

UNIVERSIDADE FEDERAL DE MINAS GERAIS

Luciana Werneck Zuccherato

**Estrutura populacional e
diversidade de variações em número
de cópias (CNVs) de genes do
sistema imune em populações
nativas da América do Sul**

BELO HORIZONTE

2012

Luciana Werneck Zuccherato

**Estrutura populacional e
diversidade de variações em número
de cópias (CNVs) de genes do
sistema imune em populações
nativas da América do Sul**

Tese apresentada ao Programa de Pós-graduação em Genética do Departamento de Biologia Geral do Instituto de Ciências Biológicas da Universidade Federal de Minas Gerais, como requisito parcial à obtenção do título de Doutor em Genética
Orientador: Prof. Dr. Eduardo Martín Tarazona-Santos

**Belo Horizonte
Departamento de Biologia Geral
Instituto de Ciências Biológicas
Universidade Federal de Minas Gerais
Novembro de 2012**

Agradecimentos

Gostaria prestar meus sinceros agradecimentos:

Ao professor Eduardo Tarazona pela orientação e atenção durante todas as conquistas científicas desde o meu ingresso no LDGH.

Ao professor Edward Hollox por ter me aceitado carinhosamente em seu laboratório e por todos os votos de confiança. Aos amigos ingleses, escoceses, australianos, italianos, alemães, gregos, árabes e brasileiros que me ensinaram que viver fora do país é uma experiência enriquecedora.

Aos colegas e amigos presentes e ausentes do LDGH: Maria Clara, Wagner, Maíra, Ferdi, Camila, Andrea, Nelson, Rodrigo, Moara, Juliana, Roxana, Fernanda, Marília, Giordano, Raquel, Laélia, Silvana Matheus e Latife pelo carinho e companheirismo.

À Cláudia Benedetto, por ter me transformado de técnica em cientista, pelo exemplo de dedicação e carinho pela Ciência e pelo mundo.

Aos professores da Pós-graduação em Genética pelos importantes ensinamentos que contribuíram para a minha formação e pela amizade e carinho recebidos. A todos os colegas do Departamento de Biologia Geral, que sempre me trataram com carinho e simpatia. Em especial aos do LBEM e Genepop, que sempre me fizeram sentir parte do laboratório. A todos os meus colegas de coração do ICB por todas as alegrias que me trouxeram no dia-a-dia e por nutrirem por mim o mesmo carinho que eu sinto por eles.

À Dani e a Lane, pelos conselhos essenciais nas minhas escolhas.

Aos meus pais, irmãos e sobrinhos por sempre me mostrarem o quanto a vida é linda! Por sempre me apoiarem nas minhas escolhas e sempre me lembrarem do quão especial eu sou. Aos meus amigos Mirinha, Elisa, Danusa e Sandra por sempre terem ficado ao meu lado nos momentos de alegrias e tristezas.

Enfim, a todos que plantaram uma semente de carinho em mim e me transformaram sempre numa pessoa melhor! Obrigada!

Índice

| | |
|---|------|
| Índice de Figuras | VI |
| Índice de Tabelas | VIII |
| Lista de Abreviaturas e Símbolos | IX |
| Resumo | XI |
| <i>Abstract</i> | XII |
| | |
| <i>I. INTRODUÇÃO</i> | 1 |
| 1. Variações estruturais no genoma e CNVs | 2 |
| 2. Extensão dos CNVs no genoma humano | 4 |
| 3. CNVs e SNPs | 5 |
| 4. Estrutura genética das populações humanas e CNVs | 6 |
| 5. Mecanismos de formação dos CNVs | 10 |
| 5.1. Recombinação Homóloga Não Alélica (NAHR) | 10 |
| 5.2. Junção de Extremidades Não Homólogas (NHEJ) | 13 |
| 5.3. <i>Fork Stalling and Template Switching</i> (FoSTeS) | 13 |
| 5.4. Elementos de Retrotransposição L1 | 14 |
| 6. Métodos de detecção e tipagem de CNVs | 15 |
| 6.1. PCR quantitativa em tempo real | 15 |
| 6.2. <i>Paralogue Ratio Test</i> (PRT) | 16 |
| 6.3. <i>Multiplex Amplification Probe Hybridization</i> (MAPH) | 16 |
| 6.4. <i>Multiplex Ligation dependent Probe Amplification</i> (MLPA) | 16 |
| 6.5. Hibridização Genômica Comparativa (CGH) | 17 |
| 6.6. Arranjos de SNPs | 18 |
| 6.7. Sequenciamento de nova geração (NGS) | 18 |
| 7. Aspectos funcionais dos CNVs | 19 |
| 8. CNVs e genes do sistema imune | 21 |
| 8.1. Beta-defensinas | 22 |

| | |
|---|-----------|
| | IV |
| 8.2. Receptores Fc gama | 25 |
| 8.3. Quimiocinas <i>CCL3L1/CCL4L1</i> | 29 |
| II. OBJETIVOS | 33 |
| 1. Objetivo geral | 34 |
| 2. Objetivos específicos | 34 |
| III. METODOLOGIA | 35 |
| 1. Amostras e populações estudadas | 36 |
| 2. Estimativas de miscigenação nos indivíduos nativos americanos | 37 |
| 3. Tipagem do número de cópias | 38 |
| 3.1. PCR quantitativa | 38 |
| 3.2. <i>Paralogue ratio Test</i> (PRT) | 39 |
| 3.2.1. Metodologia do PRT | 43 |
| A. Seleção da sequência dos iniciadores e PCR | 43 |
| B. Eletroforese capilar e estimativa da área do pico (<i>peak area call</i>) | 44 |
| C. Estimativa do número inteiro de cópias | 44 |
| i. Normalização dos dados | 44 |
| ii. Cálculos de máxima verossimilhança para determinação do número inteiro de cópias | 45 |
| iii. Clusterização | 47 |
| 4. Armazenamento dos dados | 48 |
| 5. Estatísticas populacionais dos diplótipos de CNVs e variantes alélicas dos locos de receptores Fc gama, beta-defensinas e quimiocinas <i>CCL3L1/CCL4L1</i> | 49 |
| 5.1. Distribuição de frequências dos diplótipos de CNVs e variantes alélicas e comparação em escala global | 49 |
| 5.2. Estimativas de frequências alélicas | 50 |
| 5.3. Análises de diversidade intrapopulacional e interpopulacional | 51 |
| 5.3.1. Escalonamento multidimensional | 52 |

| | |
|---|----|
| IV. RESULTADOS | 53 |
| 1. Estimativas de miscigenação | 54 |
| 2. Diversidade de variações em número de cópias de genes do sistema imune | 55 |
| 2.1. Região de beta-defensinas | 55 |
| 2.1.1. Correlação entre a técnica de PCR quantitativa e PRT | 55 |
| 2.1.2. Distribuição do número de cópias de beta-defensinas utilizando a técnica de PRT | 56 |
| 2.2. Frequência do Clado II da região promotora de <i>DEFB103</i> | 58 |
| 2.3. Região das quimiocinas <i>CCL3L1/CCL4L1</i> | 59 |
| 2.3.1. Distribuição do número de cópias | 59 |
| 2.3.2. Correlação do número de cópias entre os marcadores <i>CCL3C</i> , <i>CCL4A</i> e <i>LTR61A</i> | 60 |
| 2.4. Receptores Fc Gama | 61 |
| 2.4.1. Distribuição do número de cópias dos genes <i>FCGR3A</i> , <i>FCGR3B</i> e <i>FCGR2C</i> | 61 |
| 2.5. Frequência do alótipo HNA1a do gene <i>FCGR3B</i> e polimorfismo Q57X do gene <i>FCGR2C</i> | 64 |
| 3. Análise de diversidade e estruturação populacional | 65 |
| 3.1. Diversidade interpopulacional e Escalonamento multidimensional | 65 |
| 3.2. Diversidade intrapopulacional | 70 |
| V. DISCUSSÃO | 73 |
| Problemas metodológicos na inferência do número inteiro de cópias | 74 |
| Variação do número de cópias e variantes alélicas | 75 |
| Diversidade interpopulacional e intrapopulacional | 78 |
| VI. CONCLUSÕES E PERSPECTIVAS | 79 |
| VII. REFERÊNCIAS BIBLIOGRÁFICAS | 82 |
| VIII. ANEXOS | 95 |

Índice de Figuras

| | |
|---|----|
| Figura 1 Tipos de variações estruturais..... | 3 |
| Figura 2 Localização de 1447 CNVRs no genoma humano.. | 5 |
| Figura 3 Genes em CNVRs que apresentam altos níveis de estratificação populacional | 9 |
| Figura 4 Modelos dos quatro maiores mecanismos envolvidos na formação de CNVs..... | 11 |
| Figura 5 Modelos esquemáticos do mecanismo de NAHR entre LCRs..... | 12 |
| Figura 6 Estrutura do <i>cluster</i> de beta-defensinas na região 8p23.1..... | 23 |
| Figura 7 Mapa da localização dos genes FCGR no cromossomo 1q23. | 26 |
| Figura 8 Localização das quimiocinas CC no cromossomo 17q12..... | 30 |
| Figura 9 Visão geral da técnica de PRT..... | 40 |
| Figura 10 Ensaio de REDVR para determinar a razão entre as cópias dos genes <i>FCGR3A</i> e <i>FCGR3B</i> | 41 |
| Figura 11 Fluxograma da metodologia das técnicas de PRT e REDVR para a determinação do número inteiro de cópias do presente estudo..... | 42 |
| Figura 12 Regiões <i>Teste</i> e <i>Referência</i> selecionadas para o ensaio de PRT107A..... | 43 |
| Figura 13 Regressão linear dos valores obtidos das amostras controle do ensaio de PRT para determinação do número de cópias da região de beta-defensinas.. | 45 |
| Figura 14 Estimativa de ML para determinação do número de inteiro de cópias da região <i>CCL3L1/CCL4L1</i> | 47 |
| Figura 15 Histograma de clusterização dos resultados normalizados do loco <i>CCL3L1/CCL4L1</i> utilizando a ferramenta <i>CNVtools</i> | 48 |
| Figura 16 Estimativas de miscigenação das populações nativas da América Sul do presente estudo..... | 55 |
| Figura 17 Regressão linear dos valores obtidos com o <i>TaqMan® Copy Number Assay</i> para o gene <i>DEFB4</i> (<i>Applied Biosystems</i>) e a técnica de PRT..... | 56 |
| Figura 18 Número de cópias da região de beta-defensinas.. | 57 |

| | |
|---|----|
| Figura 19 Distribuição da frequência relativa do Clado II do gene <i>DEFB103</i> nas populações Ashaninka, Monte Carmelo, Shimaá e Quéchua em conjunto com as populações do painel HGDP-CEPH. | 58 |
| Figura 20 Distribuição do número de cópias da região <i>CCL3L1/CCL4L1</i> | 59 |
| Figura 21 Estrutura da região de quimiocinas <i>CCL3L1/CCL4L1</i> e localização dos marcadores usados no ensaio de PRT. | 60 |
| Figura 22 Matriz de dispersão dos valores normalizados dos marcadores CCL3C, CCL4A e LTR61A da região <i>CCL3L1/CCL4L1</i> | 61 |
| Figura 23 Distribuição da variação de número de cópias dos genes <i>FCGR3A</i> , <i>FCGR3B</i> e <i>FCGR2C</i> para as populações nativas peruanas do presente estudo em conjunto com as populações mundiais do painel HGDP-CEPH. | 63 |
| Figura 24 Frequência do alótipo HNA1a do gene <i>FCGR3B</i> e alelo Q do polimorfismo Q57X do gene <i>FCGR2C</i> para as populações nativas peruanas do presente estudo em conjunto com as populações mundiais do painel HGDP-CEPH. | 64 |
| Figura 25 Representação da distância genética por MDS das populações Ashaninka, Monte Carmelo, Shimaá e Quéchua do presente estudo e população europeia para o loco <i>CCL3L1/CCL4L1</i> | 66 |
| Figura 26 Representação da distância genética por MDS das populações do painel HGDP-CEPH e das populações Ashaninka, Monte Carmelo, Shimaá e Quéchua do presente estudo para a variação de número de cópias da região de beta-defensinas. | 67 |
| Figura 27 Representação da distância genética por MDS das populações do painel HGDP-CEPH e das populações Ashaninka, Monte Carmelo, Shimaá e Quéchua do presente estudo para a variação de número de cópias dos genes <i>FCGR3A</i> , <i>FCGR3B</i> e <i>FCGR2C</i> | 68 |
| Figura 28 Representação da distância genética por MDS das populações do painel HGDP-CEPH e das populações Ashaninka, Monte Carmelo, Shimaá e Quéchua do presente estudo para as variantes alélicas <i>FCGR3B</i> HNA1a/1b, <i>FCGR2C</i> Q57X e <i>DEFB103</i> Clado I/II. | 69 |
| Figura 29 Diversidade intrapopulacional estimada pelas diferenças médias par a par de números de cópias do loco <i>CCL3L1/CCL4L1</i> nas populações nativas peruanas Ashaninka, Monte Carmelo, Shimaá e Quéchua e população europeia. | 70 |
| Figura 30 Diversidade intrapopulacional estimada pelas heterozigosidade esperada da região de beta-defensinas e dos genes <i>FCGR3A</i> , <i>FCGR3B</i> e <i>FCGR2C</i> , e variantes alélicas Clado I/II de <i>DEFB103</i> , HNA1a/1b e Q57X nas populações nativas peruanas Ashaninka, Monte Carmelo, Shimaá e Quéchua e populações mundiais do painel HGDP-CEPH. | 71 |

Índice de Tabelas

| | |
|--|----|
| Tabela 1 Descrição das populações nativas Peruanas utilizadas no presente estudo..... | 36 |
| Tabela 2 Razão dos genótipos das variantes funcionais HNA1a/HNA1b e Q57X para estimativa de frequência populacional. | 50 |
| Tabela 3 Codificação dos diplótipos e genótipos para análise de parâmetros de diversidade no programa <i>Arlequin</i> | 52 |
| Tabela 4 Análise da variância molecular (AMOVA) dos CNVs das regiões e genes do sistema imune e suas variações funcionais. | 66 |
| Tabela 5 Número médio de cópias, frequência das variantes alélicas e Índices de diversidade intrapopulacional dos indivíduos nativos analisados no presente estudo..... | 72 |

Lista de Abreviaturas e Símbolos

Encontram-se abaixo as definições das abreviaturas e símbolos que foram citados no texto:

aCGH- *array Comparative Genomic Hybridization* (Arranjo de hibridização genômica comparativa)

AMOVA- *Analysis of Molecular Variance* (Análise de variância molecular)

AR- Artrite reumatoide

CEPH- *Centre d'etude de Polymorphismes Humaines*

CGH- *Comparative Genomic Hybridization* (Hibridização genômica comparativa)

CNV- *Copy Number Variation* (Variação em número de cópia)

DA- Doença de Alzheimer

DL- Desequilíbrio de ligação

DNA- *Deoxyribonucleic acid* (Ácido desoxirribonucleico)

DP- Doença de Parkinson

EM- *Expectation-Maximization* (Maximização de esperança)

FCGRs- *Fc Gamma Receptors* (Receptores Fc gama)

FISH- *Fluorescent in situ Hybridization* (Hibridização fluorescente *in situ*)

FoSTeS- *Fork Stalling and Template Switching* (Interrupção da forquilha e troca de molde)

HAART- *Highly Active Antiretroviral Therapy* (Terapia antirretroviral fortemente ativa)

HGDP- *Human Genome Diversity Panel* (Painel de Diversidade Genômica Humana)

HNA1- *Human Neutrophil Antigen-1* (Antígeno neutrófilo humano 1)

Ia- Informatividade de ancestralidade

IBD- *Identity by Descend* (Identidade por descendência)

IgG- Imunoglobulina gama

INDELS- Polimorfismos de inserção-deleção

LCLs- Linhagens celulares de linfoblastos

LCR- *Low Copy Repeat* (Repetição de baixo número de cópias)

LES- Lúpus Eritematoso Sistêmico

- LINE-** *Long Interspersed Nuclear Element* (Elemento nuclear longo intercalado)
- LTR-** *Long Terminal Repeat* (Repetição longa terminal)
- MAPH-** *Multiplex Amplification Probe Hybridization*
- MDS-** *Multidimensional Scaling* (Escalonamento multidimensional)
- MIAs-** Marcadores Informativos de Ancestralidade
- ML-** *Maximum Likelihood* (Máxima verossimilhança)
- MLPA-** *Multiplex Ligation Dependent Probe Amplification*
- NAHR-** *Non Allelic Homologous Recombination* (Recombinação homóloga não alélica)
- NK-** *Natural Killer*
- NHEJ-** *Nonhomologous End Joining* (Junção de extremidades não homólogas)
- OMIM-** *Online Mendelian Inheritance in Man*
- pb-** Pares de bases
- PCR-** *Polymerase Chain Reaction* (Reação em cadeia da polimerase)
- PRT-** *Paralogue Ratio Test*
- REDVR-** *Restriction Enzyme Digest Variant Ratio*
- RNA-** *Ribonucleic acid* (Ácido ribonucleico)
- SINE-** *Short Interspersed Nuclear Element* (Elemento nuclear intercalado curto)
- SNP-** *Single Nucleotide Polymorphism* (Polimorfismo de nucleotídeo único)
- SVA-** elemento hominídeo derivado das repetições SINE-R, VNTR e *Alu*
- DT1-** Diabetes Tipo 1
- TPRT-** *Target Site-Primed Reverse Transcription*
- VNTR-** *Variable Number Tandem Repeat* (Repetição em *tandem* de número variável)

Resumo

Regiões variáveis em número de cópias (CNVRs) são formalmente definidas como segmentos de DNA que variam de um kilobase a vários megabases de tamanho e apresentam variações em número de cópias em relação a uma sequência referência, cujo valor usual de cópias é 2. Essas regiões representam uma parte significativa da variação genética humana, onde os graus de diversidade global e os padrões de estrutura genética se assemelham ao observado em análises de polimorfismos de nucleotídeo único (*Single Nucleotide Polymorphisms*, SNPs). Podem estar presentes em regiões genômicas que envolvem genes sensíveis a dosagem, cujos impactos funcionais são responsáveis pelas associações entre variantes de número de cópias (CNVs) raras e comuns e uma variedade de doenças genéticas complexas.

No presente trabalho nos propusemos a estimar a diversidade de genes que apresentam variação de número de cópias no genoma de populações nativas, cuja amostragem é sempre sub-representada nos grandes estudos genéticos e limitada às populações disponíveis no painel comercial HGDP-CEPH. Analisamos a distribuição de frequência de CNVs multicópias do *cluster* de beta-defensinas, genes *CCL3L1/CCL4L1*, *FCGR3A*, *FCGR3B* e *FCGR2C* e as variantes alélicas *FCGR3B-HNA1a/1b*, *FCGR2C-Q57X* e Clado I e II do gene *DEFB103*, em um grande conjunto de amostras autóctones de nativos americanos peruanos das populações de Ashaninka, Monte Carmelo, Shimaá e Quéchua, utilizando o técnica de *Paralogue Ratio Test*. Os resultados foram comparados com dados disponíveis na literatura para se investigar a existência de distribuições diferenciais de diplótipos em nativos sul-americanos, e acrescentar novas informações sobre a diversidade genética de CNVs do continente. Os resultados mostraram que os ameríndios tendem a apresentar uma maior frequência de eventos de deleções no gene *FCGR3B*, um aumento da frequência do alótipo *FCGR3B-HNA1a*, assim como de duplicações do número de cópias da região *CCL3L1/CCL4L1*. Especialmente, a população Shimaá mostra uma maior frequência do diplótipo 7 cópias na região de beta-defensinas.

Estes resultados representam possíveis impactos funcionais dos genes envolvidos na resposta imunológica destas populações apresentadas no contexto evolutivo da diversidade genômica global, com futuras implicações para estudos epidemiológicos de suscetibilidade a doenças autoimunes e infecciosas.

Abstract

Copy number variable regions (CNVRs) are formally defined as a segment of DNA ranging from one kilobase to several megabases in size and with variable copy number in comparison to the usual copy number of two of a reference genome. These regions represent a significant part of human genetic variation and reveal degrees of diversity and overall patterns of genetic structure comparable to single nucleotide polymorphism (SNP) diversity. Their presence in genomic regions harboring dosage-sensitive genes may have functional impacts that account for associations between rare and common copy number variation (CNV) and a variety of complex genetic diseases.

In the present study we aimed to estimate the diversity of genes that show CNV in Native American populations whose sampling is always underrepresented in genetic studies and limited to populations from the HGDP-CEPH commercial panel. We analyzed the frequency distribution of the multicopy immune system beta-defensin region, *CCL3L1/CCL4L1*, *FCGR3A*, *FCGR3B* and *FCGR2C* genes, and the allelic variants *FCGR3B*-HNA1a/1b, *FCGR2C*-Q57X and *DEFB103*-Cladel/II in a large sample set of autochthonous Peruvian Native Americans Ashaninka, Monte Carmelo and Shimaa populations using the Parologue Ratio Test (PRT) technique. The results were compared to published data to investigate the existence of differential diplotype distributions in Native Americans from South Americans and brought up to date information concerning their CNV diversity. Amerindians tend to display an increased frequency of deletion events in the *FCGR3B* gene and of the *FCGR3B* HNA1a allotype, and duplication of *CCL3L1/CCL4L1* copy number. Especially, the Shimaa population shows a higher frequency of the seven copies diplotype at the beta-defensin region.

These results posit possible impacts to expression profiles of the involved immune genes in the presented population in the context of evolutionary global genomic diversity, with implications for future epidemiological studies on susceptibility to autoimmune and infectious diseases.