes in JK1-CTCFL(DOX) in adult tissues. Hierarchical gene clustering heatmaps for A. upregulated, and B. downregulated genes displayed the gene clusters associated with tissue-selective expression. The similarity matrix plot analysed the expression profile similarity for the C. upregulated, and D. downregulated genes in the adult tissues. |
| Figure 4.19 | The intersection of CTA genes with the upregulated genes in ESC-CTCFL(DOX) and JK1-CTCFL(DOX) cells203 |
| Figure 4.20 | CTCFL aberrant expression in undifferentiated ESCs and differentiated JK1 cells altered the regulation of developmental genes and processes |
| Figure 5.1 | The experimental design conducted in this chapter 214 |

| Figure 5.2 | The clusters of GSCs were continuously expressing the EGFP. The images were visualised via 100× and 200× magnification using a fluorescence microscope |
|------------|--|
| Figure 5.3 | The analysis pipeline for RNA-seq libraries224 |
| Figure 5.4 | Visual assessment of GSCs for transfection of TYE 563 labelled DsiRNA as transfection control. A. The non-transfected GSCs did not emit the red fluorescent due to the absence of TYE 563. B. The TYE 563 transfected GSCs emitted the red fluorescent signal. The yellow signal indicated the overlapping of the green GSCs, expressing EGFP and red fluorescent of TYE 563. The images were visualised via 100× and 200× magnification using a fluorescence microscope |
| Figure 5.5 | The preliminary assessment of gene knockdown by DsiRNAs for 24 hours. A. The relative expression of <i>Hprt1</i> in response to the positive control DsiRNA(<i>Hprt1</i>). B. The relative expression of <i>Ctcfl</i> was assessed upon the incubation with DsiRNA(<i>Ctcfl</i>)-1, DsiRNA(<i>Ctcfl</i>)-2, and DsiRNA(<i>Ctcfl</i>)-3, respectively. The experiments were conducted in a single replicate |
| Figure 5.6 | The experimental design for preparation of <i>Ctcfl</i> knockdown GSCs. The knockdown using co-transfected DsiRNA(Ctcfl)-1+2 and negative control samples were prepared in three replicates. The DsiRNA incubation was conducted for 24 and 48 hours. TYE 563 served as the transfection control |
| Figure 5.7 | Visual assessment of GSCs for transfection of TYE 563 labelled DsiRNA as transfection control. A. The non-transfected GSCs did not emit the red fluorescent due to the absence of TYE 563. B. The TYE 563 transfected GSCs emitted the red fluorescent signal. The yellow signal indicated the overlapping of the green GSCs, expressing EGFP and red fluorescent of TYE 563. The images were visualised via 100× magnification using a fluorescence microscope. |

| Figure 5.8 | Relative expression of <i>Ctcfl</i> upon knockdown treatment. The expression was compared between the negative control and DsiRNA(<i>Ctcfl</i>)-1+2 samples in three replicates. Statistical analysis was done using the Student's T-test for significant changes (** p-value < 0.01) |
|-------------|---|
| Figure 5.9 | The two-dimensional PCA plot for the transformed counts from the samples in three replicates. 'Replicate' was represented by the shape, and the 'condition' was represented by orange for knockdown GSC- <i>Ctcfl</i> (kd) and green for negative control GSC(Neg). Samples for R1 and R2 were clustered together for each condition but not R3. |
| Figure 5.10 | The differential expression of <i>Ctcfl</i> analysed by A. DESeq2 for RNA-seq using two replicates and B. qPCR using three replicates show the reduction of <i>Ctcfl</i> expression in GSC- <i>Ctcfl</i> (kd). Statistical analysis for the differential <i>Ctcfl</i> expression by qPCR was done using the Student's T-test (*** p-value < 0.001)238 |
| Figure 5.11 | The volcano plot for the distribution of coding genes and differentially expressed (DE) genes |
| Figure 5.12 | A. Relative gene expression by qPCR for the downregulated genes in the negative control, Neg and <i>Ctcfl</i> knockdown, <i>Ctcfl</i> (kd) GSCs. Statistical analysis by Student's T-test using three replicates (*** p-value < 0.001, ** p-value < 0.01, * p-value < 0.05). B. The comparison of log2 fold change of <i>Mmp3</i> , <i>Mmp9</i> , <i>Hmox1</i> , <i>Pdgfrb</i> , <i>Cepbp</i> and <i>Tmem148b</i> were determined by RNA-Seq analysis (black) and qPCR validation (blue). Spry2 was a non-DE gene |
| Figure 5.13 | A. The intersection of DE genes in GSC- <i>Ctcfl</i> (kd) with upregulated genes in ESC-CTCFL(DOX) and JK1-CTCFL(DOX). B. The intersection between the DE genes in GSC- <i>Ctcfl</i> (kd) and CTA genes (Wang <i>et al.</i> , 2016) |
| Figure 5.14 | The Venn diagrams display the intersection counts of DE genes in A. GSC- <i>Ctcfl</i> (kd), B. ESC-CTCFL(DOX) and C. JK1- |

| | CTCFL(DOX) with the landmark genes in spermatogonia (Spg), spermatocytes (Sct) and spermatids (Std)249 |
|-------------|---|
| Figure 5.15 | The radar plot for the distribution of DE genes in GSC- <i>Ctcfl</i> (kd) overlapped with the landmark genes in each cluster of P6 spermatogonia |
| Figure 5.16 | CTCFL-downstream regulation of the biological processes for spermatogonia development may poise the differentiation of spermatogonia |
| Figure 6.1 | The experimental design conducted in this chapter265 |
| Figure 6.2 | Flowchart of pipeline used for reads processing, enrichment analysis, quality assessment and visualisation of ChIP-seq libraries |
| Figure 6.3 | The top 20 predicted transcription regulators of the DE genes in ESC-CTCFL(DOX). A. The significant regulators in upregulated genes and B. downregulated genes. Yellow bars indicate the top ten transcription factors marked in the scatter plot (Appendix A – Supplementary result A16). The bars with thick outlines indicate regulators in which ChIP-seq peaks data were derived from ESCs |
| Figure 6.4 | The average relative enrichment of CTCFL, CTCF and SUZ12 to their respective binding sites against the negative control, IG3. The relative quantification was measured using the Percent Input method in two biological replicates (n = 2). CTCF and CTCFL enriched their binding sites that encoded for <i>Prss50</i> and <i>Adgrg1</i> . SUZ12 enriched its target sites near <i>Pchdh8</i> |
| Figure 6.5 | Quality assessment by ChIPQC for the ChIP-seq libraries (CTCFL, CTCF and SUZ12) and their replicates based on genomic enrichment in Chromosome 10 |
| Figure 6.6 | Hierarchical clustering based on Pearson correlation coefficient compute the similarity between the ChIP-seq libraries (CTCFL, CTCF and SUZ12) and their replicates |

| Figure 6.7 | The merged peaks of replicates corresponding to CTCF and |
|-------------|---|
| | CTCFL. The shared peaks between replicates for A. CTCF and |
| | B. CTCFL were used for downstream analysis. C. The 'true' |
| | peaks were retrieved for CTCFL after discarding the shared peaks |
| | in (DOX) and (DMSO) ChIP libraries288 |
| Figure 6.8 | CTCF and CTCFL enrichment in ESC-CTCFL(DOX). A. The intersection of CTCF and CTCFL peaks generated three classes of binding sites. The binding of CTCF and CTCFL at these sites was verified through B. the enrichment signal heatmap and C. mean read density score. D. The binding of CTCF and CTCFL at <i>Prss50</i> promoter was viewed using IGV genome browser290 |
| Figure 6.9 | Genomic distribution and binding motif of CTCF and CTCFL in ESC-CTCFL(DOX). ChIPseeker annotated the peak sets of A. CTCF and B. CTCFL peak sets. The genomic features are untranslated regions (UTR), promoter, intron, exon, downstream and distal intergenic regions. The sequence motif for CTCF and CTCFL were derived from MEME-ChIP motif discovery |
| Figure 6.10 | (Machanick and Bailey, 2011) |
| Figure 6.11 | The differential CTCF enrichment at CTCF binding sites in control ESC-CTCFL(DMSO) and CTCFL-expressing ESC-CTCFL(DOX). A. The volcano plot displays the CTCF-induced (Log FC > 0) and CTCF-reduced sites (Log FC < 0) with FDR ≤ 0.05 . B. The mean read or enrichment score of CTCF for CTCF-induced and CTCF-reduced sites between the two conditions297 |
| Figure 6.12 | Enrichment of SUZ12 in ESC-CTCFL. A. The IGV tracks show the enrichment of SUZ12 to the targeted loci near <i>Pcdh8</i> and <i>Hoxa</i> gene clusters. B. The correlation heatmap of SUZ12 enrichment in this study with SUZ12 binding in published data sets (a) (Riising <i>et al.</i> , 2014) (B) (Kloet <i>et al.</i> , 2016). C. The |

| | genomic distribution of annotated SUZ12-bound sites in ESC- CTCFL |
|-------------|---|
| Figure 6.13 | The enrichment of H3K27me3 to SUZ12 recruited sites. The enrichment intensity of H3K27me3 in ESCs was used from a published report (Riising <i>et al.</i> , 2014). A. The k-mean clustering for the enrichment of H3K27me3 to SUZ12 recruited sites resulted in two clusters: SUZ12(H3K27me3+) and SUZ12(H3K27me3-). B. The pie charts for the annotation of genomic features for both clusters by ChIPseeker |
| Figure 6.14 | The co-occupancy of CTCFL, CTCF and SUZ12. A. BARTweb analysis discovered CTCF and SUZ12 as the top transcriptional regulators associated with CTCFL-bound sites. The regulators were ranked based on Irwin-Hall P-value (p-value < 5E-02, 0.05). B. The intersections of SUZ12 enriched sites with the three classes of CTCF and CTCFL binding sites |
| Figure 6.15 | Upset plots of the intersection of A. upregulated, and B. downregulated genes in ESC-CTCFL(DOX) with the annotated genes that contained CTCFL, CTCF and SUZ12 binding sites at the promoter. (≤ 3 kb, upstream and downstream from TSS)307 |
| Figure 6.16 | The average intensity profile plots for SUZ12 enrichment to the recruited sites in control (DMSO) and aberrant CTCFL expression (DOX). |
| Figure 6.17 | The IGV tracks show the enrichment of CTCFL, CTCF and SUZ12 at the gene promoters (<i>Ctcfl, Ednra, Tgf\beta1</i> and <i>Prss50</i>), which contained the overlapping CTCF&CTCFL binding sites. These genes also displayed a reduction of SUZ12 enrichment at their promoter. The bar graph indicates the log ₂ ratio for SUZ12 (DOX) against SUZ12 (DMSO). The arrow indicates the TSS of the genes |
| Figure 6.18 | A. The IGV tracks show the co-enrichment of CTCF and CTCFL at CTCF&CTCFL sites, but not SUZ12 at the gene promoters (<i>Stra8</i> and <i>Nudt17</i>). B. The IGV tracks show the promoters that |

| | (Batf3) and the promoter that also had CTCFL-only sites (Dennd2d). The bar graph indicates the log ₂ ratio for SUZ12 (DOX) against SUZ12 (DMSO). The arrow indicates the TSS of the genes |
|-------------|---|
| Figure 6.19 | A. The IGV tracks show the enrichment of CTCF at CTCF-only sites at the gene promoters (<i>Foxa1</i> and <i>Wnt4</i>). B . The IGV tracks show the promoters that CTCFL-only sites at the gene promoter of <i>Cited1</i> and <i>Mapk12</i> . SUZ12 also enriched these promoters. Based on the log2 ratio of SUZ12 enrichment, SUZ12 occupancy at <i>Cited1</i> promoter was significantly reduced. The bar graph bar indicates the log2 ratio for SUZ12 (DOX) against SUZ12 (DMSO). The arrow indicates the TSS of the genes |
| Figure 6.20 | The enrichment score and heatmaps of H3K4me3 and H3K27me3 on A. CTCF&CTCFL binding sites and B. CTCF&CTCFL+SUZ12 binding sites in ESCs. The enrichment score was log ₂ transformed |
| Figure 6.21 | The IGV tracks show the upregulated genes associated with A. Cluster 1 and B. Cluster 2. The promoters contained CTCF&CTCFL binding sites with distinct enrichment of SUZ12, H3K4me3 and H3K27me3 |
| Figure 6.22 | Upset plots of the intersection of A. upregulated and B. downregulated genes in JK1-CTCFL(DOX) with the annotated genes that contained CTCF and/or CTCFL binding sites at the promoter (promoter region: \leq 3 kb, upstream and downstream from TSS) |
| Figure 6.23 | Upset plots of the intersection of A. upregulated and B. downregulated genes in GSC- $Ctcfl(kd)$ with the annotated genes that contained CTCF and CTCFL binding sites at the promoter (promoter region: ≤ 3 kb, upstream and downstream from TSS)323 |
| Figure 6.24 | Based on the findings in this study, two CTCFL regulating modes are proposed. In the first mode, CTCF and CTCFL might co- |

| enrich the gene promoter at CTCF&CTCFL sites and abnormally | |
|---|-----|
| activate the germline gene, irrespective of PRC2 and bivalent | |
| chromatin. The alteration of germline repression is unclear. In the | |
| second mode, abnormal CTCFL enrichment near the promoter | |
| may obstruct CTCF-PRC2 repression of the developing gene, | |
| driving to aberrant gene expression. | 333 |

LIST OF SYMBOLS

% Percentage

 $1\times$ One time dilution

10× Ten-time dilution

 $\times 10^4$ Ten raised to a power of four

bp Base pair(s)

kb Kilobase pair

μg Microgram

μl Microlitre

μM Micromolar

M Molar/molarity

mL Millilitre

mm Millimetre

mM Millimolar

ng Nanogram

nM Nanomolar

°C Degree Celsius

pM Picomolar

pmol Picomol

rpm Revolutions per minute

V Volt

w/v Mass per volume

LIST OF ABBREVIATIONS

(DMSO) DMSO treated cells

(DOX) Doxycycline treated cells,

(3×FLAG-CTCFL expressed cells)

3' UTR 3' Untranslated regions

BSA Bovine serum albumin

cDNA Complementary deoxyribonucleic acid

CDS Coding sequence

ChIP Chromatin immunoprecipitation

ChIP-seq Chromatin immunoprecipitation and sequencing

cRNA Complementary ribonucleic acid

CTA Cancer testis antigen
CTCF CCCTC-binding factor

Ctcf CCCTC-binding factor (gene)
CTCFL CCCTC-binding factor like

Ctcfl CCCTC-binding factor-like (gene)

CTCF&CTCFL DNA sequences targeted by CTCF and CTCFL

CTCF-only DNA sequences targeted by CTCF only
CTCFL-only DNA sequences targeted by CTCFL only

DE gene Differentially expressed gene

DMSO Dimethyl sulfoxide

DNA Deoxyribonucleic acid

DOX Doxycycline

dsDNA Double-stranded DNA

DsiRNA Dicer-substrate short interfering RNAs

E Embryonic day

E Efficiency

EGFP Enhanced Green Fluorescent Protein

ESCs Embryonic stem cells
ESC-CTCFL ESC-CTCFL cell line

ESC-CTCFL(DMSO) Control ESC-CTCFL cell line

(repressed CTCFL, DMSO treated cells)

ESC-CTCFL (DOX) ESC-CTCFL cell line with 3×FLAG-CTCFL expression

(doxycycline-treated cells)

FBS Fetal bovine serum

FC Fold-change

FDR False discovery rate

GO Gene ontology

GSCs Germline stem cells

GSC(Neg) Control GSCs (without *Ctcfl* knockdown)

GSC-*Ctcfl*(kd) *Ctcfl* knockdown GSCs

HF High fidelity

H3K4me3 Histone-3-lysine-4-trimethylation H3K27me3 Histone-3-lysine-27-trimethylation

JK1 JK1 testicular stromal cell line

JK1-CTCFL JK1-CTCFL cell line

JK1-CTCFL(DMSO) Control JK1-CTCFL cell line

(repressed CTCFL, DMSO treated cells)

JK1-CTCFL (DOX) JK1-CTCFL cell line with 3×FLAG-CTCFL expression

(doxycycline-treated cells)

kb Kilo base pairs kd Knockdown kDa Kilo Dalton

KEGG Kyoto Encyclopedia of Genes and Genomes

KSR Knockout serum replacement
MEFs Mouse embryonic fibroblasts

mRNA Messenger RNA

n Number/count

n Number of replicates

neg Negative control

NGS Next-generation sequencing

P Postnatal day

p Passage number

PAGE Polyacrylamide gel electrophoresis

PBS Phosphate buffered saline

PCA Principle component analysis

PCR Polymerase chain reaction

PRC2 Polycomb repressor complex 2

qPCR Quantitative polymerase chain reaction

R Replicate

Rep Replicate

RIN RNA Integrity Number

RNA Ribonucleic acid
RNA-seq RNA sequencing
RNAi RNA interference

RT-qPCR Reverse transcription-quantitative polymerase chain

reaction

SSCs Spermatogonial stem cells

SUZ12 Suppressor of zeste 12 protein

TNF Tumour necrosis factor
TSS Transcriptional start sites

ZF Zinc finger

LIST OF APPENDICES

Appendix A Supplementary results

Appendix B List of materials

Appendix C List of plasmids and oligonucleotides

Appendix D Preparation of reagents and solutions

Appendix E Bioinformatics tools

Appendix F Protocols from manufacturer

PEMPROFILAN MOLEKUL DAN PENGAWALATURAN TERHADAP FAKTOR TRANSKRIPSI TIKUS CTCFL DALAM SEL BUKAN GERMA DAN SEL STEM GERMA JANTAN

ABSTRAK

Faktor pengikat CCCTC khusus testis (CTCFL) adalah faktor transkripsi yang diekspresi dalam sel germa jantan dan penting untuk penghasilan sperma. CTCFL ialah paralog kepada faktor pengikat CCCTC (CTCF), faktor transkripsi dan protein pengatur-bina kromatin. Mereka berkongsi domain pengikat DNA jejari-zink yang serupa tetapi mempunyai dua domain hujung yang berbeza. CTCFL tidak ekspres dalam tisu somatik; walau bagaimanapun, ekspresi yang abnormal dalam sel bukan germa dikaitkan dengan kanser. Fungsi pengawalseliaan CTCFL dalam penghasilan sperma dan pertumbuhan tumor masih tidak jelas. Kajian ini bertujuan untuk menjelaskan profil pengawalseliaan CTCFL murin dalam sel bukan germa dan sel germa jantan. Ekspresi ektopik CTCFL bertanda FLAG (3×FLAG-CTCFL) daktifkan dalam sel bukan germa, iaitu sel induk embrio tikus (ESCs) dan sel stroma testis JK1 (JK1) selama 24 Jam. Untuk mengetahui peranan CTCFL dalam sel germa, ekspresi Ctcfl endogen dalam sel stem germanium jantan (GSCs) telah dikurangkan melalui gangguan RNA selama 48 jam. Perubahan global dalam transkrip telah dianalisis menggunakan microarray dan penjujukan RNA. Analisis anotasi fungsi mendapati gen dan proses perkembangan telah diganggu apabila CTCFL diekspresi dalam ESC, sebagai ESC-CTCFL(DOX) dan sel JK1, sebagai JK1-CTCFL(DOX). Pengurangan ekspresi Ctcfl dalam GSCs menindas ekspresi gen yang terlibat dalam peraturan kematian sel dan proses selular dalam spermatogonia. Mod pengawalseliaan CTCFL telah disimpulkan berdasarkan profil ikatan protein-DNA dalam genom ESC-CTCFL(DOX). Analisis awal in silico menganggarkan

perkaitan utama antara gen-gen yang mempunyai pertambahan ekspresi dalam ESC-CTCFL(DOX) dengan pengayaan CTCFL, CTCF, dan komponen-komponen Polycomb Repressor Complex 2 (PRC2), termasuk Suppressor of zeste 12 protein (SUZ12). Pengayaan CTCFL, CTCF, dan SUZ12 terhadap genom dalam ESC-CTCFL(DOX) telah dianalisis menggunakan analisis imunoprekipitasi kromatin dan penjujukan (ChIP-seq). Hasilnya mengesahkan pertindihan pengikatan CTCFL dengan tapak pengayaan CTCF dan SUZ12 yang mendominasi tapak promotor pengatur perkembangan sel. Pengikatan bersama CTCF dan CTCFL di tapak CTCF&CTCFL boleh memacu pengaktifan gen germa dalam sel bukan germa. Pengayaan SUZ12 telah dikurangkan di tapak CTCF&CTCFL berhampiran promotor gen perkembangan yang ditambah ekpresinya, menyimpulkan perubahan fungsi penindasan PRC2 oleh CTCFL. Ringkasnya, CTCFL mungkin memainkan peranan dalam peraturan perkembangan dengan mengubah suai proses dalam sel bukan germa atau dengan mengawal diferensiasi spermatogonia. Disregulasi aktiviti CTCF dan PRC2 boleh menjadi mod pengawalseliaan CTCFL dalam pertumbuhan barah dan memerlukan penyiasatan lanjut.

MOLECULAR AND REGULATORY PROFILING OF MURINE TRANSCRIPTION FACTOR CTCFL IN NON-GERM CELLS AND MALE GERMLINE STEM CELLS

ABSTRACT

Testis-specific CCCTC-binding factor-like (CTCFL) is a transcription factor expressed in male germ cells and essential for spermatogenesis. CTCFL is the paralog of CCCTC-binding factor (CTCF), a transcription factor and chromatin architectural protein. They share a similar zinc-finger DNA binding domain but different end termini. CTCFL is repressed in somatic tissues; however, the aberrant expression in the non-germ cells is associated with cancer. The regulatory functions of CTCFL in spermatogenesis and tumourigenesis are still unclear. This study aimed to elucidate the regulatory profile of murine CTCFL in non-germ cells and male germ cells. To investigate the transcriptional effects of CTCFL aberrant expression in non-germ cells, ectopic FLAG-tagged CTCFL (3×FLAG-CTCFL) was expressed in the pluripotent mouse embryonic stem cells (ESCs) and JK1 testicular stromal cell line (JK1) for 24 hours. To discover the roles of CTCFL in the germ cells, the expression of endogenous Ctcfl in male germline stem cells (GSCs) was reduced by RNA interference for 48 hours. Global changes in the transcriptome were measured by microarray and RNA sequencing, respectively. The functional annotation analysis observed altered developmental genes and processes in CTCFL expressing ESCs, ESC-CTCFL(DOX), and JK1 cells, JK1-CTCFL(DOX). Ctcfl knockdown in GSCs repressed the genes involved in cell death regulation and cellular processes in spermatogonia. The regulatory mode of CTCFL was inferred based on the genomewide protein-DNA binding profile in ESC-CTCFL(DOX). Preliminary in silico analysis highlighted the association of the upregulated genes in ESC-CTCFL(DOX)

with the enrichment of CTCFL, CTCF, and Polycomb Repressor Complex 2 (PRC2) components, including Suppressor of zeste 12 protein (SUZ12). The enrichment of CTCFL, CTCF, and SUZ12 to the genome in ESC-CTCFL(DOX) were analysed by chromatin immunoprecipitation and sequencing (ChIP-seq). The results validated the intersection of CTCFL binding with CTCF and SUZ12 enrichment sites that dominated the promoter of developmental regulators. CTCF and CTCFL co-binding at CTCF&CTCFL sites could drive the activation of germline genes in the non-germ cells. SUZ12 enrichment was reduced at the CTCF&CTCFL sites near the promoter of upregulated developmental genes, inferring alteration of PRC2 repression by CTCFL. In summary, CTCFL may play a role in developmental regulation by modifying the processes in non-germ cells and regulating spermatogonia's differentiation. CTCF and PRC2 activity dysregulation could be the CTCFL regulatory tumourigenesis further investigations. mode in and warrant

CHAPTER 1

INTRODUCTION

1.1 Background of Study

Transcription factors are one of the main components in the transcriptional machinery regulating gene transcription into mRNA. In turn, the gene expression programs convey cell-fate decisions for cellular activity and development. Transcription factors act by binding to specific DNA motifs within gene regulatory elements (Stadhouders *et al.*, 2019). The sequence-specific transcription factors may become the activator or repressor that regulate the general transcription factor complex that controls the gene transcription (Suter, 2020). Regulation of gene expression by transcription factors interplays with chromatin regulators to trigger the local alteration of chromatin structure. Transcription factors and chromatin regulators confer epigenomic and transcriptomic levels to govern the circuitry of the gene regulatory network (Wilson and Filipp, 2018).

CCCTC-binding factor (CTCF) is a transcription factor and an architectural protein for chromatin organisation (Phillips and Corces, 2009; Xiang and Corces, 2020). CTCF binds to its cognate DNA sites highly distributed across the genome in all cell types (Chen *et al.*, 2012a; Maurano *et al.*, 2015). CTCF plays diverse roles in transcription regulation and mediating higher-order chromatin organisation (Braccioli and de Wit, 2019; Herold *et al.*, 2012; Phillips and Corces, 2009; Wu *et al.*, 2020). Polycomb repressor complex (PRC2) is a chromatin regulator that controls the gene repression by catalysing the methylation of lysine 27 on Histone H3 (H3K27) (Glancy *et al.*, 2021). Both PRC2 and CTCF are essential regulators for cellular development (Arzate-Mejia *et al.*, 2018; Deevy and Bracken, 2019). They also cooperate in regulating the expression of imprinted genes (Li *et al.*, 2008). Both

are essential for the appropriate expression of developmental genes in pluripotent stem cells (Dowen *et al.*, 2014; Xu *et al.*, 2014). CTCF and PRC2 are often mutated, causing dysregulation of gene expression and driving to cancer (Debaugny and Skok, 2020; Piunti and Shilatifard, 2021).

CTCFL (CCCTC-binding factor-like) is a zinc-finger transcription factor expressed selectively in male germ cells (Loukinov *et al.*, 2002). CTCFL is the paralog of CTCF with restrictive expression in male germ cells (Klenova *et al.*, 2002; Loukinov *et al.*, 2002; Sleutels *et al.*, 2012). Hence, *Ctcfl* is a germline gene as its expression is essential for spermatogenesis. CTCFL regulates the expression of essential germline genes, which are *Gal3st1* and *Prss50* (Kosaka-Suzuki *et al.*, 2011; Suzuki *et al.*, 2010). Depletion of CTCFL results in defective spermatogenesis and sub-sterility in mice (Sleutels *et al.*, 2012).

Nevertheless, CTCFL is also a cancer-testis antigen (CTA). CTA genes have a restrictive expression in normal testis, but its derepression in the non-testicular tissues is associated with malignancies (Gibbs and Whitehurst, 2018; Kalejs and Erenpreisa, 2005). CTCFL aberrant expression beyond the male germ cells is linked with the invasive phenotype of cancer (Debaugny and Skok, 2020; Debruyne *et al.*, 2019; Janssen *et al.*, 2020; Klenova *et al.*, 2002; Soltanian and Dehghani, 2018). The reactivation of CTAs beyond the male germ cells leads to the hijacking of cell development processes by supporting the survivability of cancer cells (Gordeeva, 2018). As a transcription factor, the derepression of CTCFL in cancer may contribute to the alteration of transcription profile which drives the malignant genome reprogramming and cell transformation (Debaugny and Skok, 2020; Wang *et al.*, 2011). CTCFL is also a potential candidate for immunotherapeutic against cancer and metastasis (Loukinov, 2018).

1.2 Problem Statement

CTCFL aberrant expression in the non-testicular tissues can be detrimental as the protein is a transcription factor and a CTA, which may drive the soma-to-germline expression program for malignant reprogramming (Wang *et al.*, 2011). In addition, CTCFL can activate the expression of other CTAs and germline genes in cancer, although the derepression is speculated by several studies (Debaugny and Skok, 2020; Hillman *et al.*, 2019; Woloszynska-Read *et al.*, 2011).

A study by Pugacheva *et al.* (2015) discovered the regulatory mode of CTCFL that is common between germline and cancer cells. CTCFL can bind to the CTCF binding sites near the promoter of germline genes, forming a heterodimer with CTCF and activating the aberrant gene expression in the somatic cells. This abnormal regulatory action may predispose them to cancer. In parallel, induction of aberrant CTCFL expression by Sati *et al.* (2015) during mouse embryogenesis resulted in poor growth development. They observed that the dead pups had growth retardation and malformations in multiple somatic organs, particularly the eyes, brain, and vascular system. It is peculiar that a testis-specific transcription factor could cause grievous alterations in developing non-testicular tissues. This finding provides a clue about the dysregulations by CTCFL in somatic development by altering the Transforming Growth Beta (TGF- β) signalling pathway. Still, establishing a link between CTCFL activation of the germline programme and cancer is difficult due to the lack of study on its role in germ cells.

Most transcription factors do not independently participate in gene transcriptional regulation (Will and Helms, 2014). Besides binding to their target DNA binding sites, they also interact with other regulatory proteins such as mediators and chromatin modifiers (Nakagawa *et al.*, 2018). Albeit being discovered

for 20 years, less is known about the functional binding sites of CTCFL and its cooperation with other regulators. CTCFL binding overlaps a subset of CTCF target sites as paralogous proteins due to the similar ZF DNA binding domain (Pugacheva *et al.*, 2015). The different end termini and configuration of ZF domains modulate the selection of DNA binding and interaction with protein partners (Bergmaier *et al.*, 2018; Nishana *et al.*, 2020). Less is known about the downstream regulations, primarily when CTCF and CTCFL co-express in the tumour and male germ cells. They may interact with similar DNA sequences and other regulatory proteins but with unique downstream regulations.

1.3 Rationale of Study

Germline genes, including CTAs, involve a wide range of processes such as genomic maintenance, transcriptional regulation and meiosis in spermatogenesis, which could connect to their aberrant activation derepression outcomes with neoplastic behaviours (Gibbs and Whitehurst, 2018; Whitehurst, 2014). CTCFL participates in transcriptional regulation in tumourigenesis and spermatogenesis. The implications of CTCFL "out of context" expression may associate with the derepression of germline genes, as well as dysregulation of somatic development (Pugacheva *et al.*, 2015; Sati *et al.*, 2015). Still, there is a gap of understanding of how CTCFL functions. CTCFL regulation and its role in male germ cells require further elucidation.

The current study was conducted to investigate CTCFL regulation in both expression settings for a comprehensive understanding, based on a hypothesis that the implications of CTCFL aberrant expression in non-germ cells could be linked with its functional roles in male germ cells. In addition, this study also explored the

mechanism of transcriptional regulation of CTCFL as a transcription factor that may link CTCFL function in spermatogenesis and cancer.

1.4 Research Objectives

The main objective of this study was to elucidate the regulatory profile of murine CTCFL in non-germ cells and male germ cells. The following specific objectives were targeted in this study.

- (1) To investigate the transcriptional effects of CTCFL aberrant expression in non-germ cells by introducing the ectopic CTCFL (Chapter 4).
- (2) To discover the roles of CTCFL in male germ cells by knocking down the endogenous *Ctcfl* transcript expression (Chapter 5).
- (3) To infer the regulatory mode of CTCFL as a transcription factor based on the genome-wide protein-DNA binding profile (Chapter 6).

The first and second objectives focussed on the transcriptional regulation of CTCFL in the non-germ cells and male germ cells, respectively. For the first objective, this work utilised the recombinant construct, 3×FLAG-CTCFL, for the ectopic CTCFL expression in the non-germ cells. Ectopic expression is the expression of a gene in a cell type, developmental stage, or condition where it should not be expressed (Prelich, 2012). The ectopic CTCFL expression was conducted in pluripotent embryonic stem cells (ESCs) and JK1 testicular stromal cell lines, in which the protein is repressed. Tumour or cancer cell lines were opted out of this study as the altered genetic and epigenetic events in are confounding factors may conceal the direct effects of CTCFL aberrant expression. The ectopic CTCFL was tagged FLAG-tag epitope to allow protein probing by anti-FLAG monoclonal antibody for protein expression analysis and chromatin immunoprecipitation. The

implications of CTCFL aberrant expression in the ESCs and JK1 cells was analysed by microarray to measure the global gene expression changes.

ESCs are ideal for investigating mammalian transcriptional regulation (Rao, 2012). An exhaustive range of datasets for ESCs allows simultaneous understanding and construction of hypotheses behind the multiple layers of gene expression regulation in the cell system. Furthermore, as ESCs can differentiate into all cell lineages (the germ layers and germline cells), the cells are heavily utilised for investigating developmental process and lineage specification. Several studies have used ESCs or pluripotent stem cells for investigating the roles of CTCFL (Bergmaier *et al.*, 2018; Nishana *et al.*, 2020; Sati *et al.*, 2015; Sleutels *et al.*, 2012). Meanwhile, JK1 cells are immortalised cells originating from CD34+ enriched testicular stromal cells, particularly the peritubular myoid cells (Kim *et al.*, 2008; Seandel *et al.*, 2007). This stromal cell line can be utilised as feeder cells that support the expansion of adult spermatogonial stem cells (Kim *et al.*, 2008). Thus, JK1 cells may resemble the testicular somatic cell type, closely connected with the growth of male germ cells.

Ctcfl expression knockdown was conducted in male germline stem cells (GSCs) for the second objective. GSCs are the *in vitro* cell model representing undifferentiated spermatogonial stem cells (SSCs). GSCs were established from the SSCs of neonatal mice and continuously proliferated as undifferentiated spermatogonia (Kanatsu-Shinohara *et al.*, 2003a; Kanatsu-Shinohara *et al.*, 2003b). CTCFL has the most outstanding expression level in mitotic spermatogonia compared to other germ cell stages (Rivero-Hinojosa *et al.*, 2021). Hence, GSCs express the endogenous CTCFL and can analyse the physiological role of this transcription factor in male germ cells. The knockdown approach or gene silencing was chosen to reduce the gene expression aberrantly and elucidate the direct roles of

CTCFL in the cells. The biggest concern of using an RNAi-based system for knockdown studies is the off-target effect. This study used an improved RNAi system, the Dicer-substrate short interfering RNAs (DsiRNA), to target the *Ctcfl* at specific sites. DsiRNA has improved performance than the canonical small interfering RNA (siRNA) in the RNAi pathway by minimising the off-target effect (Snead *et al.*, 2013). RNA sequencing (RNA-seq) analysed the transcriptome profile of *Ctcfl* knockdown GSCs, and data analysis computed the differential expressed genes.

Lastly, the genome-wide CTCFL enrichment was investigated to infer the regulatory mode or the mechanism of CTCFL regulation. ChIP-seq profiled the CTCFL binding sites, and further analysis retrieved the possible direct targets and associated downstream biological processes. Genomic enrichment of its paralog, CTCF and a core subunit of PRC2, SUZ12, were also assessed. An integrative analysis was conducted to assess the association between the transcriptional regulation and the annotated protein-DNA binding profile. The analysis revealed the possible mechanism of CTCFL-dependent gene regulation in non-germ cells and germ cells. Figure 1.1 illustrates the research design in this study based on the three specific objectives.

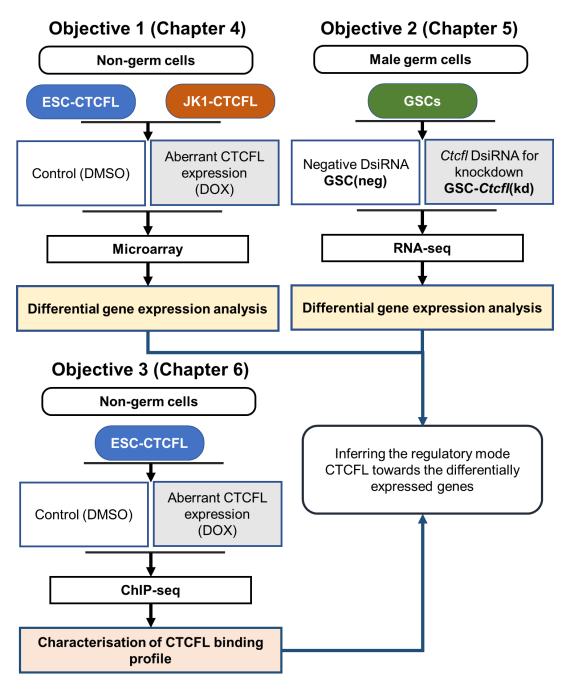


Figure 1.1 The overall research design conducted in this study was based on three specific objectives.

1.5 Thesis Organisation

The thesis is organised into chapters that cover the three specific objectives. Chapter 2 describes the biological background of transcription factors and chromatin regulators by introducing CTCF and PRC2 as the respective example. The literature review highlights the CTCFL as the protein of interest in this study. An overview of genome-wide analysis methods, including microarray, RNA-seq and ChIP-seq, was described later in this chapter.

Chapter 3 provides the essential research methods used to achieve the objectives of this study. The thesis highlighted the detailed protocol for the mammalian cell culture, quantitative polymerase chain reaction (qPCR), protein expression analysis, microarray and library preparation for RNA-seq and ChIP-seq.

Chapters 4, 5 and 6 present the experiments and findings for the three specific objectives, respectively. In Chapter 4, the aim was to investigate the implications of CTCFL aberrant expression in ESCs and JK1 cells. Transcriptomic changes were analysed using microarray to discover the differentially expressed (DE) genes. This chapter highlights the effects of CTCFL expression toward the potential dysregulation of cell development in the stem cells and somatic cells.

In Chapter 5, the aim was to elucidate the roles of CTCFL in GSCs. The expression of *Ctcfl* in GSCs was aberrantly reduced via gene knockdown followed by RNA-seq. Transcriptomic analysis exhibited alteration of gene expression, which involved cellular processes such as cell death and signalling pathways. This chapter suggests the potential role of CTCFL in the male germ cells during the early stage of spermatogenesis.

In Chapter 6, the aim was to investigate the regulatory mode of CTCFL based on its genome-wide binding profile. The enrichment of CTCFL and CTCF and a core subunit of PRC2; SUZ12 were characterised by ChIP-seq. CTCFL binding was predominant at the gene promoters that CTCF and SUZ12 occupied. Furthermore, the analysis discovered the probable alteration of the roles of CTCF and PRC2 in response to CTCFL enrichment to the overlapping binding sites.

Finally, Chapter 7 concludes the thesis by summarising the findings of this study and mapping out possible areas for future research.

CHAPTER 2

LITERATURE REVIEW

2.1 An Overview of Transcription Factors and Histone Modifications

2.1.1 Chromatin organisation

DNA stores the hereditary genetic information of an organism. DNA is a two-stranded helix in which each strand consists of four different nucleotides called adenine (A), guanine (G), cytosine (C) and thymine (T). DNA exists within the nucleus and would stretch up to two meters for human DNA. DNA is folded in a highly compacted structure known as chromatin in eukaryotic cells, including nucleic acids and histone proteins. The protein-DNA complex is called the nucleosomes in which 147 base pairs (bp) of DNA is wrapped approximately 1.65 to 1.7 times around a core protein octamer composed of two units of each of the histone proteins, H2A, H2B, H3, and H4 (Luger and Hansen, 2005; Luger *et al.*, 1997; Richmond and Davey, 2003).

Chromatin structure divides into three hierarchical states based on compaction: primary, secondary, and tertiary (Figure 2.1). The primary state consists of core chromatin that appears as beads-on-a-string, forming an eleven nm chromatin fibre. Next, a range of 20 to 50 bp of linker DNA associated with histone H1 separates the nucleosomes (Oudet *et al.*, 1975). The secondary state involves condensing the eleven nm fibre further to form a 30 nm chromatin fibre through nucleosomal interactions, but the formation remains controversial (Felsenfeld and McGhee, 1986; Maeshima *et al.*, 2010). Finally, in the tertiary level of compaction, the chromatin is folded, resulting in interchromosomal and intrachromosomal interactions of the 30 and 11 nm fibres that form the higher-order structures (Bian and Belmont, 2012; Rattner and Hamkalo, 1978).

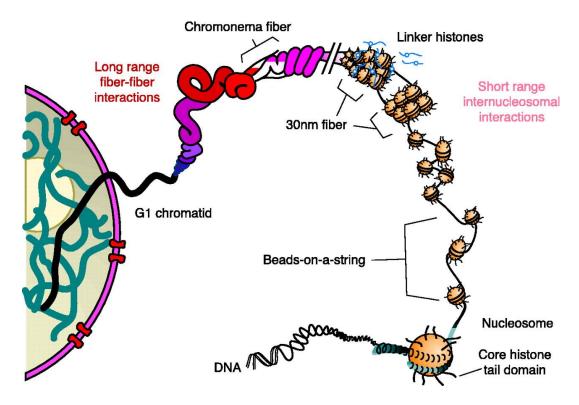


Figure 2.1 Levels of chromatin organisation. DNA and its associated proteins are called chromatin. The DNA is tightly wrapped around a histone octamer to form the nucleosome. These nucleosomes are separated by linker histones and compacted in a 30 nm chromatin fibre. Chromatin is further compacted into higher-order structures. The figure was adapted from Botchkarev *et al.* (2012)

Chromatin is functionally divided into two general states, euchromatin and heterochromatin, distributed throughout the genome. Euchromatin is less condensed, gene-rich, and often associated with active gene expression (Dillon and Festenstein, 2002). Heterochromatin is highly condensed, generally gene-poor and lacks much active gene expression. The transcriptionally inactive chromatin can be in two varieties: permanently silenced or constitutive heterochromatin and facultative heterochromatin. Constitutive chromatin localises extensively near centromeres, telomeres, and gene desert regions. Facultative chromatin is the condensed euchromatic region where the compaction is developmentally regulated (Gilbert *et al.*, 2003). The convoluted chromatin changes through histone modifications and accessibility to regulatory proteins are linked with the coordinated regulation of gene expression for cellular functions such as transcription, DNA replication and DNA repair (Blakey and Litt, 2015; Parmar and Padinhateeri, 2020).

2.1.2 Transcription factors

Gene transcriptional regulation is central to tissue-specific gene expression and stimulus-mediated gene activity in development and cellular differentiation (Holmberg and Perlmann, 2012; Hu and Gallo, 2010; Molina *et al.*, 2013; Sonawane *et al.*, 2017). Transcription factors are regulatory proteins bound to short specific DNA sequences or motifs and form a protein complex system to control the gene transcription (Lambert *et al.*, 2018; Latchman, 2010). Hence, transcription factors cooperate with DNA and coregulators, determining the transcription factor function and transcriptional activity (Cheatle Jarvela and Hinman, 2015). The protein-DNA and protein-protein interactions are the key parameters determining the three functional features of transcription factors: local chromatin access, self-sustained

remodelling, and sporadic occupancy/selective gene activation (Francois *et al.*, 2020).

Specific transcription factor binding motif localises to *cis*-regulatory elements that include the promoters, enhancers (or silencers) and insulators for precise control of gene expression (Inukai *et al.*, 2017). The promoter is a stretch of upstream DNA sequence where transcription is initiated, and this region includes the TSS to which the basal transcriptional machinery is enriched. The promoters can be tissue-specific or broadly scattered in many tissues. The tissue-specific promoters are primarily governed by transcriptional control, while the broad class is more dependent on the chromatin accessibility and epigenetic marks (Lenhard *et al.*, 2012; Mora *et al.*, 2016).

The distal *cis*-regulatory elements such as enhancers, silencers and insulators are localised far from the TSS. Enhancers and silencers are short motifs containing binding sites for transcription factors, but they have the opposite regulation (Kolovos *et al.*, 2012). The enhancers function in amplifying transcriptional regulation while the silencers suppress the gene expression. These elements possess a bifunctional role (enhancers or silencers) depending on the specific enrichment of transcription factors and epigenetic marks (Bandara *et al.*, 2021; Gisselbrecht *et al.*, 2020). Meanwhile, insulators are boundary elements that demarcate the active domains and restrict the promoters' promiscuous interaction between enhancers or silencers (Gaszner and Felsenfeld, 2006). The interactions between the promoters and the distal regulatory elements determine the transcriptional activity in housekeeping and cell or tissue-specific regulations for development and disease (Ko *et al.*, 2017; Zabidi and Stark, 2016).

The transcription factor binding specificity and functional features are modulated by multiple simultaneous protein-protein interactions (PPIs) with other transcription factors and transcriptional coregulators (Ahsendorf *et al.*, 2017; Inukai *et al.*, 2017). Most transcription factors operate as homo- or heterodimers and exhibit cooperative binding to execute multiple regulatory roles (Funnell and Crossley, 2012; Morgunova and Taipale, 2017). Transcriptional coregulators are also required to activate or suppress chromatin-dependent transcriptional signalling after transcription factors bind to cis-regulatory DNA sequences. Transcriptional coregulators (coactivators and corepressors) are diverse functional classes of chromatin-associated proteins including, but not limited to, chromatin modifiers (writers and erasers), chromatin remodellers, chromatin readers, and other scaffolding proteins (Bishop *et al.*, 2019). Transcription factors and coregulators assemble diverse regulatory complexes at the *cis*-regulatory elements. Dysregulation of these assemblies has severe implications for cell homeostasis and often leads to disease development.

2.1.3 Histone modifications

Chromatin regulators conduct the dynamic modification of chromatin architecture known as chromatin remodelling for the accessibility of transcription machinery to the regulatory elements (Morgan and Shilatifard, 2020). Remodelling occurs via covalent histone modifications by enzymes such as histone acetyltransferases (HATs), histone deacetyltransferases (HDACs) and histone methyltransferases (HMTs) (Wang *et al.*, 2007). Histone modifications include acetylation, methylation, and ubiquitination. Other modifications such as crotonylation, 5-hydroxylation, sumoylation, ADP-ribosylation, and glycosylation also may occur at the histone polypeptides and their residues (Rothbart and Strahl,

2014). Each modification plays a crucial role in mediating the functional chromatin states for the proper organisation and efficient execution of nuclear processes (Adkins *et al.*, 2013).

The best-studied modifications are methylated and acetylated lysine residues on histone H3. They are deposited on particular areas of the genome and associated with distinct states of gene transcription. Histone acetylation is generally correlated with active transcription, histone methylation associates with a variety of transcriptional states, depending on the particular lysine position that is methylated. (Gates et al., 2017). For example, histone-3-lysine-4-mono-, di- and trimethylation (H3K4me1/2/3), histone-3-lysine-36 mono- and tri-methylation (H3K36me1/3), histone-3-lysine-27-acetylation (H3K27ac) and histone-3-lysine-9-acetylation (H3K9ac) are located in actively transcribed regions and accumulated near the euchromatin. On the other hand, the heterochromatic regions are often marked by histone-3-lysine 9 (H3K9), histone-3-lysine-27 (H3K27) and histone-4-lysine-20 methylation (H4K20me) and lack of histone acetylation. Table 2.1 lists the genomic localisation of notable histone H3 modifications associated with transcription states of chromatin (Gates et al., 2017; Zhao and Garcia, 2015). Another transcription state is the poised state of chromatin by the bivalency of histone marks. Bivalent chromatin is the genomic feature co-occupied by active H3K4me3 and repressive H3K27me3, mainly at the promoters (Bernstein et al., 2006; Blanco et al., 2020). Bivalency regulates gene expression for cell fate and tissue lineage commitment in development (Harikumar and Meshorer, 2015; Jeon and Tucker-Kellogg, 2020).

Table 2.1 Transcriptional role and localisation of histone modifications (Gates $et\ al.$, 2017; Zhao and Garcia, 2015).

| Histone modification | Transcriptional role | Location of enrichment |
|----------------------|----------------------|---|
| H3K4me1 | Activation | Enhancers |
| H3K4me3 | Activation | Promoters |
| H3K27ac | Activation | Enhancers, promoters |
| Н3К9Ас | Activation | Enhancers, promoters |
| H3K36me3 | Activation | Gene bodies |
| H3K79me3 | Activation | Gene bodies |
| H3K27me3 | Repression | Promoters, gene-rich regions |
| H3K9me3 | Repression | Satellite repeats, telomeres, pericentromeres |

Histone-modifying enzymes are chromatin regulators that mediate histone modifications that frequently operate as large multiprotein complexes (DesJarlais and Tummino, 2016). The histone-modifying enzymes are grouped into writers, erasers, and readers (Biswas and Rao, 2018; Hojfeldt *et al.*, 2013). The 'writers' are enzymes that set the post-translational modification to the residues on histones. The covalently modified histone tails are recognised via specific domains of a group of proteins referred to as 'readers'. The enzymes that remove the post-translational modifications are called 'erasers'. These enzymes and protein domains modify and read the specific amino acids on the histones and form the basis of the dynamic epigenetic regulation of gene expression (Figure 2.2).

Among the writers are the Polycomb group (PcG) and Trithorax group (TrxG) proteins which play a fundamental role in regulating histone modifications connected with gene repression and activation, respectively (Blanco *et al.*, 2020; Kuroda *et al.*, 2020). The TrxG complexes display cooperativity with PcG complexes through their opposing roles in regulating gene expression and development (Kadoch *et al.*, 2016; Kuroda *et al.*, 2020). PRC2 is a PcG complex that maintains the repression of gene expression by depositing H3K27me3 at the gene regulatory elements (Gaydos *et al.*, 2014). Meanwhile, the TrxG complex deposits active H3K4me3 to oppose the PcG silencing (Tie *et al.*, 2014). Both complexes are also responsible for establishing the bivalent domains for poised chromatin and may act as the master switch that coordinates the transcriptional programming (Kuroda *et al.*, 2020).

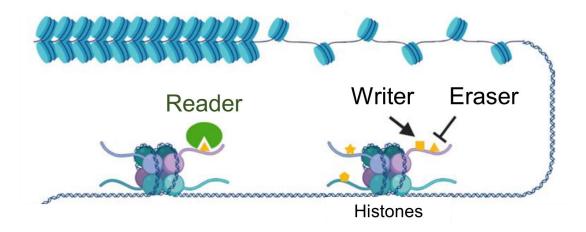


Figure 2.2 The regulators of histone modifications. Histone post-translational modifications are added or removed by specific enzymes ('writers' and 'erasers', respectively) and recognised by their binding proteins ('readers'). The image was adapted from Cao and Yan (2020).

2.2 Transcription Factor: CCCTC-binding Factor (CTCF)

2.2.1 Background

CTCF is a highly conserved zinc-finger DNA-binding protein initially discovered as a transcriptional repressor of c-myc (Filippova *et al.*, 1996; Klenova *et al.*, 1993). CTCF plays multiple roles in transcriptional activation and repression, regulation of imprinted loci and X chromosome inactivation (XCI) and organisation of higher-order chromatin by chromosomal looping and nuclear tethering (Braccioli and de Wit, 2019; Herold *et al.*, 2012; Phillips and Corces, 2009; Wu *et al.*, 2020). CTCF expresses in all stages of development across most metazoan tissues (Phillips and Corces, 2009). Notably, CTCF regulates the lineage specification during cell/tissue-type specific gene expression by controlling the transcriptional regulation and chromatin organisation (Arzate-Mejia *et al.*, 2018; Braccioli and de Wit, 2019).

CTCF comprises an N-terminal domain, a central DNA binding domain with 11 C2H2 ZFs and a C-terminal domain (Klenova *et al.*, 1993). The plethora of CTCF roles is contributed by the three domains that determine the binding site preference and regulation of chromosome organisation (Nishana *et al.*, 2020; Ong and Corces, 2014; Pugacheva *et al.*, 2020). The dynamic combination of ZFs binds to different DNA sequences and interacts with various protein factors (Hashimoto *et al.*, 2017; Nakahashi *et al.*, 2013). In addition, CTCF may interact with other proteins, RNA and susceptible to post-translational modifications that could affect interactions with DNA or other proteins, as reviewed elsewhere (Arzate-Mejia *et al.*, 2018; Braccioli and de Wit, 2019; Wu *et al.*, 2020; Zlatanova and Caiafa, 2009).

2.2.2 CTCF target sites (CTSes)

CTCF is a sequence-specific transcription factor that regulates 3D genome organisation and critical aspects of gene expression such as transcription activation

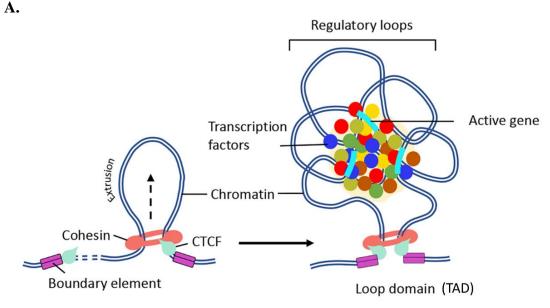
and repression, RNA splicing, and enhancer/promoter insulation. (Ong and Corces, 2014). CTCF target sites (CTSes) are highly conserved and widely spread across the mammalian genome. A vast number of more than 80,000 CTSes may be found in the mammalian genome (Chen *et al.*, 2012a; Maurano *et al.*, 2015). CTCF binding is predominant to the distal regions and gene bodies than the TSS (Handoko *et al.*, 2011; Holwerda and de Laat, 2013). Still, CTCF enrichment near TSS is crucial for forming transcription or enhancer hubs through long-distance chromatin interactions that bridge distal enhancers with target promoters (Guo *et al.*, 2012; Kubo *et al.*, 2021).

CTCF forms a complex with cohesin, a key component of chromatin and binds to the DNA to establish the DNA loop domains (Hansen *et al.*, 2017; Lee and Iyer, 2012) (Figure 2.3A). The loops facilitate the interactions between promoter and enhancer elements. A loop domain, known as topologically associated domains (TAD), is formed by multiple loops of intervening DNA between two convergent CTCF sites. TADs form the basis of chromatin regulatory segmentation that demarcate the gene regulatory elements (Merkenschlager and Nora, 2016; Xiang and Corces, 2020).

CTCF exhibits dynamic binding to the cognate genomic sites, recognised by diverse combinations of its 11 ZFs as a "multivalent protein" (Xu et al., 2018a). The systematic mapping of CTCF ZFs to specific chromatin sites enables CTCF to execute diverse functions in different contexts and cell types. CTCF core motif consists of about 20 bp, in which the CTCF ZF domain employs ZFs three to seven to bind to the 15 bp consensus sequence. Meanwhile, the ZFs nine to eleven modulate CTCF-binding stability by interacting with the second CTCF motif located at 21 to 22 bp upstream to the core motif (Guo et al., 2015; Hashimoto et al., 2017;

Nakahashi *et al.*, 2013; Persikov and Singh, 2014; Xie *et al.*, 2007; Xu *et al.*, 2018a). The core sequences can be partial palindromic, and the directionality is determined by the second motif (Guo *et al.*, 2015; Wu *et al.*, 2020; Xie *et al.*, 2007). The CTCF binding motifs are displayed in Figure 2.3B, adapted from an extensive review by Wu *et al.* (2020).

CTSes elements can exist as single non-clustered and clustered CTCF sites. CTCF clustered sites represent a group of at least two CTSes separated by less than ten kb on the genome. The vast majority of clustered CTSes colocalise with cohesin and enrich near transcription start sites to stabilise the chromatin contacts (Kentepozidou *et al.*, 2020). Meanwhile, the single CTSes function as the insulators that restrict the interactions between promoters and distal enhancers (Jia *et al.*, 2020).



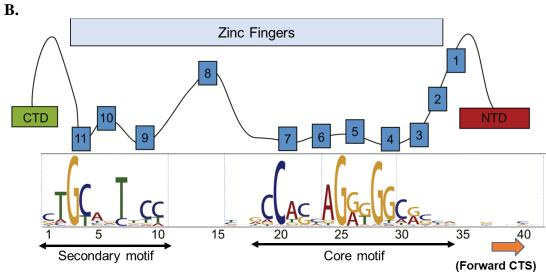


Figure 2.3 The formation of chromatin loops and CTCF binding motifs. **A.** CTCF and cohesin generate chromatin loops through loop extrusion and form TAD. The image was adapted from Pongubala and Murre (2021). **B.** The combination of CTCF ZFs binds to the core and secondary binding motifs. The image was adapted from Wu *et al.* (2020).

_

2.2.3 CTCF in the development and cancer

The CTCF-mediated interactions are also closely related to developmental gene regulation (Allahyar *et al.*, 2018; Merkenschlager and Nora, 2016). CTCF is required for embryonic development since CTCF-null embryos cannot implant (Moore *et al.*, 2012). CTCF mediates chromatin looping in the pluripotent genome and is essential to determine the stem cells' fate (Rao *et al.*, 2014). As an architectural protein, CTCF controls the expression of pluripotency genes and lineage-specifying genes by organising the chromatin state (Balakrishnan *et al.*, 2012; Dowen *et al.*, 2014; Handoko *et al.*, 2011; Nora *et al.*, 2017; Pekowska *et al.*, 2018). The pluripotency genes such as *Nanog* and *Oct4* are localised within the CTCF-anchored super-enhancer domain. The domains insulate the super-enhancers from interacting with lineage-specific master transcription factors and maintain the stem cell state (Dowen *et al.*, 2014; Justice *et al.*, 2020; Whyte *et al.*, 2013).

CTCF contributes to the lineage specification and development in multiple tissues, including brain, cardiovascular, limb, muscle, retina, immune cells, and gametes (Arzate-Mejia *et al.*, 2018). The chromatin has a flexible structure with weak organisation and high accessibility in the pluripotent stem cells (Schlesinger and Meshorer, 2019). Gain and enhancement of chromatin loops, especially at CTCF binding sites, mark the exit from pluripotency to differentiation (Pekowska *et al.*, 2018). The chromatin becomes more compartmentalised as the developmental program activates and progresses to more lineage-specific regulation (Zheng and Xie, 2019).

CTCF binding at the promoters mediates the promoter-promoter and promoter-enhancer interactions for active gene expression during cellular differentiation (Kubo *et al.*, 2021). An architectural protein, YY1 interacts with