

**Universidade de São Paulo
Escola Superior de Agricultura Luiz de Queiroz**

**Análise da estrutura populacional e do desequilíbrio de ligação de um
painel de acessos de sorgo: uma abordagem usando teoria da
coalescência**

João Ricardo Bachega Feijó Rosa

Tese apresentada para obtenção do título de Doutor
em Ciências. Área de concentração: Genética e Me-
lhoramento de Plantas

**Piracicaba
2016**

João Ricardo Bachega Feijó Rosa
Engenheiro Agrônomo

**Análise da estrutura populacional e do desequilíbrio de ligação de um
painel de acessos de sorgo: uma abordagem usando teoria da
coalescência**

versão revisada de acordo com a resolução CoPGr 6018 de 2011

Orientador:

Prof. Dr. **ANTONIO AUGUSTO FRANCO GARCIA**

Tese apresentada para obtenção do título de Doutor
em Ciências. Área de concentração: Genética e Me-
lhoramento de Plantas

Piracicaba
2016

**Dados Internacionais de Catalogação na Publicação
DIVISÃO DE BIBLIOTECA - DIBD/ESALQ/USP**

Rosa, João Ricardo Bachega Feijó

Análise da estrutura populacional e do desequilíbrio de ligação de um painel de acessos de sorgo: uma abordagem usando teoria da coalescência / João Ricardo Bachega Feijó Rosa. -- versão revisada de acordo com a resolução CoPGr 6018 de 2011. -- Piracicaba, 2016 .

83 p.

Tese (Doutorado) -- USP / Escola Superior de Agricultura "Luiz de Queiroz".

1. Passado Evolutivo 2. Fluxo Gênico 3. Recombinação 4. Associações Não-aleatórias 5. Mapeamento Associativo 6. *Sorghum bicolor* . I. Título.

DEDICATÓRIA

À minha avó materna, **Amélia** (*em memória*), por todo seu amor, carinho e dedicação, além da inesquecível convivência. À minha mãe, **Fátima**, pelo amor incondicional e exemplo de perseverança e luta. À minha irmã, **Daniela**, pela força e determinação nos momentos difíceis.

Às três, por fazerem o possível e o impossível para que eu chegasse até aqui, e por sempre acreditarem em meus sonhos:

Dedico.

À minha noiva, **Natália**, o grande amor da minha vida. Meu caminho se transformou completamente com a sua chegada. Pelo amor, carinho, companherismo, amizade, paciência, força e motivação em todos os momentos das nossas vidas:

Dedico.

Ao meu cunhado, **Gilberto**, o grande irmão que não tive em minha infância. Pelo companherismo, amizade, força e humildade:

Dedico.

AGRADECIMENTOS

À Deus, primeiramente. Enorme foi a graça concedida durante o esforço da caminhada. Força, luta, coragem, dedicação e persistência. Paciência, tranquilidade e serenidade. Faltam-me palavras para agradecer-lhe por toda honra e glória conquistadas.

Às pessoas mais importantes da minha vida: Maria Fátima (mãe), Natália (noiva), Daniela (irmã) e Gilberto (cunhado). Muito obrigado por vocês fazerem parte da minha vida e por terem estado ao meu lado, concedendo muita força, durante esta caminhada. Só vocês sabem o quanto foi difícil continuar neste caminho por uma série de razões que não cabe ser exposta neste momento. Muitíssimo obrigado! Serei eternamente grato por tudo!

Meus sinceros agradecimentos ao Departamento de Genética e ao Programa de Pós-Graduação em Genética e Melhoramento de Plantas da ESALQ/USP. Foi uma honra ter tido a oportunidade de cursar o doutorado (e o mestrado) nesta casa, que é digna de grande excelência e brilhantismo. Grande aprendizado foi-me concedido ao longo dos últimos seis anos, que certamente se eternizarão em minha memória.

Um agradecimento mais que especial ao Prof. Dr. Antonio Augusto Franco Garcia, pelo comprometimento em tornar-me um cientista. Lembro-me como se fosse hoje o nosso primeiro encontro em 2007, quando estive na ESALQ para manifestar interesse em sua orientação. Tive uma das melhores recepções da minha vida, que, de lá para cá, mudou e se transformou completamente. Faltam-me palavras para agradecer por tudo - acolhimento, abertura, orientação, ensinamentos, oportunidades, conversas extremamente produtivas, compreensão, paciência e amizade. Sinto-me honrado e satisfeito por ter feito parte do seu grupo e da sua história.

Outro agradecimento especial ao meu sogro Orlando Volpe Júnior, à minha sogra Maria Alice Masiero Volpe e aos meus cunhados Henrique Volpe e Ana Cláudia Francelin. Vocês foram a minha segunda família ao longo de todos esses anos de muito esforço e empenho.

Aos pesquisadores Dr. Jurandir Vieira Magalhães e Dr^a. Maria Marta Pastina, da Embrapa Milho e Sorgo (Sete Lagoas/MG). O apoio de vocês e a troca de ideias foram muito importantes para a utilização dos dados de sorgo disponíveis e a realização deste trabalho. Foi uma imensa honra ter tido o apoio de profissionais como vocês, que possuem competência e brilhantismo inigualáveis.

Ao Prof. Dr. Geoffrey Morris, do *Department of Agronomy, Crop Genetics & Genomics, Kansas State University*, por ter tornado público o conjunto de dados genotípicos do painel mundial de acessos de sorgo utilizado neste trabalho. Sem tais dados e informações históricas relevantes acerca da espécie, que só podem ter sido elaboradas por cientistas renomados de sorgo, não teríamos tido a motivação de realizar o presente estudo.

À Prof^a. Dr^a. Anete Pereira de Souza, da Universidade Estadual de Campinas (UNICAMP), pela concessão da bolsa de doutorado sanduíche realizado na França. Embora o projeto desenvolvido neste país não esteve diretamente relacionado a este trabalho, a oportunidade e a experiência de ter vivido em outro país, com suas crenças e seus costumes, foram incríveis. A experiência profissional adquirida pelo desenvolvimento do estudo por lá realizado foi algo inestimável para a minha carreira. Muito obrigado por ter acreditado em meu trabalho.

Ao Dr. Vincent Le Guen, do Centre de Coopération Internationale en Recherche Agronomique pour le Développement (CIRAD, Montpellier/França). Assim como aconteceu com o Prof. Dr. Antonio Augusto Franco Garcia, tive uma das melhores recepções da minha vida. Ter sido supervisionado por você durante o tempo em que

estive no CIRAD e na França foi uma honra e um grande prazer para mim. Muito obrigado por ter confiado em meu trabalho! E, além disso, tenho certeza que conquistei um grande amigo!

Ao Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPq), pela bolsa concedida durante o curso de doutorado (processo 141117/2012-5), e à Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES), pela concessão da bolsa de doutorado sanduíche (processo 3243/15-0) realizado na França.

A todos os colegas que fizeram e fazem parte do Laboratório de Genética Estatística da ESALQ/USP: (i) Aos mais *antigos*: Rodrigo, Maria Marta, Marcelo, Graciela, Gabriel, Renato e Edjane; (ii) Aos *contemporâneos*: Carina, Luciano, Guilherme, Adriana, Rodrigo Júnior, Maria Izabel e Letícia; (iii) Aos mais *novos*: Rafael, Felipe, Marianella, Letícia Lara, Danilo, Cristiane e Jhonathan; e (iv) Aos *agregados*: Amanda e Fernando. Cada um de vocês foi muito importante de alguma forma durante a minha caminhada. Alguns foram importantes por fornecer apoio pessoal e profissional, enquanto que outros foram importantes pela agradável convivência e por tornar o nosso ambiente de trabalho mais alegre e divertido.

Um agradecimento especial aos colegas mais próximos do Laboratório: Carina, Renato, Guilherme, Adriana e Felipe. Pela agradável convivência em todos os momentos e pelas conversas científicas inteligentes e extremamente produtivas. À Carina, que esteve ao meu lado desde o princípio e se tornou uma grande amiga. Ao Renato, com quem aprendi muito durante o tempo em que estive no laboratório e por quem tenho grande admiração. Ao Guilherme, pela conversas super inteligentes sobre ciência e pela competência. À Adriana, pelo exemplo de pessoa íntegra, amiga e profissional competente. Ao Felipe, pelo exemplo de dedicação e esforço em quaisquer circunstâncias.

Aos grandes amigos Fernando Guerra e Augusto Lima, com quem pude vivenciar grandes momentos durante o tempo em que estive em Piracicaba.

A todos os professores do Departamento de Genética e de Estatística e Experimentação Agronômica da ESALQ/USP. De uma forma ou de outra, vocês foram muito importantes pelos avanços alcançados durante o doutorado e pela confecção do presente trabalho.

Aos funcionários do Departamento de Genética da ESALQ/USP: Seu Antônio, Marcos, Valdir, Berdan, Léia, Macedônio, Fernandinho e Maídia, pelos grandiosos auxílios nos momentos decisivos. Aos funcionários da Biblioteca Central: Eliana, Cristina e Glória, pela ajuda na confecção final deste documento.

Finalmente, a todas as pessoas que direta ou indiretamente contribuíram para que eu chegasse até aqui e pudesse realizar este trabalho. Muito obrigado!

EPIGRAFE

“Talvez não tenhamos conseguido fazer o melhor
Mas, lutamos para o que o melhor fosse feito
Não somos o que deveríamos ser, não somos o que iremos ser
Mas, graças a Deus, já não somos mais o que éramos.”

Martin Luther King (15/01/1929 - 04/04/1968)

“To every complex problem
there is a simple solution,
and it is usually wrong.”

Henry Louis Mencken
(12/09/1880 - 29/01/1956)

“Se as coisas são inatingíveis...Ora!
Não é motivo para não querê-las.
Que tristes os caminhos se
não fora a presença distante das estrelas.”

Mário Quintana (30/07/1906 - 05/05/1994)

SUMÁRIO

Resumo	9
Abstract	10
1 Introdução	11
2 Revisão Bibliográfica	13
2.1 Desequilíbrio de Ligação	13
2.2 Teoria da Coalescência	14
2.2.1 Histórico e Conceitos Importantes	14
2.2.2 Estatísticas de Diversidade	19
2.2.3 Teoria Neutralista e Testes de Neutralidade	19
2.2.4 Modelos Evolutivos	22
2.2.4.1 Mutação	22
2.2.4.2 Estrutura Populacional e Fluxo Gênico	25
2.2.4.3 Recombinação	27
2.2.5 Máxima Verossimilhança e Inferência Bayesiana	28
2.3 Desequilíbrio de Ligação e Teoria da Coalescência	29
2.4 Sorgo	32
2.4.1 Classificação e Taxonomia	32
2.4.2 Centro de Origem e Domesticação	32
2.4.3 Importância e Melhoramento Genético	33
2.4.4 Alguns Genes de Importância	34
2.4.4.1 Exemplos Relevantes Para o Estudo	34
2.4.5 Marcadores Moleculares e Genotipagem Por Sequenciamento	36
3 Material e Métodos	39
3.1 Material Biológico e Dados Moleculares	39
3.1.1 Painel de Acessos	39
3.1.2 Genotipagem Por Sequenciamento e Imputação de Dados	39
3.1.3 Análises Prévias no Painel	40
3.1.4 Seleção das Regiões Genômicas	40
3.2 Diversidade Genética e Testes de Neutralidade	41
3.3 Modelos de Coalescência	41
3.3.1 Estrutura Populacional e Fluxo Gênico	41
3.3.2 Recombinação	43
3.3.2.1 Abordagem Completa	43
3.3.2.2 Abordagem Aproximada	43
3.4 Desequilíbrio de Ligação Via Coalescência	45
4 Resultados	47
4.1 Seleção dos Acessos de Sorgo	47

4.2	Seleção das Regiões Genômicas	47
4.3	Estatísticas de Diversidade e Testes de Neutralidade	48
4.4	Estrutura Populacional e Fluxo Gênico	50
4.5	Recombinação e Desequilíbrio de Ligação	51
5	Discussão	61
6	Considerações Finais	67
	Referências	69

RESUMO

Análise da estrutura populacional e do desequilíbrio de ligação de um painel de acessos de sorgo: uma abordagem usando teoria da coalescência

A estrutura populacional e o desequilíbrio de ligação são dois processos fundamentais para estudos evolutivos e de mapeamento associativo. Tradicionalmente, ambos têm sido investigados por meio de métodos clássicos comumente utilizados. Tais métodos certamente forneceram grandes avanços no entendimento dos processos evolutivos das espécies. No entanto, em geral, nenhum deles utiliza uma visão genealógica de forma a considerar eventos genéticos ocorridos no passado, dificultando a compreensão dos padrões de variação observados no presente. Uma abordagem que possibilita a investigação retrospectiva com base no atual polimorfismo observado é a teoria da coalescência. Assim, o objetivo deste trabalho foi analisar, com base na teoria da coalescência, a estrutura populacional e o desequilíbrio de ligação de um painel mundial de acessos de sorgo (*Sorghum bicolor*). Para tanto, análises de mutação, migração com fluxo gênico e recombinação foram realizadas para cinco regiões genômicas relacionadas à altura de plantas e maturidade (*Dw1*, *Dw2*, *Dw4*, *Ma1* e *Ma3*) e sete populações previamente selecionadas. Em geral, elevado fluxo gênico médio ($M = \frac{m}{\mu} = 41,78 - 52,07$) foi observado entre as populações considerando cada região genômica e todas elas simultaneamente. Os padrões sugeriram intenso intercâmbio de acessos e história evolutiva específica para cada região genômica, mostrando a importância da análise individual dos locos. A quantidade média de migrantes por geração (M) não foi simétrica entre pares recíprocos de populações, de acordo com a análise individual e simultânea das regiões. Isso sugere que a forma pela qual as populações se relacionaram e continuam interagindo evolutivamente não é igual, mostrando que os métodos clássicos utilizados para investigar estrutura populacional podem ser insatisfatórios. Baixas taxas médias de recombinação ($\rho_L = 2N_e r = 0,030 - 0,246$) foram observadas utilizando o modelo de recombinação constante ao longo da região. Baixas e altas taxas médias de recombinação ($\rho_r = 2N_e r = 0,060 - 3,395$) foram estimadas utilizando o modelo de recombinação variável ao longo da região. Os métodos tradicional (r^2) e via coalescência ($E[r_{\text{rhomap}}^2]$) utilizados para a estimação do desequilíbrio de ligação mostraram resultados próximos para algumas regiões genômicas e populações. No entanto, o r^2 sugeriu padrões descontínuos de desequilíbrio em várias ocasiões, dificultando o entendimento e a caracterização de possíveis blocos de associação. O método via coalescência ($E[r_{\text{rhomap}}^2]$) forneceu resultados que pareceram ter sido mais consistentes, podendo ser uma estratégia eventualmente importante para um refinamento dos padrões não-aleatórios de associação. Os resultados aqui encontrados sugerem que o mapeamento genético a partir de um único *pool* gênico pode ser insuficiente para detectar associações causais importantes para características quantitativas em sorgo.

Palavras-chave: Passado Evolutivo, Fluxo Gênico, Recombinação, Associações Não-aleatórias, Mapeamento Associativo, *Sorghum bicolor*

ABSTRACT

Analysis of population structure and linkage disequilibrium of a sorghum accession panel: an approach using coalescent theory

Population structure and linkage disequilibrium are two fundamental processes for evolution and association mapping studies. Traditionally, both have been investigated using classical methods that are commonly used. These methods certainly provided important advances for the understanding of the evolution processes of the species. However, in general, none of them uses a genealogical view to consider genetic events occurred in the past, making difficult the understanding of the variation patterns observed in the present. An approach that enables the retrospective investigation based on the actual observed polymorphism is the coalescent theory. Here, we used the coalescent theory to analyze the population structure and linkage disequilibrium of a worldwide sorghum (*Sorghum bicolor*) accession panel. To reach this purpose, analyses of mutation, migration with gene flow and recombination were performed to five genomic regions related to plant height and maturity (*Dw1*, *Dw2*, *Dw4*, *Ma1* e *Ma3*) and seven previously selected populations. In general, high average gene flow ($M = \frac{m}{\mu} = 41,78 - 52,07$) was observed between populations considering each genomic region and all the regions simultaneously. The patterns suggested a high exchange of accessions between populations and a specific evolutionary history for each genomic region, showing that the individual analysis of each locus was important. The average number of migrants per generation (M) was not symmetric between reciprocal pairs of populations, according to the specific and simultaneous analyses of the regions. This result suggests that the historical and recent evolutionary relations between populations are not equal, showing that the classical methods to investigate population structure may be unsatisfactory. Low average recombination rates ($\rho_L = 2N_e r = 0,030 - 0,246$) were observed using a constant recombination model along the region. Low and high average recombination rates ($\rho_r = 2N_e r = 0,060 - 3,395$) were estimated using a variable recombination model along the region. Both traditional (r^2) and coalescent ($E[r_{\text{rhomap}}^2]$) methods for the estimation of linkage disequilibrium showed similar results for some genomic regions and populations. However, r^2 suggested discontinuous patterns of linkage disequilibrium in several cases, making difficult the understanding and definition of the association blocks. The coalescent method ($E[r_{\text{rhomap}}^2]$) provided results that seemed to be more consistent and could be an eventually important strategy to refine the non-random association patterns. The results detected here suggest that the genetic mapping from a unique gene *pool* may be insufficient to detect important causal associations for quantitative traits in sorghum.

Keywords: Evolutionary Past, Gene Flow, Recombination, Non-random Associations, Association Mapping, *Sorghum bicolor*

1 INTRODUÇÃO

Tradicionalmente, a estrutura populacional em painéis de associação e populações naturais tem sido investigada com base em métodos clássicos bastante conhecidos. Basicamente, tais métodos consistem em agrupamentos hierárquicos (ODONG *ET AL.*, 2011), inferência bayesiana (STRUCTURE) (PRITCHARD *ET AL.*, 2000), componentes principais (ZHU *ET AL.*, 2008), F_{ST} (WRIGHT, 1931) e suas extensões (NATH and GRIFFITHS, 1996). Certamente, grandes contribuições foram geradas pelo uso de tais abordagens. No entanto, em geral, todos eles utilizam as frequências alélicas dos locos para capturar níveis de diferenciação que podem ser o resultado de eventos de fluxo gênico recentes, pouco considerando eventos migratórios ocorridos no passado (WILSON and RANNALA, 2003; HÜBNER *ET AL.*, 2012). Além disso, esses métodos não fornecem uma estimativa de fluxo gênico com a qual é possível conhecer a magnitude da relação histórica entre populações, dificultando o entendimento do passado evolutivo. Assim, uma visão genealógica da estrutura populacional e do fluxo gênico é de grande motivação, sobretudo para fins de mapeamento associativo, que é o foco do presente trabalho.

Da mesma forma, o desequilíbrio de ligação presente em populações tem sido analisado com base em medidas relativas bastante conhecidas, tais como D' (LEWONTIN, 1964) e r^2 (HILL and ROBERTSON, 1968). Normalmente, tais medidas apresentam grandes variações em suas estimativas (NORDBORG and TAVARÉ, 2002), dificultando uma clara interpretação dos padrões de associação (EVANS and CARDON, 2005). Além disso, D' e r^2 apresentam limitações por não considerarem associações entre múltiplos locos simultaneamente e por não permitirem explicações que mostrem as possíveis razões de locos estarem preferencialmente associados ao longo das gerações (PRITCHARD and PRZEWORSKI, 2001; SLATKIN, 2008). Dessa forma, uma visão genealógica do desequilíbrio de ligação, com atenções voltadas ao mapeamento associativo, é bastante desejável (NORDBORG and TAVARÉ, 2002; ZÖLLNER and PRITCHARD, 2005).

Uma das estratégias que podem ser utilizadas para se obter uma visão genealógica de processos evolutivos é a teoria da coalescência (KINGMAN, 1982a,b,c). A abordagem proposta por esta teoria possibilita investigar diferentes genealogias caminhando-se do presente em direção ao passado, num processo retrospectivo (FU and LI, 1999). Tais genealogias podem ser investigadas de forma a considerar processos evolutivos importantes, tais como mutação, fluxo gênico e recombinação (HUDSON, 1990; PRITCHARD and PRZEWORSKI, 2001; HEIN *ET AL.*, 2004; MARJORAM and TAVARÉ, 2006; NORDBORG, 2007), fornecendo informações relevantes acerca do passado. No contexto do mapeamento associativo, a teoria da coalescência poderá ser importante para se ter maior compreensão dos padrões de associação entre locos, o que, em última análise, será fundamental para o entendimento das relações entre genótipo e fenótipo.

Em plantas domesticadas, o uso da coalescência parece ser bastante restrito. Em sorgo (*Sorghum bicolor*), tal abordagem foi utilizada para estimar taxas de recombinação populacional em regiões genômicas específicas (HAMBLIN *ET AL.*, 2005) ou ao longo de todo o genoma (MORRIS *ET AL.*, 2013). Contudo, embora esses trabalhos tenham gerado resultados importantes e fornecido uma ideia dos padrões de associação, nenhum deles direcionou atenções para o cálculo do desequilíbrio de ligação a partir da recombinação via coalescência. Veremos mais adiante que o desequilíbrio de ligação pode ser diretamente obtido a partir das taxas de recombinação. Além disso, até o nosso conhecimento, nenhum trabalho foi publicado para o sorgo considerando a estrutura populacional e o fluxo gênico estimados via coalescência.

Neste contexto, o objetivo deste trabalho foi analisar, com base na teoria da coalescência, a estrutura populacional e o desequilíbrio de ligação presentes em um painel mundial de acessos de sorgo. Para alcançar tal objetivo, análises de mutação, migração com fluxo gênico e recombinação foram realizadas para diferentes regiões genômicas e populações que constituem o presente painel. Buscou-se com isso um maior entendimento dos processos evolutivos ocorridos no passado, que podem ser cruciais para o mapeamento genético de características quantitativas importantes. Acredita-se que os resultados aqui gerados possam fornecer informações relevantes para estudos futuros de mapeamento associativo em sorgo.

2 REVISÃO BIBLIOGRÁFICA

Basicamente, a revisão bibliográfica está estruturada em três partes. Na primeira parte, serão abordados conceitos gerais e importantes sobre o desequilíbrio de ligação, que é um dos objetivos do presente trabalho. Não abordaremos os métodos clássicos utilizados para inferir estrutura populacional por acreditarmos que tais métodos estão amplamente discutidos na ciência. Na segunda parte, serão abordados conceitos importantes e os modelos evolutivos da teoria da coalescência, incluindo estrutura populacional (migração com fluxo gênico) e recombinação. Por fim, na terceira parte será apresentado um panorama geral sobre o sorgo, relatando desde a sua origem e domesticação até as estratégias adotadas pelo melhoramento genético, descrevendo alguns genes importantes que serão utilizados no presente estudo. Concluímos a revisão bibliográfica com uma rápida descrição sobre a genotipagem por sequenciamento, que também será utilizada neste trabalho.

2.1 Desequilíbrio de Ligação

Desequilíbrio de ligação (DL) é qualquer desvio das frequências genotípicas em relação às frequências esperadas sob independência, indicando associação preferencial entre alelos de diferentes locos em uma população (LEWONTIN and KOJIMA, 1960). Apesar do termo ter tido grande aceitação por parte da comunidade científica, DL não corresponde necessariamente à ligação física (FLINT-GARCIA *ET AL.*, 2003; TEMPLETON, 2006; ALLENDORF and LUIKART, 2007; HAMILTON, 2009; HEDRICK, 2010). DL pode ocorrer tanto entre locos ligados quanto entre locos não ligados (PRITCHARD and PRZEWORSKI, 2001; ARDLIE *ET AL.*, 2002), enquanto que ligação física está relacionada a alelos que cosegregam no mesmo cromossomo ao longo das gerações (FLINT-GARCIA *ET AL.*, 2003; GUPTA *ET AL.*, 2005). Assim, é fundamental que esses conceitos não sejam confundidos, já que abordagens importantes baseiam-se no DL por ligação física, como é o caso dos estudos de mapeamento genético.

Grande parte da teoria sobre DL é descrita para pares de locos bialélicos (SLATKIN *ET AL.*, 1996). Embora o DL esteja fundamentado no conceito proposto por LEWONTIN and KOJIMA (1960), designado por D , outras medidas, como D' (LEWONTIN, 1964) e r^2 (HILL and ROBERTSON, 1968), têm sido muito utilizadas. No entanto, devido à obtenção frequente de amostras pequenas (FLINT-GARCIA *ET AL.*, 2003), a segunda delas tem sido mais amplamente utilizada, como pode ser verificado nos estudos de mapeamento associativo em humanos (PRITCHARD and PRZEWORSKI, 2001; SLATKIN, 2008) e plantas (FLINT-GARCIA *ET AL.*, 2003; GUPTA *ET AL.*, 2005). De qualquer forma, espera-se que o DL seja melhor compreendido no contexto de pares de locos multialélicos e múltiplos locos bialélicos (KIM *ET AL.*, 2008), o que certamente fornecerá informações mais refinadas para estudos de mapeamento genético (McVEAN, 2007a).

Vários fatores populacionais afetam o DL. Mutação, deriva genética, recombinação, estrutura populacional, fluxo gênico, endogamia e seleção podem ser considerados os principais fatores que influenciam no aumento ou na diminuição do DL entre locos (ARDLIE *ET AL.*, 2002; FLINT-GARCIA *ET AL.*, 2003; RAFALSKI and MORGANTE, 2004; GUPTA *ET AL.*, 2005; SLATKIN, 2008). Mutação, deriva genética, estrutura populacional, fluxo gênico, endogamia e seleção atuam no surgimento de associações preferenciais entre locos; já a recombinação atua na sua redução ao longo das gerações, algo que é proporcional à distância genética entre os mesmos (HARTL and CLARK, 2007; HEDRICK, 2010). Assim, quanto menor for essa distância, menor será a chance de recombinação e, com ela, a quebra do DL entre locos (PRITCHARD and PRZEWORSKI, 2001). No contexto do mapeamento associativo, o

objetivo é detectar locos intimamente ligados e que estejam em DL, sugerindo uma mesma genealogia ao longo das gerações (YU and BUCKLER, 2006; ZHU *ET AL.*, 2008; ABDURAKHMONOV and ABDUKARIMOV, 2008).

De maneira geral, o DL ocorre de forma aleatória no genoma e depende da espécie e população sob estudo (ORAGUZIE *ET AL.*, 2007). Assim, o estudo da extensão do DL é muito importante, e possibilita determinar a resolução de mapeamento e a densidade de marcadores necessários na identificação de genes de interesse (ORAGUZIE *ET AL.*, 2007). Nesse sentido, se o DL permanecer em curtas distâncias no genoma, uma alta resolução de mapeamento será esperada, mas uma elevada quantidade de marcadores será necessária. Em contrapartida, se o DL se estender a maiores distâncias no genoma, a resolução de mapeamento tenderá a ser baixa, mas uma quantidade menor de marcadores será requerida (ZHU *ET AL.*, 2008; ROSA, 2011).

Apesar dessas informações serem relevantes para a realização de estudos de mapeamento associativo, é desejável obter uma visão do DL sob o ponto de vista genealógico (NORDBORG and TAVARÉ, 2002; ZÖLLNER and VON HAESLER, 2000; McVEAN, 2002; McVEAN and CARDIN, 2005; McVEAN, 2007a), considerando as possíveis relações evolutivas entre sequências de DNA. Nesse sentido, a teoria da coalescência, que realiza inferências do presente em direção ao passado (HUDSON, 1990; FU and LI, 1999; WAKELEY, 2009), poder ser bastante útil para a caracterização do DL em painéis de associação, principalmente para espécies que possuem poucas informações de genealogias. Mais do que isso, pode ser consideravelmente útil em relação às medidas tradicionais do DL, tais como D' e r^2 , que têm mostrado grande variação nos padrões de associação (NORDBORG and TAVARÉ, 2002; EVANS and CARDON, 2005).

2.2 Teoria da Coalescência

2.2.1 Histórico e Conceitos Importantes

A teoria da coalescência foi formalmente proposta por John Frank Charles Kingman (Isaac Newton Institute of Mathematical Sciences, Cambridge University, Inglaterra) no início da década de 80. Três trabalhos seminais foram publicados em 1982 com os princípios teóricos sobre esta abordagem (KINGMAN, 1982a,b,c), que passou a ter enorme representatividade para os estudos envolvendo genética de populações e evolução. Com os avanços obtidos na área da biologia molecular e computação, a teoria da coalescência passou a ser amplamente utilizada em trabalhos recentes de evolução (KINGMAN, 2000). Por ser uma abordagem sofisticada do ponto de vista matemático e estatístico, a teoria da coalescência tem fornecido resultados bastante precisos e importantes em estudos evolutivos. Nos dias de hoje, é considerada por vários cientistas como a moderna (PALCZEWSKI and BEERLI, 2013) e mais popular (MARJORAM and TAVARÉ, 2006) genética de populações, a qual revolucionou o entendimento sobre os processos histórico-evolutivos (HOLSINGER and WEIR, 2009).

Em linhas gerais, o modelo de coalescência é um processo estocástico com o qual genealogias históricas e desconhecidas podem ser modeladas a partir de um conjunto de sequências de determinada população (FU and LI, 1999; ROSENBERG and NORDBORG, 2002; NORDBORG and TAVARÉ, 2002; MARJORAM and TAVARÉ, 2006). Ao contrário dos métodos clássicos utilizados na genética de populações, o método da coalescência é baseado em inferências retrospectivas, partindo-se do presente em direção ao passado (FU and LI, 1999; HEIN *ET AL.*, 2004; MARJORAM and TAVARÉ, 2006). Assim, é possível desvendar os processos genéticos do passado que relacionam sequências de gerações presentes, compreendendo a história evolutiva ocorrida ao longo das gerações (HUDSON, 1990; FU and

LI, 1999; WAKELEY, 2009). Portanto, a teoria da coalescência sugere que os fenômenos podem ser explicados e compreendidos com base nos eventos que já ocorreram e não naqueles que ainda poderão ocorrer (MARJORAM and JOYCE, 2011).

Para melhor compreender os conceitos e as motivações da teoria da coalescência, considere um conjunto de n sequências de DNA correspondentes à determinada região do genoma (Figura 2.1). Essas sequências representam uma amostra aleatória de indivíduos coletados recentemente em determinada população. Inicialmente, a história evolutiva que explica os padrões de variação envolvendo quatro sítios polimórficos é totalmente desconhecida. Assim, é de interesse desvendar esta história, sendo de grande motivação o uso da teoria da coalescência.

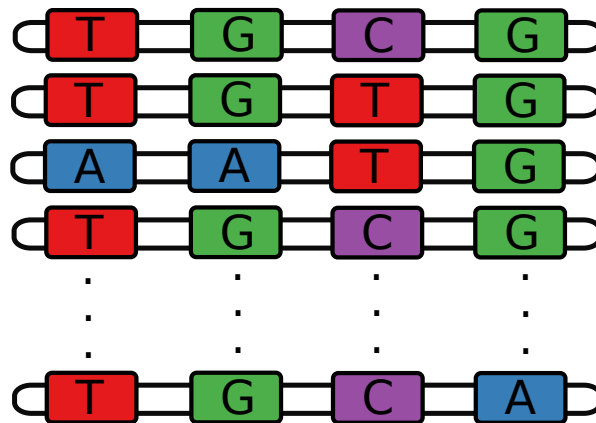


Figura 2.1: Amostra populacional contendo n sequências de DNA e quatro sítios polimórficos para determinada região do genoma

Para o modelo de coalescência proposto por KINGMAN (1982a,b,c), no caso de haver mutação, considera-se que os padrões de polimorfismo observados podem ser explicados pelo *modelo de sítios infinitos* (KARLIN and MCGREGOR, 1967; KIMURA, 1969; WATTERSON, 1975), com o qual um número infinito de locos, sujeitos à mutação, pode ocorrer em posições contínuas no genoma. Neste caso, cada um dos quatro sítios polimórficos é resultado de um único evento de mutação ocorrido ao longo do processo evolutivo, gerando uma forma variante nunca antes observada (GRIFFITHS and TAVARÉ, 1995; WAKELEY, 2009). Trata-se de modelo apropriado para dados de sequências, cuja probabilidade de ocorrer mutações repetidas para uma mesma posição no genoma é negligenciável (FEARNHEAD and DONNELLY, 2001).

A história evolutiva do segundo sítio polimórfico da amostra (Figura 2.1) é mostrada na Figura 2.2. Este sítio será utilizado para demonstração dos conceitos de coalescência, embora as inferências de genealogia aconteçam para todos os sítios simultaneamente, considerando os dados de sequências.

Ao retornar ao passado, é possível detectar ancestrais comuns, ou eventos *coalescentes*, envolvendo duas sequências quaisquer da amostra até alcançar o chamado *ancestral comum mais recente* (MRCA, do inglês *Most Recent Common Ancestor*), envolvendo todas as sequências da amostra (FU and LI, 1999; HEIN ET AL., 2004; NORDBOG and TAVARÉ, 2002; WAKELEY, 2009). Um evento coalescente corresponde à fusão de duas sequências quaisquer presentes no tempo t em uma única sequência qualquer presente no tempo $t - 1$, que teve, numa visão prospectiva, seu DNA duplicado para transmitir a informação genética para sequências de gerações atuais. Para uma amostra contendo n sequências de DNA, $n - 1$ eventos de coalescência devem estar presentes caminhando-se do presente em direção ao passado, sendo que o último evento coalescente corresponderá sempre ao MRCA

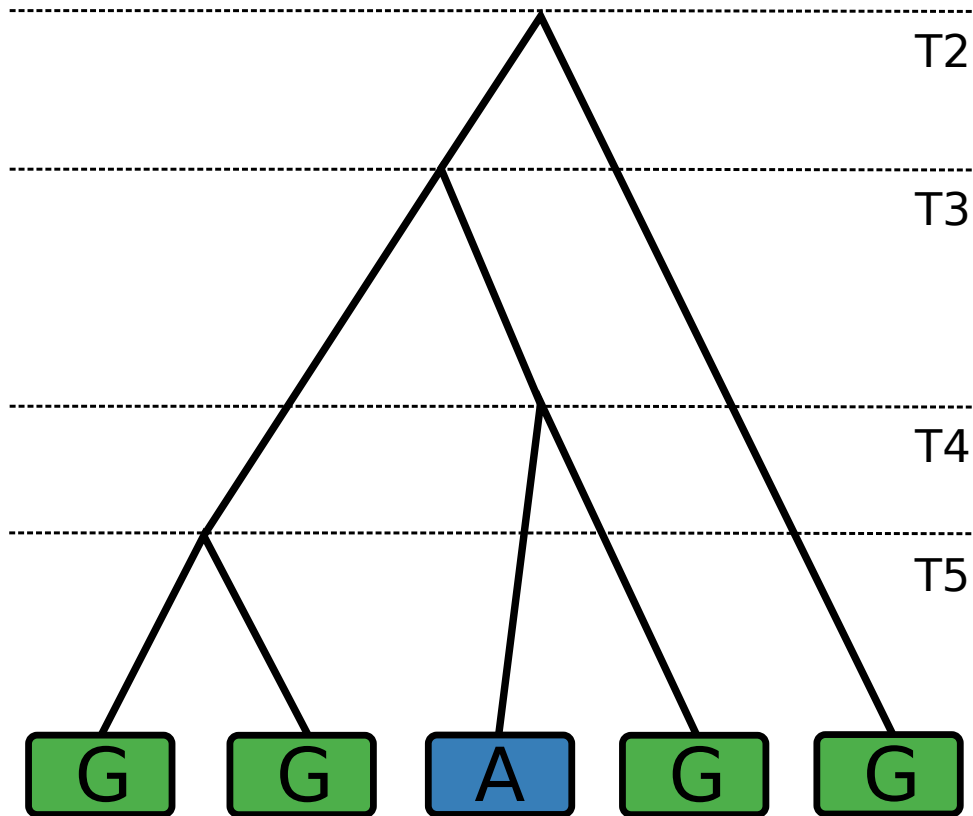


Figura 2.2: Possível história evolutiva do segundo sítio polimórfico (Figura 2.1) com base na teoria da coalescência. Cinco sequências de DNA estão representadas na geração atual, sendo que quatro delas possuem o nucleotídeo G (primeira, segunda, quarta e quinta) e uma única sequência possui o nucleotídeo A (terceira). Caminhando-se do presente em direção ao passado, o primeiro evento de coalescência acontece entre a primeira e a segunda sequências (à esquerda) que possuem o nucleotídeo G, no intervalo de tempo T5. O segundo evento de coalescência acontece entre a terceira (nucleotídeo A) e a quarta (nucleotídeo G) sequências no intervalo de tempo T4. O terceiro evento de coalescência acontece entre as sequências ancestrais obtidas a partir do primeiro e do segundo eventos de coalescência, no intervalo de tempo T3. E, por fim, o quarto e último evento de coalescência acontece entre a sequência ancestral das quatro primeiras sequências e a quinta sequência (nucleotídeo G), no intervalo de tempo T2. O último evento de coalescência corresponde ao ancestral comum mais recente da amostra (MRCA). Note que a terceira sequência é a única que possui o nucleotídeo A, mostrando que um evento de mutação deve ter ocorrido no ramo em que tal sequência se fundiu para formar o segundo evento de coalescência. O cálculo de uma probabilidade fornece o quanto que esta história evolutiva explica a variação dos dados da geração atual

(NORDBORG, 2007; WAKELEY, 2009). É importante notar que os eventos coalescentes correspondem aos eventos de divergência numa visão prospectiva da história evolutiva das populações (FU and LI, 1999). No exemplo do segundo sítio polimórfico mostrado na Figura 2.2, quatro eventos coalescentes podem ser observados a partir de uma amostra contendo cinco sequências de DNA. Esses eventos estão situados nas posições exatas das linhas tracejadas mostradas horizontalmente no gráfico.

O primeiro evento coalescente ocorrerá entre duas das n sequências presentes na amostra. Uma vez que a genealogia histórica é totalmente desconhecida, o modelo de coalescência levará em consideração a possibilidade de todas as n sequências estarem envolvidas com o primeiro evento coalescente. Assim, considerando todos os possíveis pares envolvendo n sequências, um total de $\frac{n(n-1)}{2}$ combinações será possível para o primeiro evento coalescente, e o mesmo raciocínio pode ser utilizado para os demais eventos de coalescência. No modelo proposto, a probabilidade de duas sequências quaisquer coalescerem em alguma geração do passado é $\frac{1}{2N_e}$, sendo N_e o tamanho efetivo populacional. A probabilidade de duas sequências quaisquer não coalescerem em alguma geração

do passado é $1 - \frac{1}{2N_e}$. Dessa forma, a probabilidade de ocorrência do primeiro evento coalescente é $\frac{n(n-1)}{2} \frac{1}{2N_e}$, levando em consideração todas as possíveis combinações envolvendo pares de sequências. É importante notar que todos esses pares de sequências são *equiprováveis* neste modelo de coalescência, algo que nem sempre acontece em cenários mais complexos (HUDSON, 1990; WAKELEY, 2009).

Os eventos de coalescência são estatisticamente modelados seguindo um processo de Poisson (distribuição discreta), com taxa de 1 evento de coalescência ($\lambda = 1$) para cada par de sequências e intervalo de tempo (T_2, T_3, \dots, T_n). Os tempos de coalescência, que correspondem aos intervalos entre coalescentes sucessivos, são estatisticamente modelados seguindo uma distribuição exponencial (distribuição contínua) (HUDSON, 1990). Assim, a abordagem de coalescência é proposta de forma modificada em relação aos modelos que explicam uma população finita (HEIN ET AL., 2004; MARJORAM and TAVARÉ, 2006), na qual o tempo é modelado de forma discreta (em gerações) e não contínua. Neste caso, os coalescentes e os tempos de coalescência são modelados pelas distribuições binomial e geométrica, respectivamente (HUDSON, 1990; WAKELEY, 2009). Verifica-se que a inferência de certa genealogia a partir da teoria da coalescência baseia-se em diferentes modelos probabilísticos.

As diferentes formas de se obter eventos de coalescência determinam o número de genealogias (NG) ou configurações possíveis envolvendo sequências de gerações atuais. A equação 2.1 mostrada a seguir ilustra o processo de Poisson com o qual o número de genealogias possíveis pode ser modelado (WAKELEY, 2009):

$$NG = \prod_{i=2}^n \binom{i}{2} = \binom{2}{2} \times \binom{3}{2} \times \dots \times \binom{n}{2} \quad (2.1)$$

em que o produtório ($\prod_{i=2}^n$) indica que as genealogias dos diferentes tempos de coalescência são modeladas independentemente (WAKELEY, 2009).

Como foi dito, cada tempo de coalescência pode ser explicado por uma distribuição exponencial. Um modelo para todos os tempos de coalescência, até o momento em que o MRCA da amostra é alcançado, segue, portanto, uma distribuição exponencial conjunta. Para este caso, a função densidade conjunta é obtida pelo produto das funções correspondentes aos intervalos de coalescência, assumindo que os mesmos sejam independentes (HUDSON, 1990; HEIN ET AL., 2004; WAKELEY, 2009). Portanto, a forma pela qual o $(n-1)$ -ésimo evento coalescente ocorreu no tempo $t-1$ não influenciará na forma pela qual o $(n-2)$ -ésimo evento ocorrerá no tempo $t-2$, e assim sucessivamente. As funções de densidade exponencial para cada tempo de coalescência e todos os tempos conjuntamente são apresentadas pelas equações 2.2 e 2.3, respectivamente (WAKELEY, 2009):

$$f_{T_i}(t_i) = \binom{i}{2} e^{-\binom{i}{2}t_i} \quad (2.2)$$

$$f_{T_2, \dots, T_n}(t_2, \dots, t_n) = \prod_{i=2}^n f_{T_i}(t_i) \quad (2.3)$$

$$= f_{T_2}(t_2) \times f_{T_3}(t_3) \times \dots \times f_{T_n}(t_n) \quad (2.4)$$

$$= \binom{2}{2} e^{-\binom{2}{2}t_2} \times \binom{3}{2} e^{-\binom{3}{2}t_3} \times \dots \times \binom{n}{2} e^{-\binom{n}{2}t_n} \quad (2.5)$$

em que $T_i(t_i)$ é o tempo ou intervalo de coalescência quando i sequências estão presentes.

É possível obter a esperança ($E[T_i]$) e a variância ($Var[T_i]$) dos tempos de coalescência (equações 2.6 e 2.7, respectivamente):

$$E[T_i] = \frac{4N_e}{i(i-1)} \quad (2.6)$$

$$\text{Var}[T_i] = \left(\frac{4N_e}{i(i-1)} \right)^2 \quad (2.7)$$

Claramente, o tempo para ocorrência de um evento coalescente (T_i) é função direta do tamanho efetivo populacional e do tamanho da amostra ($i = 2, 3, \dots, n$). Assim, os primeiros eventos coalescentes acontecerão, em média, muito rapidamente no passado, e os coalescentes mais antigos demorarão cada vez mais à medida que o MRCA se aproxima (WAKELEY, 2009). Para modelos de coalescência mais complexos, em que há alteração do tamanho efetivo populacional, outras situações podem ocorrer para a história evolutiva da amostra. Por exemplo, para uma população que sofreu expansão do seu tamanho efetivo ao longo das gerações, espera-se que os primeiros eventos coalescentes tenham ocorrido há muitas gerações no passado, e os coalescentes mais antigos tenham ocorrido em menor tempo em relação aos primeiros.

Duas medidas importantes do tamanho de uma genealogia são: (i) o tempo para alcançar o MRCA (T_{MRCA} , equação 2.8); e (ii) o comprimento total de todos os ramos que constituem a genealogia (T_{TOTAL} , equação 2.9):

$$T_{\text{MRCA}} = \sum_{i=2}^n T_i \quad (2.8)$$

$$T_{\text{TOTAL}} = \sum_{i=2}^n iT_i \quad (2.9)$$

Para ambas as medidas, é possível calcular a esperança (equações 2.10 e 2.11) e a variância (equações 2.12 e 2.13), considerando que diferentes genealogias podem ser formadas:

$$E[T_{\text{MRCA}}] = \sum_{i=2}^n E[T_i] = \sum_{i=2}^n \frac{2}{i(i-1)} \quad (2.10)$$

$$E[T_{\text{TOTAL}}] = \sum_{i=2}^n iE[T_i] = \sum_{i=2}^n i \frac{2}{i(i-1)} = 2 \sum_{i=1}^{n-1} \frac{1}{i} \quad (2.11)$$

$$\text{Var}[T_{\text{MRCA}}] = 8 \sum_{i=2}^n \frac{1}{i^2} - 4 \left(1 - \frac{1}{n} \right)^2 \quad (2.12)$$

$$\text{Var}[T_{\text{TOTAL}}] = 4 \sum_{i=1}^{n-1} \frac{1}{i^2} \quad (2.13)$$

KINGMAN (1982a,b,c) mostrou que uma população finita e muito grande (tendendo ao infinito), que segue o modelo proposto por FISHER (1930) e WRIGHT (1931), pode ser aproximadamente explicada pelo modelo de coalescência. Neste caso, tal modelo assume que as populações apresentam tamanho populacional finito e constante ao longo das gerações, bem como se reproduzem por meio de panmixia, que corresponde aos cruzamentos aleatórios que acontecem entre indivíduos de uma mesma geração (HEIN ET AL., 2004; WAKELEY, 2009). A coalescência também leva em consideração os princípios da teoria neutralista da evolução molecular (KIMURA, 1955; HUDSON,

1990; Fu and Li, 1999; Nordborg, 2007), com a qual alelos de determinado loco não apresentam diferenças em seus desempenhos e, portanto, são seletivamente neutros. Assim, apenas deriva genética é responsável pelas alterações alélicas ao longo das gerações (Hahn, 2008; Wakeley, 2009), de modo que seleção natural não seria força evolutiva preponderante nas populações (Hahn, 2008; Wakeley, 2009).

No entanto, os pressupostos de Fisher (1930) e Wright (1931) e da teoria neutralista nem sempre podem ser utilizados para explicar os padrões de variação. Certas populações podem apresentar tamanho efetivo populacional variável ao longo das gerações, subdivisão geográfica, recombinação, seleção e outros processos que alteram a estrutura das genealogias e causam desvios dos pressupostos iniciais. Assim, é comum o uso de testes estatísticos que podem detectar esses desvios e sugerir, pelo menos em princípio, cenários evolutivos mais complexos. São os chamados testes de neutralidade, que são baseados em estatísticas de diversidade que capitalizam a variação genética. Tais estatísticas serão apresentadas a seguir, e os testes de neutralidade mais conhecidos serão abordados na sequência.

2.2.2 Estatísticas de Diversidade

Historicamente, três importantes estatísticas foram propostas para explicar padrões da variação genética (Wakeley, 2009). A primeira delas é o número de sítios segregantes ou polimórficos (S) presentes em uma amostra de sequências (Watterson, 1975). A segunda é o número de diferenças nucleotídicas envolvendo pares de sequências (π) presentes nesta amostra (Tajima, 1983). E, por fim, a terceira é o número de sítios polimórficos (η_i) presentes na amostra de acordo com o número de cópias de suas bases nucleotídicas (i cópias de certa base e $n - i$ cópias da outra base), correspondendo à frequência nucleotídica dos sítios polimórficos (Wakeley, 2009). Para esta última, normalmente é impossível saber quais bases são mutantes (ξ_i) e quais são ancestrais (ξ_{n-i}) para os sítios, levando em consideração a ocorrência de um único evento de mutação para cada sítio (*modelo de sítios infinitos*) ao longo da história da amostra (Hein et al., 2004; Wakeley, 2009).

Embora S , π e η_i estejam diretamente relacionadas (as duas primeiras podem ser diretamente representadas a partir da terceira), cada qual sumariza a informação genética contida em uma amostra de diferentes formas. Assim, todas essas estatísticas têm sido historicamente utilizadas em diversos estudos de genética de populações, fornecendo evidências importantes acerca da variação. Além disso, pelo fato dos seus valores serem facilmente obtidos até mesmo para amostras grandes e com muitos sítios polimórficos, a motivação para o seu uso é ainda maior. Contudo, todas elas sumarizam a variação genética para um único valor, o que pode ser aceitável e desejável para certos casos, mas pode ser insuficiente quando o objetivo é compreender de que forma a variação genética foi estruturada diante de cenários evolutivos mais complexos (Kuhner, 2009; Wakeley, 2009).

2.2.3 Teoria Neutralista e Testes de Neutralidade

Os primeiros modelos desenvolvidos para explicar a variação genética presente em populações foram propostos por R. Fisher, S. Wright e J. Haldane em meados da década de 30 (Fisher, 1930; Wright, 1931). As ideias sugeridas por esses cientistas, com base nos modelos propostos, eram basicamente que as variações genéticas presentes nas populações aconteceriam em decorrência de eventos de mutação e seleção natural, sendo este último o fator que deveria, de fato, governar a evolução. Na época, esses cientistas revolucionaram pensamentos duvidosos e contrários a respeito do darwinismo e do mendelismo, mostrando que a seleção natural e a genética deveriam ser

encaradas conjuntamente (FUTUYMA, 2005).

No final dos anos 50, Motoo Kimura e colaboradores, com base em dados de proteínas, propuseram que grande parte da variação genética presente em populações finitas deveria ser neutra (KIMURA, 1955). Uma vez que mutações deletérias deveriam ser rapidamente eliminadas por seleção natural e mutações favoráveis mostrariam-se raras na população, esses cientistas concluíram que a grande maioria da variação não deveria exercer influências sobre o valor adaptativo dos indivíduos da população. Assim, alelos de locos neutros não mostrariam diferenças em seus desempenhos (HAHN, 2008) e, portanto, seriam seletivamente neutros, de modo que maiores ou menores proporções de um ou outro alelo não acarretariam maiores ou menores adaptações do ponto de vista evolutivo.

A teoria neutralista da evolução molecular, que passou a ser aplicada anos mais tarde, tornou-se um dos principais alicerces dos geneticistas e evolucionistas na tentativa de explicar a variação genética observada dentro e entre populações (HAHN, 2008). De maneira geral, o neutralismo fornece a base teórica para o entendimento da variação ao longo do genoma, possibilitando a elaboração de hipóteses testáveis e abordagens estatísticas que promovem a distinção entre seleção e deriva genética (KREITMAN, 2000; NIELSEN, 2001; HAHN, 2007).

Neste contexto, três diferentes testes estatísticos foram historicamente desenvolvidos para investigar desvios do modelo padrão baseado em neutralismo. O primeiro e mais conhecido teste de neutralidade foi proposto por TAJIMA (1989) no final da década de 80, sendo aqui designado como estatística D . Anos mais tarde, Fu and Li (1993) propuseram dois outros testes de neutralidade, designados por D^* e F^* . Embora D , D^* e F^* apresentem suas particularidades metodológicas, todas essas estatísticas foram propostas a partir de uma mesma motivação, sendo baseadas nas estatísticas de diversidade S , π e η_i .

A estatística D pode ser representada pela seguinte equação:

$$D = \frac{\pi - \frac{S}{a_1}}{\sqrt{\widehat{\text{Var}}[\pi - \frac{S}{a_1}]}} \quad (2.14)$$

em que $a_1 = \sum_{i=1}^{n-1} i$. Uma vez que a esperança das estatísticas π e S são iguais, respectivamente, a $E[\pi] = \theta$ e $E[S] = \theta a_1$, a diferença $\pi - \frac{S}{a_1}$ será igual a zero se ambas as estatísticas, cada qual com a sua forma de explicar a variação, forem iguais. Assim, no modelo padrão baseado em neutralismo, a estatística D é igual a zero, sugerindo que outros processos evolutivos, tais como estrutura populacional e seleção, podem não explicar os padrões de variação de certa amostra de sequências (WAKELEY, 2009). O parâmetro θ corresponde à taxa de mutação ao longo da história da amostra, que será abordada na seção 2.2.4 (Modelos Evolutivos) a seguir.

TAJIMA (1989) sugeriu que a estatística D poderia se aproximar de uma distribuição do tipo beta. Embora D não esteja distribuído exatamente de acordo com o modelo beta, valores de D acima de 2 e abaixo de -2 são estatisticamente significativos a 5% de nível de significância (WAKELEY, 2009). Isso significa que valores de D iguais a 2, -2 e entre ambos ($-2 \leq D \leq 2$) não são estatisticamente significativos e diferentes de zero a 5%, estando, portanto, de acordo com o modelo padrão baseado em neutralismo (HEIN ET AL., 2004).

Considerando todo o genoma, HEIN ET AL. (2004) argumentam que valores negativos de D (ramos externos longos) sugerem crescimento populacional ou severo *bottleneck*, enquanto que valores positivos (ramos internos longos) sugerem subdivisão geográfica ou leve *bottleneck*. Em contrapartida, considerando regiões específicas do genoma, valores negativos de D indicam seleção direcional e valores positivos indicam seleção balanceadora ou recente mistura populacional.

As estatísticas D^* e F^* (Fu and Li, 1993) podem ser representadas pelas seguintes equações:

$$D^* = \frac{\frac{S}{a_1} - \frac{n-1}{n}\eta_1}{\sqrt{\widehat{\text{Var}}[\frac{S}{a_1} - \frac{n-1}{n}\eta_1]}} \quad (2.15)$$

$$F^* = \frac{\frac{\pi}{a_1} - \frac{n-1}{n}\eta_1}{\sqrt{\widehat{\text{Var}}[\frac{\pi}{a_1} - \frac{n-1}{n}\eta_1]}} \quad (2.16)$$

em que η_1 é o número de *singletons* ou sítios polimórficos nos quais uma de suas bases nucleotídicas possui apenas uma cópia na amostra de sequências. Ao invés de S e π estarem presentes em uma única medida de neutralidade, como é o caso da estatística D , ambas as estatísticas são consideradas separadamente na abordagem proposta por Fu and Li (1993). Dessa forma, os desvios do modelo padrão referentes às estatísticas D^* e F^* são, respectivamente, baseados nas diferenças entre S e π com quantidades proporcionais ao número de *singletons* presentes na amostra. Portanto, D^* e F^* podem fornecer valores diferentes e, de certo modo, complementares ao valor da estatística D em estudos de genética de populações (WAKELEY, 2009).

No entanto, a estatística D parece ser mais poderosa para detectar desvios do modelo padrão baseado em neutralismo, devendo ser preferencialmente utilizada. SIMONSEN ET AL. (1995), ao utilizarem as estatísticas D , D^* e F^* para casos em que populações sofreram seleção, redução do tamanho populacional e subdivisão, verificaram que apenas a estatística D mostrou valores significativos e diferentes de zero. Contudo, esses valores foram observados somente quando processos evolutivos estiveram ocorrendo por muitas gerações, embora a detecção de tais processos tenha ocorrido considerando um passado relativamente recente. Portanto, é possível que a estatística D apresente valores significativos e diferentes de zero em situações mais drásticas de desvios do modelo padrão baseado em neutralismo.

Embora as estatísticas D , D^* e F^* apresentem comportamento diferentes, as respostas aos desvios do modelo padrão neutro são as mesmas: os valores se tornam negativos quando há excesso de polimorfismos de baixa e alta frequências e deficiência de polimorfismos de média frequência (WAKELEY, 2009). A diferença é que apenas polimorfismos com frequências baixas e altas ao extremo contribuem para valores negativos das estatísticas D^* e F^* , enquanto que isso não acontece para D (WAKELEY, 2009). Para esta última, a contribuição dos sítios polimórficos pode variar de positiva ou negativa simplesmente pelo tamanho da amostra.

Outro importante teste de neutralidade foi desenvolvido por RAMOS-ONSINS and ROZAS (2002). Denominado de estatística R_2 , este teste foi proposto para detectar desvios do modelo baseado em neutralismo especificamente causados por crescimento populacional. Os resultados mostraram que R_2 foi superior para detectar desvios em amostras pequenas e na presença de recombinação. Recentemente, RAMÍREZ-SORIANO ET AL. (2008), ao estudarem diferentes testes na presença de eventos demográficos e recombinação, concluíram que as estatísticas D e R_2 forneceram os resultados mais confiáveis, sendo, portanto, recomendadas.

Embora os testes de neutralidade possam sugerir desvios do modelo padrão baseado em neutralismo, evidências pouco conclusivas são normalmente fornecidas acerca do processo evolutivo (Fu and Li, 1999; WAKELEY, 2009). Assim, modelos evolutivos específicos podem ser investigados com base na teoria da coalescência, utilizando-se abordagens mais precisas na exploração da variação genética e na estimação dos parâmetros de interesse. Veremos a seguir alguns desses modelos, que também serão investigados no presente trabalho.

2.2.4 Modelos Evolutivos

Como foi dito, vários modelos genético-populacionais podem ser assumidos e investigados com base na teoria da coalescência (HUDSON, 1990; PRITCHARD and PRZEWORSKI, 2001; MARJORAM and TAVARÉ, 2006; NORDBORG, 2007). Esses modelos consideram forças evolutivas importantes, tais como mutação (constante e variável), fluxo gênico, recombinação e seleção (FU and LI, 1999; ROSENBERG and NORDBORG, 2002). Em geral, esses modelos apresentam grande complexidade inferencial e demandam enorme tempo computacional, dificultando o uso da abordagem por parte dos pesquisadores. Contudo, essa complexidade dependerá da quantidade de dados disponíveis, da abordagem e do número de parâmetros a ser estimado, que, muitas vezes, é proporcional à complexidade do modelo evolutivo.

Iniciaremos a seguir com o modelo de mutação, que é considerado o mais simples. Alguns conceitos importantes serão inicialmente apresentados e, na sequência, o processo de inferência será mostrado. Depois disso, uma visão global dos modelos de fluxo gênico e recombinação será apresentada, visto que ambos serão aqui investigados. Tais investigações foram escolhidas com o propósito de fornecer informações relevantes para estudos futuros de mapeamento associativo em sorgo, que é a motivação do presente estudo.

2.2.4.1 Mutação

A extensão mais simples do modelo proposto por KINGMAN (1982a,b,c) é aquela que considera apenas mutação (FU and LI, 1999; NORDBORG and TAVARÉ, 2002; WAKELEY, 2009). Para este caso, genealogias históricas são normalmente amostradas e eventos de mutação são simultaneamente incorporados por entre seus ramos, modelando-se genealogia e mutação de maneira independente (HUDSON, 1990). Esta estratégia é totalmente válida pelo fato de que eventos de mutação são seletivamente neutros em populações de FISHER (1930) e WRIGHT (1931), não interferindo na estrutura das genealogias (WAKELEY, 2009). Se os estados alélicos influenciam diretamente a reprodução, não é possível distinguir entre sequências idênticas por descendência ou idênticas por estado (NORDBORG, 2007).

Neste modelo, o objetivo é estimar o parâmetro designado por θ , que corresponde ao número esperado de mutações entre duas sequências quaisquer até o instante em que ambas se coalescem no passado (HEIN ET AL., 2004). θ pode ser obtido a partir dos estimadores $\hat{\theta} = 4N_e\mu$, para o caso de espécies diploides, ou $\hat{\theta} = 2N_e\mu$, para o caso de espécies haploides (ou diploides quando os indivíduos são homozigotos para todos os locos).

Há métodos que estimam θ como sendo a taxa de mutação por sítio (KUHNER, 2006) e métodos que o estimam como sendo a taxa de mutação por região cromossômica, considerando os sítios polimórficos (WATTERSON, 1975). Neste contexto, é fundamental saber aquilo que de fato está sendo estimado, pois, no primeiro o caso, μ é a taxa de mutação por geração, indivíduo e para cada sítio, enquanto que, no segundo caso, μ é a taxa de mutação por geração, indivíduo e região cromossômica (FU and LI, 1999; NORDBORG, 2007; NORDBORG and TAVARÉ, 2002; HEIN ET AL., 2004; McVEAN and CARDIN, 2005; WAKELEY, 2009). De qualquer forma, normalmente o interesse está no parâmetro θ , pois é muito difícil distinguir N_e e μ a partir dos dados sem informações históricas prévias.

O parâmetro θ pode ser estimado com base em diferentes métodos. Dentre eles, tem-se: (i) estimadores clássicos baseados no método dos momentos (WATTERSON, 1975; TAJIMA, 1983; FELSENSTEIN, 1992); (ii) estimadores baseados em máxima verossimilhança (GRIFFITHS and TAVARÉ, 1994a,b,c, 1995; KUHNER ET AL., 1995; FU, 1998); e (iii) estimadores baseados em inferência bayesiana (BEERLI, 2006; KUHNER and SMITH, 2007; KUHNER, 2009). O

estimador clássico proposto por WATTERSON (1975), comumente designado por θ_W , tem sido amplamente utilizado em estudos evolutivos até os dias de hoje. No entanto, apesar de θ_W ser facilmente obtido a partir de um conjunto de sequências, trata-se de um estimador que capitaliza apenas parte da variação dos dados. Isso é reflexo do método dos momentos, que, indiscutivelmente, possui propriedades inferiores quando comparado ao método da máxima verossimilhança (MYUNG, 2003). Assim, é desejável obter estimativas mais precisas e confiáveis para θ , e os estimadores baseados em máxima verossimilhança e inferência bayesiana podem trazer resultados bastante satisfatórios.

As principais abordagens de máxima verossimilhança para estimação de θ foram desenvolvidas por dois grupos: (i) Robert Griffiths e colaboradores (GRIFFITHS and TAVARÉ, 1994a,b,c, 1995); e (ii) Mary Kuhner e colaboradores (KUHNER ET AL., 1995; FELSENSTEIN ET AL., 1999). Na primeira abordagem, a estimação de θ é baseada na amostragem de importância (IS, do inglês *Importance Sampling*), com o uso do procedimento de Monte Carlo para simulação de genealogias por meio da integração computacional. Na segunda abordagem, a estimação de θ também é baseada na amostragem de importância, mas com o uso do procedimento de Monte Carlo via Cadeias de Markov (MCMC, do inglês *Markov Chain Monte Carlo*) para simulação de genealogias por meio da integração computacional. Segundo WAKELEY (2009) e HEY and NIELSEN (2007), Kuhner e colaboradores foram os precursores no uso do MCMC para simulação de genealogias e estimação do parâmetro θ no contexto de coalescência, mostrando que o uso de verossimilhança é possível para investigar modelos de mutação complexos.

Ambas as abordagens utilizam princípios semelhantes para a realização das inferências. Pelo fato dos conceitos coalescentes serem baseados em cadeias de Markov, a primeira abordagem também pode ser considerada como uma forma de MCMC. Contudo, talvez a grande diferença esteja na forma como a simulação é sequencialmente realizada a partir de uma distribuição de probabilidades. Enquanto a primeira abordagem utiliza Monte Carlo para simular genealogias *independentes*, a segunda utiliza Monte Carlo para simular genealogias *correlacionadas* entre si (FELSENSTEIN ET AL., 1999; WAKELEY, 2009). Assim, espera-se que a segunda forneça resultados mais precisos, visto que há tendência de que genealogias altamente prováveis sejam continuamente simuladas para explicar a variação dos dados.

Para compreender o processo de inferência da segunda abordagem (KUHNER ET AL., 1995; FELSENSTEIN ET AL., 1999), considere novamente a amostra de n sequências da Figura 2.1. É possível modelar o passado evolutivo dessas sequências com base na teoria da coalescência, utilizando-se os modelos das equações 2.1 e 2.2. Levando em consideração o tamanho efetivo populacional, é possível reescrever a equação 2.3 de forma a representar explicitamente N_e (FELSENSTEIN ET AL., 1999):

$$P(G' | N_e) = \prod_{i=2}^n \frac{2}{4N_e} \exp\left(\frac{-i(i-1)}{4N_e} u_i \right) \quad (2.17)$$

sendo G' o conjunto de genealogias possíveis na história, i o número de sequências ou linhagens do intervalo de coalescência, e u_i o comprimento do intervalo contendo i sequências.

A equação 2.17 corresponde à distribuição a priori do modelo de coalescência proposto inicialmente por KINGMAN (1982a,b,c). A partir deste e de um modelo evolutivo que considera mutações ao longo da história (FELSENSTEIN, 1981), é possível obter a distribuição de probabilidade dos dados D dado a genealogia, representada por $P(D | G', \mu)$. Combinando ambas as probabilidades, $P(G' | N_e)$ e $P(D | G', \mu)$, é possível obter a distribuição a posteriori dos dados representada pela equação 2.18 (FELSENSTEIN ET AL., 1999):

$$P(D | N_e, \mu) = \int_{G'} P(G' | N_e) P(D | G', \mu) \quad (2.18)$$

com a qual deve-se integrar para todas as genealogias possíveis (FELSENSTEIN, 1988).

As probabilidades $P(G' | N_e)$ e $P(D | G', \mu)$ podem ser facilmente calculadas. Em princípio, a probabilidade $P(D | N_e, \mu)$ também pode ser facilmente calculada pelo produto entre $P(G' | N_e)$ e $P(D | G', \mu)$, considerando coalescência e mutação como histórias independentes. No entanto, é complicado, ou totalmente inviável, obter $P(D | N_e, \mu)$ porque, normalmente, enorme quantidade de genealogias pode ser obtida para a história da amostra, visto que o passado evolutivo é totalmente desconhecido. Esta quantidade é representada pela integral $\int_{G'}$ na equação 2.18, sendo computacionalmente inviável integrar para todas as possíveis histórias com o objetivo de encontrar aquela que melhor explica a variação dos dados. É exatamente por conta disso que a segunda abordagem foi proposta no contexto do MCMC, para integrar sobre genealogias importantes e *correlacionadas*, o que é computacionalmente viável.

No entanto, a integração proposta inicialmente por KUHNER ET AL. (1995) não é realizada diretamente a partir da probabilidade $P(D | N_e, \mu)$. Como já foi dito, é muito difícil distinguir N_e e μ a partir dos dados, de modo que o foco permanece na combinação de ambos, ou seja, no parâmetro θ (KUHNER ET AL., 1995). Assim, alterando a escala de tempo para considerar θ , é possível reescrever a distribuição a priori do modelo de coalescência (equação 2.17) da seguinte forma (FELSENSTEIN ET AL., 1999):

$$P(G | \theta) = \prod_{i=2}^n \frac{2}{\theta} \exp\left(-\frac{i(i-1)}{\theta} v_i\right) \quad (2.19)$$

sendo $G = \mu G'$ e $v_i = \mu u_i$. G corresponde ao conjunto das genealogias G' com escala modificada. v_i corresponde ao comprimento do intervalo u_i com i sequências também com escala modificada (FELSENSTEIN ET AL., 1999).

Da mesma forma, é possível escrever a distribuição a posteriori (equação 2.18) em função de θ pela seguinte equação (FELSENSTEIN ET AL., 1999; STEPHENS, 2007):

$$P(D | \theta) = \int_G P(G | \theta) P(D | G, \theta) \quad (2.20)$$

sendo $P(D | G, \theta)$ a distribuição dos dados a partir dos modelos de coalescência e mutação com escala modificada. Note que o objetivo é estimar o parâmetro θ a partir dos dados, sendo que a verossimilhança, $L(\theta)$, é proporcional à probabilidade $P(D | \theta)$.

A simulação de genealogias importantes proposta por KUHNER ET AL. (1995) é realizada a partir da distribuição escalada $P(D | \theta)$, e as razões para o uso de tal estratégia são claramente descritas em FELSENSTEIN ET AL. (1999). No entanto, para melhor direcionar o processo em função da grande quantidade de histórias (\int_G), KUHNER ET AL. (1995) propuseram a atribuição de um valor inicial para θ . A partir deste valor inicial (θ_0), o processo de simulação via MCMC é realizado com base na distribuição *corrigida* de $P(D | \theta)$. Isso significa que o MCMC visitará um espaço menor de histórias possíveis para as sequências dos dados, condicionado ao valor de θ_0 inicialmente atribuído e, portanto, à distribuição de $P(G | \theta_0)$. Claramente, o valor de θ mais provável estará condicionado ao valor de θ_0 e poderá ser insatisfatório se este não estiver próximo ao seu valor real (KUHNER ET AL., 1995, 1998; KUHNER, 2006, 2009; WAKELEY, 2009). Então, outros valores iniciais para θ_0 , bem como repetições e maior quantidade de iterações poderão ser importantes para a obtenção de resultados mais precisos.

Com base no conjunto de genealogias simuladas via MCMC, outros possíveis valores de θ podem ser investigados. É possível calcular a razão entre as probabilidades de se obter cada genealogia simulada dado o valor investigado de θ e o valor inicialmente atribuído para θ_0 , e obter uma média a partir da soma de todas as razões das genealogias para obtenção de uma superfície de razão de verossimilhanças. Assim, tem-se que:

$$\frac{L(\theta)}{L(\theta_0)} \approx \frac{1}{r} \sum_{i=1}^r \frac{P(D | G_i)P(G_i | \theta)}{P(D | G_i)P(G_i | \theta_0)} \quad (2.21)$$

$$\approx \frac{1}{r} \sum_{i=1}^r \frac{P(G_i | \theta)}{P(G_i | \theta_0)} \quad (2.22)$$

em que r é o número de genealogias simuladas via MCMC. Para cada valor de θ investigado, certo valor de $\frac{L(\theta)}{L(\theta_0)}$ será obtido, de modo que, ao final, um gráfico de razão de verossimilhanças e diferentes valores θ poderá ser construído (KUHNER ET AL., 1995).

Embora a distribuição a posteriori (equação 2.20) apresente a complexidade do modelo de coalescência com mutação, pode ser difícil, apenas analisando a mesma, compreender de que forma tal complexidade está estruturada no referido modelo. Com base nisso, WAKELEY (2009), ao invés de representar a quantidade enorme de genealogias possíveis por \int_G , desmembrou esta integração de modo a mostrar que algumas variáveis estão envolvidas com um processo discreto e outras com um processo contínuo no modelo de coalescência:

$$\begin{aligned} L(\theta) &= \int_G P(G | \theta)P(D | G, \theta) \\ &= \sum_{\Lambda} \int_T P(D | \Lambda, T, \theta)P(\Lambda)f(T)d(T) \\ &= \sum_{\Lambda} \sum_{U \in \Omega_{\Lambda}} \int_T P(D | U, \Lambda)P(U | \Lambda, T, \theta)P(\Lambda)f(T)d(T) \\ &= \sum_{\Lambda} \sum_{U \in \Omega_{\Lambda}} \int_T P(D | U, \Lambda)P(U, \Lambda | \theta) \end{aligned}$$

sendo G correspondente às possíveis genealogias envolvendo as sequências da amostra. G leva em consideração todas as possíveis configurações Λ de genealogias, todos os eventos de mutação U ($U \in \Omega_{\Lambda}$) na história da amostra, e todos os possíveis tempos de coalescência T (T_2, T_3, \dots, T_n), proporcionais aos eventos de coalescência. Dessa forma, G , que pode ser designada conjuntamente por $G = (\Lambda, U, T)$, é o resultado de diferentes processos de Poisson (Λ e U) e distribuições exponenciais independentes (T) caminhando-se do presente em direção ao passado (WAKELEY, 2009; KUHNER, 2009).

2.2.4.2 Estrutura Populacional e Fluxo Gênico

Embora grande parte dos estudos envolvendo populações estruturadas se refira ao fluxo gênico como sendo a taxa de migração, acredita-se que a primeira designação é mais apropriada. Isso porque eventos de migração podem ser meramente geográficos, de forma que a contribuição genética e evolutiva dos migrantes não seja necessariamente garantida. Portanto, neste trabalho, o termo fluxo gênico será sempre utilizado para se referir ao presente modelo.

Tradicionalmente, o fluxo gênico entre populações tem sido estimado com base na estatística F_{ST} (NIELSEN and WAKELEY, 2001; PEARSE and CRANDALL, 2004) ou suas extensões (NATH and GRIFFITHS, 1996). Embora

avanços tenham sido alcançados no entendimento das relações entre populações com o uso desta medida, trata-se de uma estatística limitada (KUHNER, 2009; SOUSA and HEY, 2013). Isso é justificável basicamente por quatro motivos: (i) um mesmo valor da estatística F_{ST} pode ser obtido para cenários bastante diferentes do ponto de vista evolutivo (SOUSA and HEY, 2013); (ii) um único valor de F_{ST} é estimado considerando que todas as populações possuem o mesmo tamanho efetivo e que há um número infinitamente grande de populações; (iii) o fluxo gênico entre diferentes populações é inferido de forma simétrica, ou seja, tendo-se sempre o mesmo valor (BEERLI and FELSENSTEIN, 2001); e, por fim, (iv) a estatística F_{ST} é baseada no método dos momentos, que claramente não capitaliza toda a variação. Neste contexto, acredita-se que os métodos de fluxo gênico baseados em coalescência podem ser promissores para um maior entendimento das relações evolutivas envolvendo populações.

O modelo de coalescência proposto por KINGMAN (1982a,b,c) pode ser estendido para considerar fluxo gênico entre populações (HUDSON, 1990; BEERLI and FELSENSTEIN, 1999, 2001). Caminhando-se do presente em direção ao passado, além dos eventos de coalescência que deverão ser inferidos para cada população, eventos de migração com fluxo gênico também poderão ser modelados entre essas populações (BEERLI and FELSENSTEIN, 2001). Assim, modelos de coalescência com fluxo gênico tendem a ser mais complexos que modelos de coalescência padrões, que apenas consideram mutação.

Em um modelo de coalescência com fluxo gênico, o interesse está em estimar os parâmetros M ou M . M corresponde à razão $\frac{m}{\mu}$, sendo m a taxa de fluxo gênico por geração e μ a taxa de mutação por geração e sítio. Em contrapartida, M pode corresponder à $4N_e m$ (diploides) ou $2N_e m$ (haploides ou diploides homocigotos). O produto de M pela taxa de mutação θ , que é conjuntamente estimada neste modelo, resulta no valor de M (BEERLI, 1998; BEERLI and FELSENSTEIN, 1999, 2001). Então, M e M estão relacionados e correspondem às taxas de fluxo gênico ou ao número mínimo de eventos migratórios que acontecem entre populações de forma a resultar em contribuição genética. A diferença é que M depende diretamente da mutação, enquanto que M é independente deste processo (BEERLI, 1998).

Diversas abordagens têm sido propostas para investigar o fluxo gênico entre populações via modelos de coalescência (TAKAHATA, 1988; TAKAHATA and SLATKIN, 1990; NOTOHARA, 1990; NATH and GRIFFITHS, 1993, 1996; BAHLO and GRIFFITHS, 2000; BEERLI and FELSENSTEIN, 1999, 2001; NIELSEN and WAKELEY, 2001; WILSON and RANNALA, 2003; HEY and NIELSEN, 2004, 2007; HEY, 2010). Tais abordagens forneceram avanços importantes no entendimento dos conceitos de coalescência envolvendo populações estruturadas. Contudo, apenas algumas delas parecem ser muito utilizadas até os dias de hoje, o que é de se esperar devido ao desenvolvimento de modelos mais realistas para o entendimento dos processos evolutivos. Duas dessas abordagens mais utilizadas foram propostas por NIELSEN and WAKELEY (2001) e BEERLI and FELSENSTEIN (2001), que serão apresentadas a seguir.

NIELSEN and WAKELEY (2001) desenvolveram métodos de máxima verossimilhança e inferência bayesiana para estimar seis parâmetros populacionais, a saber: (i) tempo de divergência entre certa população ancestral e duas populações recentes geradas a partir da primeira; (ii) taxa de mutação para cada uma das três populações consideradas; e (iii) taxas de fluxos assimétricas entre as duas populações recentes. Os autores propuseram o uso do MCMC para simular genealogias com elevadas probabilidades e *correlacionadas* entre si. No entanto, tal abordagem foi proposta apenas para estimar o fluxo gênico entre duas populações e sem considerar o uso simultâneo de múltiplos locos, o que pode ser pouco real para muitos casos.

Em contrapartida, BEERLI and FELSENSTEIN (2001) utilizaram os mesmos métodos de inferência para

estimar taxas de mutação e fluxo gênico envolvendo n populações, considerando um único loco ou múltiplos locos simultaneamente. Com a presença de n populações em uma amostra, um total de n^2 parâmetros será estimado, correspondentes às n taxas de mutação para cada população e $n(n - 1)$ taxas de fluxo gênico entre populações. Assim, o modelo leva em consideração a existência de diferentes taxas de mutação para as populações e taxas de fluxo gênico assimétricas entre as mesmas, o que pode ser um cenário real. Este método apresenta inferência similar com a abordagem proposta por BAHLO and GRIFFITHS (2000).

Atualmente, cenários mais complexos têm sido considerados para a estimação de fluxo gênico entre populações. BEERLI and PALCZEWSKI (2010) propuseram a investigação de diferentes modelos com fluxo gênico e a comparação dos mesmos com base em suas verossimilhanças, de forma a verificar aquele que melhor explica o passado evolutivo envolvendo populações. PALCZEWSKI and BEERLI (2013) desenvolveram método aproximado para estimação de altas taxas de fluxo gênico em situações em que populações trocaram muitos alelos entre si. Além disso, taxas de fluxo gênico têm sido investigadas em modelos que consideram, pelo menos em princípio, outros processos, tais como recombinação (BECQUET and PRZEWORSKI, 2007; LOHSE *ET AL.*, 2011; MAILUND *ET AL.*, 2012) e seleção (SOUSA *ET AL.*, 2013).

2.2.4.3 Recombinação

Outra força evolutiva importante que pode ser modelada via teoria da coalescência é a recombinação (HUDSON, 1983; GRIFFITHS and MARJORAM, 1996). Ao contrário da mutação, a recombinação normalmente não deixa vestígios de ocorrência a partir de uma visualização inicial de um conjunto de sequências amostradas (FU and LI, 1999; NORDBORG and TAVARÉ, 2002). Assim, modelos de coalescência com recombinação tendem a ser mais complexos, algo que se torna ainda mais desafiador e intensivo pelo fato de que nenhuma configuração de genealogia pode ser trivialmente modelada diante deste cenário (HEIN *ET AL.*, 2004). Isso acontece porque um evento de recombinação é modelado de forma que uma única sequência mais recente é “quebrada” em posição aleatória para formar duas sequências mais antigas (HEIN *ET AL.*, 2004; McVEAN and CARDIN, 2005; NORDBORG, 2007), sendo o contrário de um evento de coalescência. Então, o processo evolutivo deste modelo é melhor representado por um gráfico e não por uma genealogia, tendo-se, portanto, tanto eventos de coalescência quanto eventos de recombinação (HEIN *ET AL.*, 2004; NORDBORG, 2007).

O primeiro modelo de coalescência com recombinação foi proposto por HUDSON (1983). Por contemplar eventos de coalescência e recombinação modelados aleatoriamente no passado, GRIFFITHS and MARJORAM (1996) propuseram a nomenclatura “gráfico de recombinação ancestral” (ARG, do inglês *Ancestral Recombination Graph*) anos mais tarde. Como alternativa ao modelo de HUDSON (1983), WU and HEIN (1999) desenvolveram uma forma de simular ARGs percorrendo as sequências amostradas de um lado para outro, caracterizando uma complexa estrutura não-Markoviana (McVEAN and CARDIN, 2005). Embora este método tenha tido sua importância, a primeira abordagem é considerada mais simples e útil (HEIN *ET AL.*, 2004).

Em um modelo de coalescência com recombinação, o interesse está em estimar os parâmetros C ou ρ . C corresponde à razão $\frac{r}{\mu}$, sendo r a taxa de recombinação por geração e para cada intersítio (KUHNER *ET AL.*, 1995). ρ pode ser igual a $4N_e r$ (diploides) ou $2N_e r$ (haploides ou diploides homocigotos). C e ρ correspondem às taxas de recombinação populacional (KUHNER *ET AL.*, 1995), sendo o número esperado de recombinações entre duas sequências quaisquer até o instante em que ocorre um evento de coalescência. É importante dizer que a taxa de

mutação θ também é estimada no presente modelo, com eventos de mutação sendo incorporados nas genealogias após a geração dos ARGs (HUDSON, 1990; McVEAN and CARDIN, 2005; NORDBORG, 2007).

A inferência de recombinação via coalescência é claramente difícil e desafiadora (WALL, 2000; STUMPF and McVEAN, 2003; NORDBORG, 2007). Estimadores baseados em momentos foram propostos para obtenção de ρ a partir das estimativas de θ (HUDSON, 1987; WAKELEY, 1997). Esses estimadores apenas utilizam parte da variação disponível, apresentando viés e elevada variância para as estimativas. Em contrapartida, estimadores baseados em máxima verossimilhança e inferência bayesiana utilizam toda a informação disponível nos dados, embora apresentem restrições por não possuírem nenhuma expressão analítica ou numérica para o cálculo de suas verossimilhanças (McVEAN and CARDIN, 2005). Certamente, isso dificulta o entendimento de suas propriedades (FU and LI, 1999) e a construção de métodos de simulação eficientes para a estimação de tal processo (McVEAN and CARDIN, 2005).

Normalmente, a quantidade de ARGs que podem ser simulados a partir de certa distribuição é enorme. A simulação sobre todos eles, muitos dos quais possuem probabilidades irrisórias para as verossimilhanças (McVEAN and CARDIN, 2005), é ingênua e intratável do ponto de vista computacional. Com base nisso, KUHNER ET AL. (2000) e FEARNHEAD and DONNELLY (2001) desenvolveram métodos de Monte Carlo mais eficientes para a estimação de recombinação. Tais métodos estimam um único valor de recombinação para certa região genômica, com base em abordagem completa. Contudo, ambos podem ser ineficientes para grandes quantidades de dados, algo que se tornou uma realidade no final da década de 90.

Como alternativa, vários métodos de coalescência foram desenvolvidos para estimar recombinação a partir de muitos dados (WALL, 2000; HUDSON, 2001; McVEAN ET AL., 2002; FEARNHEAD and DONNELLY, 2002; LI and STEPHENS, 2003; McVEAN ET AL., 2004; McVEAN and CARDIN, 2005; AUTON and McVEAN, 2007). Em essência, todos esses métodos baseiam-se em alguma estratégia de aproximação na tentativa de reduzir a complexidade da abordagem completa. O mais recente deles (AUTON and McVEAN, 2007) foi proposto para considerar taxas de recombinação variáveis e pontos quentes de recombinação (*hotspots*) entre pares de locos. Tal método é baseado em um tipo específico de MCMC no contexto bayesiano (rjMCMC, do inglês *reversible jump* MCMC) (GREEN, 1995), com o qual ARGs de elevadas probabilidades são simulados a partir de buscas em várias dimensões de probabilidade, aumentando as chances de se obter estimativas mais precisas de recombinação. Acredita-se que a abordagem proposta por AUTON and McVEAN (2007) é bastante promissora, principalmente para grandes quantidades de dados.

Um aspecto interessante é que eventos de recombinação podem ser gerados tanto a partir de trocas entre cromátides homólogas (*crossing-overs*) quanto conversão de genes. Todos os métodos descritos anteriormente são baseados no primeiro caso, o que pode ser pouco realista do ponto de vista biológico. FRISSE ET AL. (2001) e PTAK ET AL. (2004) desenvolveram abordagem para estimar recombinação baseada tanto em *crossing-overs* quanto conversão de genes. De forma mais ampla, WALL (2004) propôs a estimação de recombinação tanto com base em *crossing-overs*, para comparar com abordagens anteriores, quanto com base em ambos, *crossing-overs* e conversão gênica.

2.2.5 Máxima Verossimilhança e Inferência Bayesiana

Trabalhos recentes apresentam comparações entre os métodos de máxima verossimilhança e inferência bayesiana na estimação de parâmetros via coalescência. BEERLI (2006), ao utilizar o método proposto por BEERLI and FELSENSTEIN (2001) e sua versão bayesiana desenvolvida posteriormente, detectou diferenças importantes na es-

timação de fluxo gênico entre populações. Em contrapartida, KUHNER and SMITH (2007) não detectaram diferenças substanciais entre tais métodos na estimação de crescimento populacional e recombinação.

BEERLI (2006) menciona que a abordagem de máxima verossimilhança apresenta dificuldades para cenários complexos, com os quais grande quantidade de parâmetros é normalmente estimada. Considera que tais dificuldades também podem ocorrer quando altas taxas de fluxo gênico e menor variabilidade são esperadas entre populações. Isso porque o processo de simulação via MCMC não se comporta de forma apropriada em função da grande quantidade de eventos, dificultando a busca e a convergência para estimativas que maximizam a verossimilhança. Essas dificuldades devem ter estimulado o desenvolvimento de abordagem aproximada em trabalho recente (PALCZEWSKI and BEERLI, 2013), muito embora a inferência bayesiana possa fornecer resultados confiáveis para situações complexas (BEERLI, 2006).

2.3 Desequilíbrio de Ligação e Teoria da Coalescência

A associação não-aleatória, ou preferencial, entre alelos de diferentes locos em uma população fornece um profundo entendimento sobre a sua história evolutiva (McVEAN, 2002). Apesar dos estudos relacionados ao DL serem de longa data na genética de populações, sua relação com a teoria da coalescência apenas foi claramente apresentada por trabalhos recentes (ZÖLLNER and VON HAESELER, 2000; PRITCHARD and PRZEWORSKI, 2001; NORDBORG and TAVARÉ, 2002; McVEAN, 2002; McVEAN and CARDIN, 2005; McVEAN, 2007a,b).

ZÖLLNER and VON HAESELER (2000) utilizaram os conceitos da teoria da coalescência para estudar o DL em populações de tamanho constante e com crescimento exponencial, a partir de dados simulados de SNPs. PRITCHARD and PRZEWORSKI (2001) também utilizaram dados simulados para investigar DL considerando diferentes cenários demográficos via coalescência, tais como tamanho populacional constante, crescimento exponencial e estrutura de população. PRITCHARD and PRZEWORSKI (2001) e NORDBORG and TAVARÉ (2002) apresentaram as perspectivas do DL via coalescência para estudos de mapeamento associativo, em dois artigos que esclareceram a importância de uma visão genealógica dos padrões de associação para tal objetivo.

Apesar de esses artigos terem mostrado ideias e resultados interessantes, a relação formal do DL com a teoria da coalescência foi recentemente apresentada por McVEAN (2002). Utilizando-se a tradicional medida r^2 do DL (HILL and ROBERTSON, 1968), este autor sugeriu uma interpretação genealógica dos padrões de associação a partir do cálculo do seu valor esperado, que, embora não possua expressão analítica possível, pode ser aproximado por uma razão de esperanças, como mostrado na equação 2.23 (OHTA and KIMURA, 1969a,b, 1971):

$$E[r_{AB}^2] \approx \frac{E[D_{AB}^2]}{E[f_A(1-f_A)f_B(1-f_B)]} \quad (2.23)$$

sendo A e B dois locos bi-alélicos fisicamente ligados em determinado cromossomo.

Desenvolvendo o numerador da equação 2.23, é possível mostrar que:

$$\begin{aligned}
E[D_{AB}^2] &= E[(f_{AB} - f_A f_B)(f_{AB} - f_A f_B)] \\
&= E[f_{AB}^2 - 2f_{AB}f_A f_B + f_A f_B f_A f_B] \\
&= E[f_{AB}^2] - 2E[f_{AB}f_A f_B] + E[f_A f_B f_A f_B] \\
&= F_{A(ij)B(ij)} - 2F_{A(ij)B(ik)} + F_{A(ij)B(kl)} \\
&\quad \vdots \\
\text{Var}[D_{AB}] &= \frac{\text{Cov}[t_{A(ij)}, t_{B(ij)}] - 2\text{Cov}[t_{A(ij)}, t_{B(ik)}] + \text{Cov}[t_{A(ij)}, t_{B(kl)}]}{E[T_A T_B]}
\end{aligned}$$

em que $F_{A(ij)B(ij)}$, $F_{A(ij)B(ik)}$ e $F_{A(ij)B(kl)}$ são os coeficientes de identidade para dados de sequências. Esses coeficientes correspondem às probabilidades de duas sequências (i) i e j serem idênticas para os locos A e B ; (ii) i e j serem idênticas para o loco A , e i e k serem idênticas para o loco B ; e (iii) i e j serem idênticas para o loco A , e k e l serem idênticas para o loco B . Além disso, esses coeficientes podem ser expressos na forma de covariância dos tempos de coalescência envolvendo os locos A e B ($\text{Cov}[t_A, t_B]$), conforme descrito na última expressão da equação. T_A e T_B correspondem aos tempos dos locos A e B para alcançar o MRCA de toda a amostra.

Da mesma forma, desenvolvendo o denominador da equação 2.23, tem-se que:

$$E[f_A(1 - f_A)f_B(1 - f_B)] = \frac{E[t]^2 + \text{Cov}[t_{A(ij)}, t_{B(kl)}]}{E[T_A T_B]}$$

em que $E[t]$ é o tempo de coalescência esperado para um par de cromossomos.

Assim, reescrevendo a equação 2.23 a partir das derivações do seu numerador e denominador, é possível mostrar que:

$$E[r_{AB}^2] \approx \sigma_d^2 = \frac{\text{Cov}[t_{A(ij)}, t_{B(ij)}] - 2\text{Cov}[t_{A(ij)}, t_{B(ik)}] + \text{Cov}[t_{A(ij)}, t_{B(kl)}]}{E[t]^2 + \text{Cov}[t_{A(ij)}, t_{B(kl)}]} \quad (2.24)$$

Para uma amostra de tamanho finito, algumas modificações devem ser incorporadas de forma a considerar que as sequências i , j , k e l podem não ser totalmente distintas. Dessa forma, $E[r_{AB}^2]$, ou simplesmente σ_d^2 (OHTA and KIMURA, 1969a), é mostrado de forma modificada pela equação 2.25 (HUDSON, 1985; McVEAN, 2002):

$$\sigma_d^2 = \frac{[n^2 - 2(n-1)]C_{ij,ij} - 2(n-2)^2C_{ij,ik} + (n-2)(n-3)C_{ij,kl} + nE[t]^2}{n(n-1)E[t]^2 + 2C_{ij,ij} + 4(n-2)C_{ij,ik} + (n-2)(n-3)C_{ij,kl}} \quad (2.25)$$

em que os C 's correspondem às abreviações das covariâncias dos tempos de coalescência entre os locos A e B . É possível notar o tamanho da amostra (n) nos termos da equação 2.25.

McVEAN (2002) mostrou que as covariâncias dos tempos de coalescência são funções diretas de processos evolutivos importantes, tais como recombinação e fluxo gênico. No caso da recombinação, essas ideias

já haviam sido propostas por trabalhos anteriores de coalescência (GRIFFITHS, 1981; HUDSON, 1983; KAPLAN and HUDSON, 1985), mostrando a motivação de se estudar tal processo do ponto de vista genealógico (equações 2.26, 2.27 e 2.28).

$$\text{Cov}[t_{A(ij)}, t_{B(ij)}] = \frac{18 + \rho}{18 + 13\rho + \rho^2} \quad (2.26)$$

$$\text{Cov}[t_{A(ij)}, t_{B(ik)}] = \frac{6}{18 + 13\rho + \rho^2} \quad (2.27)$$

$$\text{Cov}[t_{A(ij)}, t_{B(kl)}] = \frac{4}{18 + 13\rho + \rho^2} \quad (2.28)$$

Utilizando-se tais covariâncias na equação 2.24, é possível expressar σ_d^2 em função de ρ (OHTA and KIMURA, 1971; WEIR and HILL, 1986; WEIR ET AL., 1990; McVEAN, 2002):

$$\sigma_d^2 = \frac{10 + \rho}{22 + 13\rho + \rho^2} = \frac{10 + \rho}{(2 + \rho)(11 + \rho)} \quad (2.29)$$

assumindo que a amostra populacional é finita e muito grande (McVEAN, 2002). A expressão 2.29 foi desenvolvida a partir de uma população que apresenta recombinação, deriva genética e baixa mutação (WEIR ET AL., 1990). No equilíbrio entre recombinação e deriva genética, σ_d^2 , ou $E[r^2]$ como é comumente utilizado, torna-se (SVED, 1968, 1971; WEIR ET AL., 1990; NIELSEN and SLATKIN, 2013):

$$E[r^2] \approx \sigma_d^2 = \frac{1}{1 + \rho} = \frac{1}{1 + 4N_e r} \quad (2.30)$$

também assumindo que a amostra populacional é finita e muito grande (McVEAN, 2002).

É importante dizer que o DL das equações 2.29 e 2.30 não foi inicialmente proposto utilizando os princípios da teoria da coalescência (WAKELEY, 2009). Tais equações foram obtidas a partir dos conceitos da clássica genética de populações. Porém, é possível investigar o DL a partir da recombinação populacional estimada via coalescência.

No caso do modelo de coalescência com fluxo gênico, McVEAN (2002) mostrou que apenas uma das covariâncias dos tempos de coalescência (equação 2.24) é função da taxa de migração com fluxo gênico ($\text{Cov}[t_{A(ij)}, t_{B(ij)}] = \frac{1}{4M^2}$, sendo M igual à $4N_e m$ ou $2N_e m$), sendo que as demais covariâncias ($\text{Cov}[t_{A(ij)}, t_{B(ik)}]$ e $\text{Cov}[t_{A(ij)}, t_{B(kl)}]$) são iguais a zero (McVEAN, 2002).

Até o momento, nenhuma abordagem do DL ainda foi proposta para considerar vários processos simultaneamente. Acredita-se que tal abordagem poderá ser desenvolvida em breve no contexto de coalescência, já que estudos recentes têm sido propostos para modelar processos conjuntamente (LOHSE ET AL., 2011; MAILUND ET AL., 2012).

Estudos do DL baseados na teoria da coalescência podem ser uma estratégia interessante para compreender o contexto histórico-evolutivo de gerações e os seus padrões de associação. Em última análise, tais estudos podem proporcionar um mapeamento genético mais refinado para espécies de interesse (NORDBORG and TAVARÉ, 2002; ROSENBERG and NORDBORG, 2002; ZÖLLNER and PRITCHARD, 2005; McVEAN, 2007a; KIMMEL ET AL., 2008).

2.4 Sorgo

2.4.1 Classificação e Taxonomia

O gênero *Sorghum* é pertencente a tribo Andropogoneae, subfamília Panicoidae, família Poaceae (Gramineae). Dentre as várias classificações propostas para este gênero, as mais compreensíveis e interessantes (SINGH and LOHITHASWA, 2006) foram apresentadas por SNOWDEN (1936), GARBER ET AL. (1950), HARLAN and DE WET (1972a) e DE WET (1978).

SNOWDEN (1936) classificou o gênero *Sorghum* nas categorias de Eusorghum e Para-sorghum, sendo a primeira dividida em Arundinaceae, Spontanea (espécies ou raças selvagens), Sativa (raças cultivadas) e Halepensis. GARBER ET AL. (1950) dividiram o gênero *Sorghum* em seis subgêneros com base em dados citotaxonômicos, a saber: Eusorghum, Chaetosorghum, Heterosorghum, Sorghastrum, Parasorghum e Stiposorghum. HARLAN and DE WET (1972a) elaboraram uma classificação simplificada com base em características morfológicas, com as quais os indivíduos selvagens foram agrupados em seis raças espontâneas (arundinaceum, aethiopicum, virgatum, verticilliflorum, propinquum e shattercane) e os indivíduos cultivados foram agrupados em cinco raças básicas (bicolor, guinea, caudatum, kafir e durra) e dez raças híbridas (guinea-bicolor, caudatum-bicolor, kafir-bicolor, durra-bicolor, guinea-caudatum, guinea-kafir, guinea-durra, kafir-caudatum, durra-caudatum e kafir-durra). DE WET (1978) desenvolveu uma classificação dentro do subgênero Eusorghum (ou sorghum) com a qual reconheceu a existência de três diferentes espécies, a saber: *Sorghum halepense* (L.) Pers, *Sorghum propinquum* (Kunth) Hitchc. e *Sorghum bicolor* (L.) Moench.

O sorgo [*Sorghum bicolor* (L.) Moench] é uma espécie autógama e diploide ($2n = 20$), sendo subdividida em três subespécies: *S. bicolor* ssp. *bicolor*, *S. bicolor* ssp. *drummondii* e *S. bicolor* ssp. *verticilliflorum*. Todos os grupos comerciais de sorgo cultivado, tais como sorgo granífero, forrageiro e vassoura, são classificados dentro da subespécie bicolor (SINGH and LOHITHASWA, 2006; DE WET and HARLAN, 1971). A classificação das cinco raças básicas e das dez raças híbridas intermediárias para o sorgo cultivado é bastante utilizada nos dias de hoje. Detalhes históricos e morfológicos sobre as espécies relatadas, a espécie cultivada *S. bicolor* ssp. *bicolor* e as raças descritas podem ser encontrados em DOGGETT ET AL. (1988).

2.4.2 Centro de Origem e Domesticação

A região da Etiópia, localizada na parte leste da África, é considerada como o centro de origem do sorgo (SINGH and LOHITHASWA, 2006). Esta região possui enorme quantidade de espécies e variedades da raça durra, que está entre as raças básicas cultivadas do sorgo (SINGH and LOHITHASWA, 2006). Segundo SMITH and FREDERIKSEN (2000) e MORRIS ET AL. (2013), a domesticação inicial do sorgo aconteceu no Sudão, África, por volta de 10.000 anos atrás, a partir dos indivíduos selvagens e cultivados provenientes da Etiópia, com difusão para outras regiões da África e países da Ásia entre 8.000 e 1.500 anos atrás. Essas regiões compreendem países como Nigéria, Niger e Mali, na parte oeste da África, Uganda e Kênia, na parte leste da África, e Índia e China, no continente asiático. Da parte oeste da África, o sorgo foi distribuído para os Estados Unidos e outras regiões do mundo por meio do tráfico de escravos há mais de 200 anos (QUINBY, 1974, 1975; ROSENOW ET AL., 1983; ROONEY and SMITH, 2000; SMITH and FREDERIKSEN, 2000), em meados do século 19 (SINGH and LOHITHASWA, 2006). Mesmo antes de 1900, o cultivo do sorgo nos Estados Unidos já apresentava-se de forma bastante extensiva (ROSENOW ET AL., 1983; SINGH

and LOHITHASWA, 2006).

Em função da origem e difusão antigas, a adaptação do sorgo para climas e práticas culturais específicas foi refletida em termos de variações morfológicas e fisiológicas envolvendo as cinco raças básicas de sorgo domesticado (HARLAN and DE WET, 1972b). Essas variações foram o resultado de eventos demográficos importantes, tais como redução do tamanho populacional, estrutura de população, seleção e migração (HAMBLIN ET AL., 2006). Mais especificamente, as raças cultivadas de sorgo foram originadas de processos de seleção disruptiva e domesticação ocorridos na região centro-leste da África, a partir de indivíduos selvagens representantes da espécie ancestral *Sorghum arundinaceum*. A seleção natural em espécies selvagens e a seleção artificial praticada pelo homem, para características relacionadas ao sorgo cultivado, resultaram em populações diversificadas com considerável fluxo gênico entre si, explicando muito das variações raciais e do processo evolutivo da espécie (SINGH and LOHITHASWA, 2006).

2.4.3 Importância e Melhoramento Genético

O sorgo é o quinto cereal mais importante do mundo, permanecendo atrás do trigo, arroz, milho e cevada. Representa grande importância para alimentação humana em diversos países da África, Ásia e América Central e para alimentação animal nos Estados Unidos, Austrália e América do Sul. O sorgo granífero possui enorme capacidade de tolerar condições ambientais adversas, produzindo de forma satisfatória em locais de grande estresse hídrico. Normalmente, esta condição ambiental é desfavorável para o crescimento e a produção de outras culturas, fazendo com que o sorgo mostre o seu potencial e adquira o seu espaço enquanto cultura (SANTOS ET AL., 2005; SINGH and LOHITHASWA, 2006).

Atualmente, o maior produtor mundial de sorgo é os Estados Unidos, com produção anual superior a 10 milhões de toneladas (15% da produção mundial) na safra 2014/2015 (USDA, 2015). O Brasil é o décimo maior produtor mundial, com produção anual de 2 milhões de toneladas, permanecendo atrás de México, Nigéria, Índia, Sudão, Argentina, Etiópia, China e Austrália. A importância do sorgo é recente à agricultura brasileira, sendo explorado de forma intensiva a partir de 1970, com a utilização de quatro tipos cultivados de sorgo: granífero, forrageiro, sacarino e vassoura. Os dois primeiros se destacam tanto pela área cultivada quanto pelo nível tecnológico adotado (SANTOS ET AL., 2005).

Em geral, o melhoramento de sorgo tem sido realizado para promover adaptação às condições específicas dos locais de interesse. Isso reflete na obtenção de cultivares com níveis aceitáveis de produtividade, estabilidade e qualidade dos grãos de acordo com o ambiente desejado. Além disso, essas cultivares precisam mostrar resistência ou tolerância a certos fatores bióticos, tais como doenças e pragas, e fatores abióticos, tais como altas temperaturas e estresse hídrico (condições de seca). Ainda, é preciso mostrar qualidade nutricional para atender às exigências das alimentações humana e animal (SINGH and LOHITHASWA, 2006).

Nos Estados Unidos, o melhoramento do sorgo passou por uma série de fases importantes, a saber (KLEIN ET AL., 2008): (i) introdução de um número pequeno de cultivares oriundos da África entre 1874 e 1908; (ii) seleção de plantas de maturação precoce e baixa estatura (insensíveis ao fotoperíodo) a partir de populações heterogêneas entre 1904 e 1936; (iii) melhoramento de plantas visando a colheita mecanizada entre 1930 e 1940; (iv) implementação da produção de sementes híbridas de 1946 até os dias de hoje; e (v) conversão de acessos tropicais exóticos de maturação tardia e alta estatura (sensíveis ao fotoperíodo) em acessos temperados de ma-

turação precoce e baixa estatura a partir de 1963 (STEPHENS *ET AL.*, 1967; QUINBY, 1974; ROSENOW *ET AL.*, 1983; KLEIN *ET AL.*, 2008), de forma a introduzir variabilidade exótica de genótipos tropicais produtivos no germoplasma temperado (STEPHENS *ET AL.*, 1967; CASA *ET AL.*, 2008). O uso deste germoplasma é utilizado até os dias de hoje no melhoramento, e cada uma dessas etapas exerceu grande influência na produção de sorgo nas regiões semi-áridas do mundo (KLEIN *ET AL.*, 2008).

A conversão de acessos tropicais em acessos temperados, conhecida como SCP (*Sorghum Conversion Program*) (STEPHENS *ET AL.*, 1967), foi uma estratégia importante estabelecida pelo USDA (*United States Department of Agriculture*) em colaboração com a *Texas University* (CASA *ET AL.*, 2008). Em linhas gerais, vários cruzamentos envolvendo cultivares tropicais (não-adaptadas) e temperadas (adaptadas) foram realizados, e os indivíduos das populações obtidas foram submetidos a vários ciclos de retrocruzamentos com as cultivares tropicais (ROSENOW *ET AL.*, 1997). Dessa forma, grande parte do genoma (~ 90%) dos indivíduos gerados correspondeu ao genoma da cultivar tropical, com a presença de alelos desejáveis de cultivares temperadas. Em torno de 840 – 850 linhagens convertidas ou parcialmente convertidas foram desenvolvidas para constituir um novo germoplasma e, assim, uma nova fonte de variabilidade para o melhoramento do sorgo (CASA *ET AL.*, 2008; KLEIN *ET AL.*, 2008).

O melhoramento genético de sorgo no Brasil inclui a obtenção de híbridos visando obter resistência a fatores bióticos (doenças e pragas) e abióticos (tolerância ao estresse hídrico e à toxidez por alumínio, aumento da eficiência de uso de fósforo no solo e acamamento), além de características particulares tanto do local como do tipo de sorgo cultivado (SANTOS *ET AL.*, 2005). A obtenção de híbridos resistentes a fatores bióticos e abióticos também tem sido objetivo do melhoramento de países da África e da Ásia (MORRIS *ET AL.*, 2013).

2.4.4 Alguns Genes de Importância

Nos últimos 20 anos, diversos estudos de mapeamento genético têm sido realizados em sorgo (MACE and JORDAN, 2010). Esses estudos tiveram como objetivo a investigação de características quantitativas importantes, tais como altura (LIN *ET AL.*, 1995; PEREIRA and LEE, 1995; BROWN *ET AL.*, 2008), maturidade (CRASTA *ET AL.*, 1999; CHANTEREAU *ET AL.*, 2001; HART *ET AL.*, 2001; FELTUS *ET AL.*, 2006), tolerância ao estresse hídrico (TUINSTRAL *ET AL.*, 1996; CRASTA *ET AL.*, 1999; SUBUDHI *ET AL.*, 2000; TAO *ET AL.*, 2000; XU *ET AL.*, 2000; KEBEDE *ET AL.*, 2001; HARRIS *ET AL.*, 2007), tolerância ao alumínio (MAGALHAES *ET AL.*, 2004), restauração da fertilidade (KLEIN *ET AL.*, 2001; JORDAN *ET AL.*, 2010) e resistência a doenças e pragas (TAO *ET AL.*, 2003). Com base nesses e outros trabalhos, MACE and JORDAN (2010) reuniram um grupo de 35 genes de maior efeito com o objetivo de posicioná-los conjuntamente em um mapa genético consenso. Esses autores tiveram como proposta a integração das informações disponíveis para promover direcionamentos importantes em programas de melhoramento. Dentre os genes descritos e estudados, MACE and JORDAN (2010) consideraram alguns importantes relacionados à maturidade e altura de plantas, que são características correlacionadas e amplamente relatadas em sorgo.

Devido à esta importância, os genes de maturidade e altura de plantas foram escolhidos para a realização do presente estudo, e portanto serão descritos a seguir.

2.4.4.1 Exemplos Relevantes Para o Estudo

Os genes relacionados à maturidade e altura de plantas são conhecidos de longa data em sorgo (QUINBY, 1974, 1975). Estudos genéticos clássicos dessas características identificaram quatro locos relacionados à matu-

ridade (*Ma1*, *Ma2*, *Ma3* e *Ma4*) e quatro locos relacionados à altura de plantas (*Dw1*, *Dw2*, *Dw3* e *Dw4*), algo que certamente gerou importante progresso e estimulou o desenvolvimento de trabalhos subsequentes. Os três primeiros locos para maturidade foram identificados por QUINBY and KARPER (1945), e o quarto loco foi identificado por QUINBY (1966). Todos os locos relacionados à altura de plantas foram inicialmente descritos por QUINBY and KARPER (1954), embora um estudo mais antigo já havia abordado com propriedade sobre o controle genético dessa característica (KARPER, 1932).

Para ambas as características, a interação de dominância entre os alelos resulta na obtenção de plantas com maturação tardia e estatura elevada, como é o caso das cultivares tropicais exóticas. Nos Estados Unidos, pesquisadores e produtores selecionaram, ao longo de anos, plantas com maturação precoce e baixa estatura (alelos mutantes), como forma de obter cultivares adaptadas ao clima temperado e à colheita mecanizada (KLEIN ET AL., 2008).

Em geral, os alelos mutantes dos locos de maturidade e altura de plantas surgiram após a entrada do sorgo nos Estados Unidos. Para ambas as características, três locos foram inicialmente identificados como tendo alelos recessivos, a saber: (i) *ma1*, *ma2* e *ma3*; e (ii) *dw1*, *dw2* e *dw3*. Em geral, as cultivares convertidas para maturação precoce e baixa estatura apresentam as três mutações para os genes *Ma* e *Dw*. O alelo recessivo *ma4* foi identificado posteriormente com a chegada de uma variedade específica da África, que mostrou adaptação temperada com apenas uma mutação. O alelo recessivo *dw4* já estava presente em certos indivíduos mesmo antes da sua entrada nos Estados Unidos, e o alelo dominante *Dw4* foi verificado em apenas um indivíduo que apresentava características típicas de cultivares tropicais (QUINBY, 1975).

Dentre os locos de maturidade, *Ma1* é o que tem maior impacto no período de florescimento e crescimento (KLEIN ET AL., 2008; MURPHY ET AL., 2011). Grande proporção da variação no florescimento foi explicada por este loco em estudo anterior (LIN ET AL., 1995), mostrando que, possivelmente, o surgimento do alelo recessivo *ma1* foi o principal fator na conversão de cultivares tropicais em temperadas. Em estudo posterior, dois locos adicionais para maturidade (*Ma5* e *Ma6*) foram identificados por ROONEY and AYDIN (1999). A motivação surgiu pelo fato de que outros locos poderiam estar envolvidos no controle do fotoperíodo além dos quatro locos iniciais, algo que já havia sido sugerido por trabalhos anteriores (QUINBY ET AL., 1973; QUINBY, 1974; LIN ET AL., 1995).

Trabalhos recentes mostraram que os genes *Ma1*, *Ma3* e *Ma6* estão situados nos cromossomos 6 (40,30 Mb), 1 (60,92 Mb) e 6 (0,70 Mb) do genoma do sorgo, respectivamente. Os dois primeiros foram anteriormente clonados e possuem localização bem conhecida (LIN ET AL., 1995; CHILDS ET AL., 1997; MURPHY ET AL., 2011), enquanto que a localização do último ainda pode ser incerta (MURPHY ET AL., 2014). Um estudo recente relatou que o gene *Ma5* está localizado no cromossomo 1 (MURPHY ET AL., 2014), muito embora não haja evidências suficientes para suportar tal afirmação. MACE and JORDAN (2010) relatam que o gene *Ma5* está localizado no cromossomo 2, tendo-se, portanto, evidências pouco conclusivas. Também não há evidências claras das posições genômicas dos genes *Ma2* e *Ma4*, embora haja informações de que o último está localizado no cromossomo 10 (MACE and JORDAN, 2010).

Um estudo recente mostrou que o loco *Ma1* é responsável pela codificação da proteína regulatória de falsa resposta 37 (PRR37, do inglês *Pseudoresponse Regulator Protein 37*) (MURPHY ET AL., 2011). Esses autores avaliaram um germoplasma de sorgo granífero contendo alelos recessivos para *Ma1*, e detectaram variações significativas de PRR37. Essas variações podem estar envolvidas com a expressão e inibição de mecanismos

relacionados ao fotoperíodo em dias curtos e longos, o que é decisivo para a adaptação da espécie em temperaturas diferentes.

O loco *Ma3* foi investigado e clonado em estudo realizado por CHILDS *ET AL.* (1997). Esses autores genotiparam dois marcadores moleculares fisicamente ligados a *Ma3*, em população de linhagens recombinantes, e detectaram fortes associações com o gene *PHYB*. Este gene é responsável pela codificação do fitocromo 123-kD em plantas, e os indivíduos mutantes (*ma3*) apresentaram fenótipos parecidos com plantas que não produziram fitocromo 123-kD. Esses resultados foram decisivos para a conclusão de que o loco *Ma3* é um tipo de *PHYB*.

Com relação aos locos relacionados à altura de plantas, acredita-se que todos eles estejam em cromossomos diferentes no genoma do sorgo (KLEIN *ET AL.*, 2008). Os locos *Dw1*, *Dw2* e *Dw3* estão localizados nos cromossomos 9 (57,20 Mb), 6 (42,20 Mb) e 7 (58,61 Mb), respectivamente (LIN *ET AL.*, 1995; BROWN *ET AL.*, 2008). A localização do loco *Dw4* ainda parece ser incerta, mas acredita-se que ele esteja situado no cromossomo 6 (6,60 Mb), com base em estudo recente de mapeamento associativo (MORRIS *ET AL.*, 2013). O loco *Dw2* encontra-se fisicamente ligado ao loco *Ma1* no cromossomo 6 (KLEIN *ET AL.*, 2008), e o loco *Dw3* é o único que foi clonado em trabalho anterior (MULTANI *ET AL.*, 2003).

2.4.5 Marcadores Moleculares e Genotipagem Por Sequenciamento

Os marcadores moleculares têm sido amplamente utilizados em sorgo para diversas finalidades. Dentre elas, tem-se: (i) construção de mapas genéticos e mapeamento de QTLs para características quantitativas importantes - MACE and JORDAN (2010) apresentam uma ampla revisão sobre esses estudos; (ii) análise do desequilíbrio de ligação (HAMBLIN *ET AL.*, 2004, 2005, 2007; BOUCHET *ET AL.*, 2012; MORRIS *ET AL.*, 2013) e mapeamento associativo (CASA *ET AL.*, 2008; SUKUMARAN *ET AL.*, 2012; MORRIS *ET AL.*, 2013); (iii) análise da diversidade genética e do fluxo gênico entre populações (MORDEN *ET AL.*, 1989, 1990; ALDRICH *ET AL.*, 1992; ALDRICH and DOEBLEY, 1992; AGRAMA and TUINSTRAN, 2004; CASA *ET AL.*, 2005; FOLKERTSMA *ET AL.*, 2005; DEU *ET AL.*, 2006, 2008; DE ALENCAR FIGUEIREDO *ET AL.*, 2008; MUTEGI *ET AL.*, 2011; SAGNARD *ET AL.*, 2011; BOUCHET *ET AL.*, 2012); e (iv) estimativa de taxas de recombinação (HAMBLIN *ET AL.*, 2005; MORRIS *ET AL.*, 2013) e seleção (HAMBLIN *ET AL.*, 2005, 2006; CASA *ET AL.*, 2006; BOUCHET *ET AL.*, 2012). A realização desses estudos contemplou o uso de diversos tipos de marcadores moleculares, tais como aloenzima, RFLP (*Restricted Fragment Length Polymorphism*), RAPD (*Random Amplified Polymorphic DNA*), SSR (*Simple Sequence Repeat*), DAiT (*Diversity Arrays Technology*) e SNP (*Single Nucleotide Polymorphism*).

Nos últimos anos, a descoberta de marcadores moleculares e a genotipagem de centenas de indivíduos têm sido realizadas a partir das novas plataformas de sequenciamento para muitos estudos (VARSHNEY *ET AL.*, 2009; DAVEY *ET AL.*, 2011). Com base em um único procedimento de sequenciamento (DAVEY *ET AL.*, 2011), essas plataformas possibilitam a geração de grandes quantidades de dados a uma velocidade bastante superior em comparação com as metodologias tradicionais (VARSHNEY *ET AL.*, 2009), mostrando eficiência, rapidez e custo reduzido. Dentre às diferentes técnicas desenvolvidas a partir desta abordagem, tem-se a genotipagem por sequenciamento (GBS, do inglês *Genotyping-By-Sequencing*), proposta por ELSHIRE *ET AL.* (2011). Trata-se de uma técnica simples, rápida, eficiente, específica e de baixo custo, que se baseia na redução da complexidade do genoma através do uso de enzimas de restrição específicas e sensíveis à metilação (ELSHIRE *ET AL.*, 2011). Muitos fragmentos são gerados a partir dos cortes realizados por essas enzimas, possibilitando o sequenciamento e a geração de dezenas de milhões

de *reads* (sequências curtas de DNA), nas quais polimorfismos podem ser detectados posteriormente (SANSALONI, 2012).

Em função dessas e outras vantagens, a técnica de GBS tem sido utilizada em estudos de diversas espécies de plantas e animais. No entanto, o GBS apresenta desvantagens importantes, tais como a distribuição não uniforme das leituras realizadas e a alta quantidade de dados perdidos (BEISSINGER ET AL., 2013). Segundo DAVEY ET AL. (2011), POLAND ET AL. (2012) e POLAND and RIFE (2012), esta última é decorrência da baixa cobertura do genoma exercida pela técnica, que pode ocorrer tanto por questões de amostragem quanto por motivos biológicos. Contudo, diversas estratégias de imputação têm sido desenvolvidas e podem ser utilizadas para recuperar informações perdidas (MARCHINI and HOWIE, 2010). No caso do mapeamento e seleção genômica, isso não chega a ser um problema, mas, para alguns tipos de estudo, a baixa cobertura pode ter implicações na qualidade dos resultados.

3 MATERIAL E MÉTODOS

Como foi dito na Introdução (seção 1), o objetivo deste trabalho foi analisar, com base na teoria da coalescência, a estrutura populacional e o desequilíbrio de ligação em um painel mundial de acessos de sorgo. Para alcançar tal objetivo, análises de mutação, migração com fluxo gênico e recombinação foram realizadas para diferentes regiões genômicas e populações que constituem este painel.

Nesta seção, descreveremos inicialmente o painel de acessos e os dados moleculares gerados a partir do procedimento de genotipagem, incluindo a estratégia adotada para a seleção das regiões genômicas de interesse. Na sequência, mostraremos as estatísticas de diversidade e os testes de neutralidade utilizados, que são métodos clássicos inicialmente investigados para se ter uma visão global da variação dos dados. Por fim, apresentaremos os modelos de coalescência, incluindo os modelos de estrutura populacional (migração com fluxo gênico) e recombinação (abordagens completa e aproximada), e o desequilíbrio de ligação via coalescência.

3.1 Material Biológico e Dados Moleculares

3.1.1 Painel de Acessos

O painel de acessos do presente estudo é constituído por 356 indivíduos de sorgo oriundos de diferentes continentes e regiões geográficas do mundo, tais como África, Ásia, América do Norte, América Central e América do Sul. Países historicamente importantes para a origem e o melhoramento desta espécie estão presentes, tais como Etiópia, Índia, Nigéria, África do Sul, Sudão, Uganda e Estados Unidos. Outros países importantes também estão presentes, mas com menor quantidade de acessos. É o caso de Kênia, Chad, Congo, Egito, Mali, Niger, Senegal, Burkina Faso, Tanzânia, entre outros. Este painel é parte de uma coleção elaborada para estudos de mapeamento associativo e diversidade genética, que foram recentemente publicados por CASA *ET AL.* (2008) e MORRIS *ET AL.* (2013). As localidades dos acessos foram coletadas no *Germplasm Resources Information Network (GRIN)*, pertencente ao *Agricultural Research Service (ARS)* do *United States Department of Agriculture (USDA)*.

Do total de 356 indivíduos, 114 correspondem a linhagens tropicais exóticas e 242 correspondem a linhagens convertidas. Em geral, as linhagens tropicais exóticas (mais antigas) foram selecionadas devido à sua importância histórica para os programas de melhoramento por serem fontes de resistência ou tolerância a estresses bióticos e abióticos, além da diversidade genética e geográfica presentes nos seus centros de origem (África). As linhagens convertidas (mais recentes) foram selecionadas de acordo com a insensibilidade ao fotoperíodo (SCP) e o desempenho como cultivares em diferentes localidades do mundo (CASA *ET AL.*, 2008).

3.1.2 Genotipagem Por Sequenciamento e Imputação de Dados

A genotipagem de marcadores por sequenciamento foi realizada para os 356 indivíduos do painel de acessos. As etapas foram realizadas pelo *Institute for Genomic Diversity* pertencente à *Cornell University* (Ithaca, NY, EUA), seguindo o método proposto por ELSHIRE *ET AL.* (2011) em plataforma HiSeq 2000 (Illumina[®] Inc., San Diego, CA, EUA). Inicialmente, o DNA genômico dos indivíduos foi digerido pela enzima *ApeKI*, que reconhece sítio contendo cinco bases (uma delas é degenerada) e possui sensibilidade à metilação. A partir desta digestão, bibliotecas de 96 e 384-plex foram construídas para os indivíduos, e as sequências obtidas via Illumina foram mapeadas contra o genoma de referência do sorgo (PATERSON *ET AL.*, 2009). Tal mapeamento foi realizado com base

no método proposto por LI ET AL. (2009).

Para a descoberta dos SNPs, utilizou-se o *pipeline* TASSEL-GBS implementado na versão 3.0 do software TASSEL. Este *pipeline* foi descrito detalhadamente por GLAUBITZ ET AL. (2014) com base em versão mais recente do software. *Tags* de GBS contendo 64 pares de bases (pb) foram alinhadas contra o genoma de referência do sorgo (linhagem *BTx623*), e aquelas que estiveram presentes, no mínimo, dez vezes (10×) nas bibliotecas foram mantidas para o mapeamento. Com base no arquivo *HapMap* gerado ao final do processo, um total de 265,487 SNPs foram descobertos para os 356 indivíduos (e as repetições utilizadas como controle). A imputação de dados perdidos foi realizada via software NPUTE (ROBERTS ET AL., 2007), com janelas de SNPs variando de 5 a 150 pb para cada cromossomo. As janelas que mostraram as maiores acurácias foram selecionadas e utilizadas para imputação.

Os dados moleculares gerados com os procedimentos anteriores foram gentilmente disponibilizados pelo Prof. Dr. Geoff Morris, do *Department of Agronomy, Crop Genetics & Genomics* pertencente à *Kansas State University* (Manhattan, KS, EUA).

3.1.3 Análises Prévias no Painel

Para reduzir a complexidade das análises de coalescência devido à grande quantidade de locais, decidiu-se utilizar apenas os países Etiópia (ETI), Índia (IND), Nigéria (NIG), África do Sul (AFS), Sudão (SUD), Uganda (UGA) e Estados Unidos (EUA). No entanto, para não reduzir o painel inicialmente proposto, análises de componentes principais foram previamente realizadas na tentativa de agrupar os acessos dos demais países naqueles aqui mais representativos. Tais análises foram realizadas utilizando-se os dados totais de SNPs, obtidos via GBS e NPUTE, e o painel de acessos com os 356 indivíduos. Essas análises foram realizadas no pacote *pcaMethods* (STACKLIES ET AL., 2007), disponível no software R (R CORE TEAM, 2015).

3.1.4 Seleção das Regiões Genômicas

Sete locos genômicos importantes para a história de melhoramento do sorgo foram selecionados para investigação no presente estudo. Estes locos correspondem a genes relacionados com altura de plantas (*Dw1*, *Dw2*, *Dw3* e *Dw4*) e maturidade (*Ma1*, *Ma3* e *Ma6*). Vários estudos têm evidenciado suas posições estimadas no genoma do sorgo (LIN ET AL., 1995; CHILDS ET AL., 1997; MULTANI ET AL., 2003; BROWN ET AL., 2008; MURPHY ET AL., 2011; MORRIS ET AL., 2013; MURPHY ET AL., 2014). No entanto, as posições que delimitam o início e o fim desses genes são desconhecidas ou pouco precisas até o momento, sendo necessário a adoção de estratégias alternativas para a investigação.

Assim, a partir das posições estimadas desses locos no genoma, o DL baseado na medida r^2 (HILL and ROBERTSON, 1968) foi estimado entre todos os pares de SNPs a cada 5.000 pb, caminhando-se simultaneamente tanto para esquerda quanto para direita. Uma vez que SNPs com elevado DL tendem a mostrar uma mesma genealogia evolutiva, acredita-se que SNPs adjacentes possam fornecer evidências importantes para os genes de interesse. O número 5.000 foi utilizado como referência ao tamanho médio de um gene em sorgo, com base no comprimento de sequências preditas relacionadas à maturidade desta espécie (BHOSALE ET AL., 2012).

Variante adjacentes às posições estimadas foram selecionadas enquanto houve DL entre os SNPs mais extremos, delimitando assim as regiões de interesse. A presença de DL foi detectada com base na correção de Bonferroni (WEIR ET AL., 1990) para múltiplos testes, a partir das probabilidades obtidas pelo teste exato de Fisher

(FISHER, 1935). Essas análises foram realizadas utilizando-se o software R (R CORE TEAM, 2015).

3.2 Diversidade Genética e Testes de Neutralidade

Estatísticas de diversidade genética e neutralidade foram estimadas para cada região genômica de interesse, a saber: (i) número de sítios polimórficos (S) considerando todos os acessos (WATTERSON, 1975); (ii) diferença nucleotídica entre pares de sequências (π) envolvendo todos os acessos (TAJIMA, 1983); (iii) taxa de mutação populacional (θ_w) para cada sítio e todos os sítios (WATTERSON, 1975); e (iv) estatísticas D de Tajima (TAJIMA, 1989) e R_2 de Ramos-Onsins-Rozas (RAMOS-ONSINS and ROZAS, 2002) como testes de neutralidade para investigar desvios do modelo padrão baseado em neutralismo (WAKELEY, 2009). A estatística R_2 foi especificamente utilizada para detectar desvios causados pela alteração do tamanho efetivo populacional (expansão ou redução). Essas análises foram realizadas utilizando-se o pacote *pegas* (PARADIS, 2010), disponível no software R (R CORE TEAM, 2015).

3.3 Modelos de Coalescência

3.3.1 Estrutura Populacional e Fluxo Gênico

Análises de coalescência com fluxo gênico foram realizadas para sete populações de diferentes localidades utilizando-se a versão bayesiana da abordagem proposta por BEERLI and FELSENSTEIN (2001). Tais localidades foram apresentadas na seção 3.1.3. Esta abordagem baseia-se no modelo de isolamento geográfico com fluxo gênico (WRIGHT, 1931), que é fundamentado nos pressupostos do modelo de FISHER (1930) e WRIGHT (1931) estruturado. No total, 49 parâmetros populacionais foram simultaneamente estimados, correspondendo às sete taxas de mutação das populações e às 42 taxas de fluxo gênico assimétricas envolvendo todos os possíveis pares de populações. Uma matriz de migração com fluxo gênico (NORDBORG, 2007) \mathcal{P} foi obtida contendo todos os parâmetros populacionais estimados, conforme representação mostrada a seguir (BEERLI and FELSENSTEIN, 2001). Tal representação foi modificada em relação à matriz original de BEERLI and FELSENSTEIN (2001), que não manteve os índices (subscritos) das taxas de fluxo gênico de acordo com os conceitos da teoria de matrizes, dificultando o entendimento.

$$\mathcal{P} = \begin{pmatrix} \theta_1 & \mathcal{M}_{12} & \mathcal{M}_{13} & \dots & \mathcal{M}_{1n} \\ \mathcal{M}_{21} & \theta_2 & \mathcal{M}_{23} & \dots & \mathcal{M}_{2n} \\ \dots & \dots & \dots & \dots & \dots \\ \mathcal{M}_{n1} & \mathcal{M}_{n2} & \dots & \mathcal{M}_{n,n-1} & \theta_n \end{pmatrix}$$

Na matriz \mathcal{P} , $\theta_1, \theta_2, \dots, \theta_n$ as taxas de mutação correspondentes às n populações e $\mathcal{M}_{12}, \mathcal{M}_{13}, \dots, \mathcal{M}_{n,n-1}$ as taxas de fluxo gênico assimétricas envolvendo pares de populações. O parâmetro \mathcal{M}_{ij} é equivalente a razão $\frac{m_{ij}}{\mu}$, sendo m_{ij} a taxa de fluxo gênico por geração que ocorre da população i para a população j e μ a taxa de mutação por sítio e geração. O parâmetro θ_i é igual a $2N_e^{(i)}\mu$, sendo $N_e^{(i)}$ o tamanho efetivo da população i em um modelo de FISHER (1930) e WRIGHT (1931). Oportunamente, o valor de θ do presente modelo será referido como θ_{LM} .

Para a estimação da matriz de parâmetros \mathcal{P} , genealogias ancestrais importantes e *correlacionadas* foram simuladas com base no procedimento de MCMC. A partir de uma distribuição a priori logarítmica para θ e

linear para \mathcal{M} (KUHNER and SMITH, 2007), e do modelo evolutivo F84 apresentado por FELSENSTEIN and CHURCHILL (1996), o MCMC foi realizado por meio do algoritmo Metropolis-Hastings (METROPOLIS ET AL., 1953; HASTINGS, 1970). Assim, certa genealogia simulada no estado G_t foi obtida a partir de pequenas alterações nos ramos ancestrais da genealogia simulada no estado G_{t-1} (BEERLI and FELSENSTEIN, 2001), com base nos coalescentes condicionais (FELSENSTEIN ET AL., 1999). A transição de uma genealogia para outra ao longo do MCMC, indicando progresso no processo de simulação, ocorreu sempre que a probabilidade da genealogia G_t foi superior à probabilidade da genealogia G_{t-1} .

Inicialmente, 20 cadeias curtas (iniciais), cada qual com 20.000 genealogias, foram percorridas para obter estimativas iniciais e direcionadas dos parâmetros. Com base nessas estimativas, 10 cadeias longas (finais), cada qual com 400.000 genealogias, foram percorridas, e as estimativas da última cadeia longa foram utilizadas para as inferências populacionais. Adicionalmente, este mesmo procedimento foi realizado considerando quatro cadeias independentes e em paralelo (MCMCMC ou MC³, do inglês *Metropolis-Coupled Markov Chain Monte Carlo*) (GEYER, 1991; KUHNER ET AL., 2000) com diferentes temperaturas (1,0; 1,2; 1,5 e 4,0) (MORRELL ET AL., 2003), sendo que, ao final de cada cadeia curta e longa, permutações foram realizadas com o objetivo de aumentar a busca e eventualmente obter estimativas mais precisas dos parâmetros. Para todos os processos, as primeiras 40.000 e 400.000 genealogias das cadeias curtas e longas, respectivamente, foram descartadas como período de aquecimento (*burnin*), correspondendo a 10% do total. O procedimento foi repetido por 10 vezes para o caso em que não houve MC³ e 5 vezes para o caso em que houve MC³, sendo realizado para cada região gênômica de interesse. Em função do tempo computacional, o MC³ não foi realizado para a análise de todas as regiões simultaneamente. Um valor máximo de 10.000 eventos de migração foi utilizado para cada processo realizado.

A função de verossimilhança do presente modelo pode ser representada pela seguinte equação (BEERLI and FELSENSTEIN, 2001):

$$L(\mathcal{P}) = \int_{\mathcal{G}} P(\mathcal{G} | \mathcal{P})P(D | \mathcal{G}) \quad (3.1)$$

sendo \mathcal{G} o espaço de genealogias simuladas e \mathcal{P} a matriz que contém os parâmetros correspondentes a θ e \mathcal{M} . Os parâmetros de \mathcal{P} foram estimados para cada genealogia simulada, sendo estratégia utilizada pela abordagem bayesiana no presente contexto.

Na expressão de verossimilhança da equação 3.1, $P(\mathcal{G} | \mathcal{P})$ foi calculado com base no seguinte modelo (BEERLI and FELSENSTEIN, 2001):

$$P(\mathcal{G} | \mathcal{P}) = \prod_{j=1}^T \left[\exp \left(-u_j \left[\sum_{i=1}^s \frac{k_{ij}(k_{ij}-1)}{\theta_i} + k_{ij} \sum_{z \neq i}^s \mathcal{M}_{zi} \right] \right) \times \left(\delta_j \mathcal{M}_{w_j v_j} + (1 - \delta_j) \frac{2}{\theta_v} \right) \right]$$

O termo exponencial é a probabilidade de que em cada tempo j e em seu comprimento de genealogia u_j não aconteçam nem eventos de fluxo gênico nem eventos de coalescência. O segundo termo da expressão indica que esses eventos podem supostamente acontecer. Tais eventos correspondem ao fluxo gênico (\mathcal{M}), que ocorre da população w_j para a população v_j , e ao coalescente, que ocorre para a população v_j . A variável indicadora δ_j apresenta o valor de 1 quando um evento de fluxo gênico acontece no tempo j ou o valor de 0 para o caso contrário. k_{ij} corresponde ao número de sequências presentes na população i durante o intervalo de tempo j .

As análises do modelo de coalescência com fluxo gênico foram realizadas utilizando a versão Linux do pacote LAMARC (KUHNER, 2006, 2009).

3.3.2 Recombinação

3.3.2.1 Abordagem Completa

Estimativas de recombinação constante via coalescência foram obtidas para cada população e região genômica de interesse. Tais estimativas consideram a ocorrência de uma única taxa de recombinação histórica entre todos os SNPs presentes em certa região. Para tanto, utilizou-se a abordagem de máxima verossimilhança proposta por KUHNER *ET AL.* (2000), que utiliza o procedimento do MCMC e o algoritmo Metropolis-Hastings (METROPOLIS *ET AL.*, 1953; HASTINGS, 1970) para simular gráficos de recombinação ancestral ou ARGs (GRIFFITHS and MARJORAM, 1996), contemplando tanto eventos de coalescência quanto eventos de recombinação (HEIN *ET AL.*, 2004). A partir de dois valores iniciais para a mutação ($\theta_{0(1)}$: valor de θ_W e $\theta_{0(2)}$: valor de 0,01 atribuído pelo software utilizado) e um único valor inicial para a recombinação (ρ_0 : valor de 0,01 atribuído pelo software utilizado), ARGs foram simulados utilizando a estratégia dos coalescentes condicionais, baseada no modelo F84 proposto por FELSENSTEIN and CHURCHILL (1996). De maneira independente, eventos de mutação foram incorporados nos ARGs amostrados (HUDSON, 1990; McVEAN and CARDIN, 2005), de modo que, ao final, estimativas de mutação e recombinação (únicas e constantes) foram obtidas a partir do conjunto de ARGs simulados. Essas estimativas representam os valores de mutação e recombinação que podem ser pontos de máxima verossimilhança. Desconsiderou-se a ocorrência de múltiplas recombinações e interferência entre intervalos adjacentes, bem como processos de conversão gênica.

No caso da mutação, o valor de $\theta_{Lp} = 2N_e\mu$ por sítio foi estimado para cada região e população. No caso da recombinação, os valores de $\frac{r}{\mu_L}$ e ρ_L entre dois possíveis sítios também foram estimados para cada região e população. Utilizou-se a mesma quantidade de iterações estabelecida para o modelo com fluxo gênico, pois, embora a quantidade de parâmetros seja aqui menor, trata-se de modelo complexo e difícil de ser investigado (HEIN *ET AL.*, 2004; McVEAN and CARDIN, 2005; NORDBOG, 2007). Então, 20 cadeias curtas de 20.000 ARGs e 10 cadeias longas de 400.000 ARGs foram percorridas, com o descarte (*burnin*) das primeiras 40.000 e 400.000 cadeias, respectivamente. Quatro cadeias independentes e em paralelo, com diferentes temperaturas (1,0; 1,2; 1,5 e 4,0) (MORRELL *ET AL.*, 2003), também foram adicionalmente utilizadas para aumentar o espaço de ARGs simulados.

Apenas duas repetições foram realizadas considerando ausência e presença de MC³. Isso foi decidido basicamente por três motivos: (i) a abordagem de máxima verossimilhança é normalmente mais intensiva que a bayesiana do ponto de vista computacional; (ii) dois valores iniciais para θ foram considerados; e (iii) as estimativas foram obtidas para cada região e população. Um valor máximo de 1.000 recombinações foi utilizado para cada processo.

As análises de coalescência para taxa de recombinação constante também foram realizadas na versão Linux do pacote LAMARC (KUHNER, 2006, 2009).

3.3.2.2 Abordagem Aproximada

Estimativas de recombinação variável para cada região genômica e população também foram obtidas com base em abordagem aproximada de verossimilhança. Ao contrário de se estimar um único valor constante

de recombinação para toda a região, como foi realizado pela abordagem completa (seção 3.3.2.1), estimativas variáveis envolvendo pares de SNPs foram calculadas utilizando-se o método de AUTON and McVEAN (2007). Tal método é baseado no rjMCMC (GREEN, 1995), que consiste em simular ARGs de alta probabilidade com base em diferentes dimensões, aumentando o espaço da busca.

O método de AUTON and McVEAN (2007) foi utilizado por meio de três etapas. Na primeira delas, a taxa de mutação populacional foi estimada a partir de uma versão modificada do método de WATTERSON (1975), que é baseada no modelo de sítios finitos. Tal modelo considera que determinado sítio polimórfico pode ter ocorrido a partir de vários eventos históricos de mutação. Então, com base em uma amostra de n sequências de comprimento L e região contendo S sítios polimórficos, foi possível obter θ_W^* para cada sítio e geração (equação 3.2):

$$\theta_W^* = \left(\sum_{i=1}^{n-1} \frac{1}{i} \right)^{-1} \ln \left(\frac{L}{L-S} \right) \quad (3.2)$$

Na segunda etapa, a região em estudo foi dividida em conjuntos envolvendo pares de SNPs, reduzindo a complexidade das análises. Na terceira etapa, o método de verossimilhança completa proposto por FEARNEHEAD and DONNELLY (2001) foi utilizado para cada conjunto de SNPs, possibilitando a simulação de muitas histórias e, assim, a estimação das taxas de mutação e recombinação. Considerou-se a possibilidade de mutação reversa e utilizou-se um intervalo de diferentes taxas de recombinação constante. Ao final, um conjunto de verossimilhanças completas e pré-calculadas foi obtido, possibilitando o cálculo da verossimilhança composta ou aproximada (L_C) de acordo com a equação 3.3 (AUTON and McVEAN, 2007):

$$L_C(\rho) = s^{-1} \sqrt{\prod_{i,j} P(D_{ij} | \rho_{r(ij)})} \quad (3.3)$$

em que a raiz com índice dependente de S corresponde à uma correção para que a verossimilhança seja confiavelmente calculada. Esta correção foi sugerida por AUTON and McVEAN (2007) por meio da análise de vários cenários com dados simulados. D_{ij} e $\rho_{r(ij)}$ correspondem aos dados e às taxas de recombinação populacional entre os SNPs i e j , respectivamente. $\rho_{r(ij)}$ é igual à $2N_e r_{r(ij)}$ (neste caso, os indivíduos são homocigotos para todos os locos), sendo $r_{r(ij)}$ a taxa de recombinação por indivíduo e geração entre os SNPs i e j .

Os parâmetros de simulação utilizados com o método de FEARNEHEAD and DONNELLY (2001) foram: (i) taxa de mutação inicial calculada para cada sítio com base no estimador de WATTERSON (1975); (ii) taxas de recombinação constante variando de 0 a 100 para cada conjunto de SNPs; e (iii) posições investigadas a cada 1 Kb de distância. Os parâmetros utilizados com o método de AUTON and McVEAN (2007) foram: (i) taxas de mutação e recombinação obtidas a partir do método de FEARNEHEAD and DONNELLY (2001); (ii) 10.000.000 de iterações percorridas; (iii) penalizações variando de 0 a 50, conforme sugerido pelos desenvolvedores, e 10 repetições para cada uma das penalizações; (iv) reamostragem de ARGs a cada 100 processos de iteração; e (v) 1.000.000 de cadeias descartadas em cada processo referentes ao período de aquecimento (*burnin*). As estimativas médias de recombinação variável foram selecionadas a partir da penalização que apresentou a menor variação ao longo das 10 repetições.

As análises foram realizadas na versão 2.2 do pacote LDhat utilizando os programas `complete` (FEARNEHEAD and DONNELLY, 2001) e `rhomap` (AUTON and McVEAN, 2007).

3.4 Desequilíbrio de Ligação Via Coalescência

Até o momento, nenhuma medida do DL foi proposta considerando múltiplos processos estimados via coalescência. Nesse sentido, as estimativas de recombinação variável obtidas via `rhomap` (AUTON and McVEAN, 2007) foram utilizadas para o cálculo do valor esperado do DL ($E[r^2]$), considerando cada população e região genômica de interesse. A abordagem do $E[r^2]$ escolhida (equação 2.30) foi aquela em que a população encontra-se em equilíbrio entre recombinação e deriva genética (SVED, 1968, 1971; WEIR ET AL., 1990). Isso porque a outra abordagem disponível do $E[r^2]$ (equação 2.29) considera a presença adicional de (baixa) mutação, que, em essência, é importante apenas por revelar algo do processo evolutivo (NORDBORG and TAVARÉ, 2002). Então, com as atenções voltadas ao mapeamento associativo, onde a recombinação é o fator de importância para o DL (PRITCHARD and PRZEWORSKI, 2001; NORDBORG and TAVARÉ, 2002), decidiu-se utilizar o $E[r^2]$ da equação 2.30.

Assim, $E[r^2]$, ou especificamente $E[r_{\text{rhomap}}^2]$, pôde ser calculado pela seguinte equação:

$$E[r_{\text{rhomap}}^2] \approx \frac{1}{1 + \rho_{r(ij)}} = \frac{1}{1 + 2N_e r_{r(ij)}} \quad (3.4)$$

em que $\rho_{r(ij)}$ e $r_{r(ij)}$ correspondem às taxas de recombinação populacional e por indivíduo entre os SNPs i e j , respectivamente. Note, portanto, que a equação 3.4 é equivalente à equação 2.30.

O DL via medida r^2 (HILL and ROBERTSON, 1968) também foi calculado para os mesmos pares de SNPs, e a correlação de Spearman entre r^2 e $E[r_{\text{rhomap}}^2]$ foi investigada utilizando o software R (R CORE TEAM, 2015). O DL via coalescência também foi calculado a partir da estimativa única e constante de recombinação obtida pelo método de KUHNER ET AL. (2000), obtendo-se $E[r_{\text{LAMARC}}^2]$. Esta quantidade foi adicionalmente calculada como forma de inferir sobre o DL de alta ordem envolvendo múltiplos locos simultaneamente dentro de certa região.

4 RESULTADOS

4.1 Seleção dos Acessos de Sorgo

Do total de 265.487 SNPs gerados com a técnica de GBS, 247.806 (93,34%) foram selecionados e utilizados nas análises prévias de componentes principais. O processo de imputação via NPUTE identificou 17.681 (6,66%) sítios como sendo monomórficos (único nucleotídeo) para o presente painel, os quais, portanto, foram excluídos das análises. As maiores acurácias obtidas estiveram entre 96,70% e 97,38% de acordo com o cromossomo, a partir de janelas variando de 26 a 34 SNPs.

O gráfico das análises de componentes principais é mostrado na Figura 4.1. É possível observar certo agrupamento dos indivíduos pertencentes a populações como ETI, IND, NIG e AFS. Contudo, os dados totais de SNPs não sugeriram a formação de grupos claros de acordo com as populações. Para os EUA, por exemplo, os indivíduos se agruparam com acessos de todas as demais populações, de forma que os componentes principais pouco sugeriram padrões que possibilitassem a inclusão dos indivíduos de populações menos representativas nas sete populações selecionadas. Portanto, 58 indivíduos pertencentes às populações não selecionadas foram excluídos do painel inicial, restando, assim, um subpainel de 298 para as análises.

4.2 Seleção das Regiões Genômicas

De um total de sete locos genômicos inicialmente selecionados, cinco foram mantidos para as análises de coalescência. Apenas 14 e 4 SNPs foram observados para os locos *Dw3* e *Ma6* (respectivamente) enquanto houve DL a partir das suas posições estimadas no genoma (resultados não mostrados). No nosso entendimento, essas quantidades são insuficientes para as inferências, e portanto ambos os locos foram desconsiderados do presente trabalho.

Os padrões do DL envolvendo pares de SNPs adjacentes às posições dos locos *Dw1*, *Dw2*, *Dw4*, *Ma1* e *Ma3* são mostrados na Figura 4.2. Os gráficos referentes a todos os locos são mostrados até o instante em que associações significativas estiveram presentes ou queda pronunciada do DL foi observada, com base na correção de Bonferroni para múltiplos testes.

Para o loco *Dw1*, uma região adjacente de 70 quilo bases (Kb) foi selecionada tanto para direita quanto para esquerda da sua posição estimada (cromossomo 9: 57,2 Mb), totalizando 140 Kb. Nesta região, 89 SNPs foram observados e 900 (22,98%) associações de pares foram estatisticamente significativas quanto ao desequilíbrio de um total de 3.916 associações. Para o loco *Dw2*, uma região de 75 Kb foi selecionada a partir da sua posição estimada (cromossomo 6: 42,2 Mb), totalizando 150 Kb. Ao todo, 85 SNPs foram detectados e 682 (19,10%) associações mostraram estar em DL a partir de 3.570 combinações possíveis envolvendo pares de SNPs. Para o loco *Dw4*, uma região total de 140 Kb foi selecionada ao redor da sua posição estimada (cromossomo 6: 6,6 Mb), contendo 54 SNPs e 653 (45,63%) associações em DL de um total de 1.431 pares de SNPs. Para o importante loco de maturidade *Ma1*, uma região total de 160 Kb foi selecionada adjacientemente à sua posição estimada (cromossomo 6: 40,3 Mb), contendo 39 SNPs e 269 (36,30%) associações preferenciais a partir de 741 associações totais. E, por fim, uma região total de 120 Kb foi selecionada em torno da posição estimada do loco *Ma3* (cromossomo 1: 60,92 Mb), contendo 52 SNPs e 235 (17,72%) associações estatisticamente significativas de um total de 1.326 combinações.

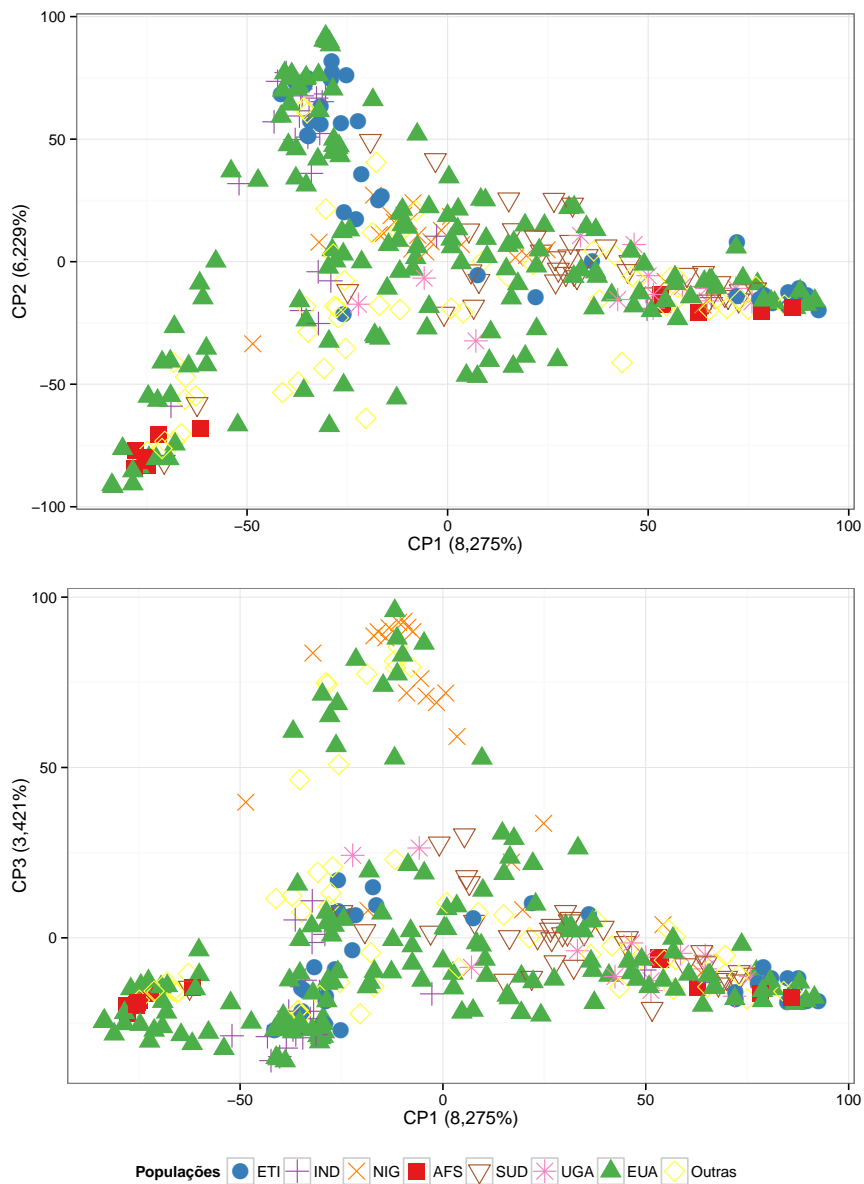


Figura 4.1: Gráfico das análises de componentes principais utilizando 247.806 SNPs e 356 indivíduos do painel mundial de acessos de sorgo. Os três primeiros componentes principais foram utilizados para inferir sobre os padrões de agrupamento. As populações referentes aos países Etiópia, Índia, Nigéria, África do Sul, Sudão, Uganda e Estados Unidos são indicadas na legenda como ETI, IND, NIG, AFS, SUD, UGA e EUA, respectivamente. Os indivíduos oriundos das demais populações estão representados na legenda por Outras

4.3 Estatísticas de Diversidade e Testes de Neutralidade

As estatísticas de diversidade genética e os testes de neutralidade são mostrados na Tabela 4.1. Tais quantidades são apresentadas para se ter uma visão inicial dos padrões de variação. É possível observar variabilidade genética entre as regiões genômicas de interesse, com base na quantidade de sítios polimórficos (S), nas estimativas das taxas de mutação para toda a região e por sítio ($\theta_{W(\text{Região})}$ e $\theta_{W(\text{Sítio})}$) e na diversidade nucleotídica (π). A região *DwI*, localizada no cromossomo 9, foi a que apresentou o maior polimorfismo ($S = 89$) e os maiores valores das taxas de mutação ($\theta_{W(\text{Região})} = 14,189$ e $\theta_{W(\text{Sítio})} = 1,148 \times 10^{-4}$). Em contrapartida, a região que apresentou o menor polimorfismo ($S = 39$) e as menores estimativas de mutação foi a *Mal* ($\theta_{W(\text{Região})} = 6,217$ e

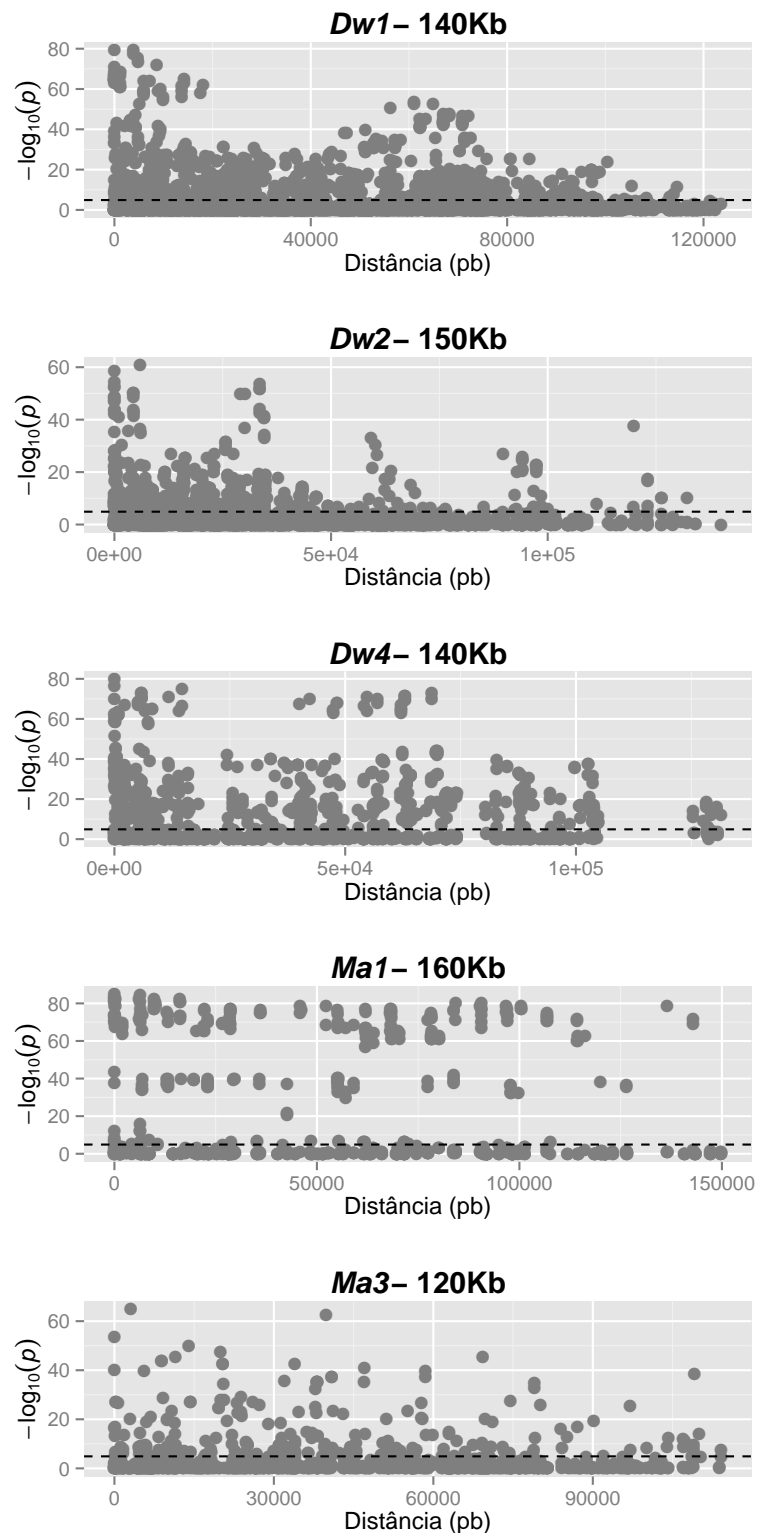


Figura 4.2: Queda do desequilíbrio de ligação, $-\log_{10}(p)$, em função da distância física entre SNPs em pares de bases (pb). A probabilidade obtida a partir do teste exato de Fisher é mostrada em escala logarítmica, e o limiar correspondente ao procedimento de Bonferroni para múltiplos testes é indicado pela linha horizontal pontilhada. Os pares de SNPs estatisticamente significativos (em DL) estão indicados acima da linha do Bonferroni. No título de cada gráfico, são mostradas as extensões em Kb determinadas para cada região genômica contendo o respectivo loco

$\theta_{W(\text{Sítio})} = 0,415 \times 10^{-4}$), situada no cromossomo 6. Um aspecto interessante é que os locos fisicamente ligados *Dw2* e *Mal* apresentaram grandes variações de polimorfismo e das taxas de mutação por região e para cada sítio. Embora o tamanho definido para as regiões possa influenciar nos padrões de variação, *Mal* foi a região que teve a maior quantidade de pares de bases e a menor variação. Esta influência pode ter sido aqui pouco expressiva.

Com base na comparação de pares envolvendo os 298 indivíduos do subpainel utilizado, as maiores diferenças entre nucleotídeos de um mesmo sítio polimórfico foram apresentadas pelas regiões *Mal* ($\pi = 0,293$) e *Dw4* ($\pi = 0,238$). A região *Mal* foi a única que apresentou valor significativo para a estatística *D*, indicando desvios do modelo padrão baseado em neutralismo. O valor positivo de *D* sugeriu uma menor proporção de polimorfismos de baixa e alta frequências para a região *Mal*. Um valor positivo e elevado de *D* também foi observado para a região *Dw4*, embora este não tenha sido estatisticamente significativo. Nenhuma região genômica apresentou valor significativo para a estatística de crescimento populacional R_2 .

Tabela 4.1: Estatísticas de diversidade genética e testes de neutralidade para as cinco regiões genômicas considerando o subpainel de 298 indivíduos de sorgo

Região	Cr.	S	Tamanho (pb)	$\theta_{W(\text{Região})}$	$\theta_{W(\text{Sítio})}$	π	<i>D</i>	R_2
<i>Dw1</i>	9	89	123.587	14,189	1,148	0,151	-0,164 ^{ns}	0,075 ^{ns}
<i>Dw2</i>	6	85	140.029	13,551	0,968	0,124	-0,673 ^{ns}	0,062 ^{ns}
<i>Dw4</i>	6	54	131.451	8,609	0,655	0,238	1,432 ^{ns}	0,119 ^{ns}
<i>Mal</i>	6	39	149.808	6,217	0,415	0,293	2,348 ^{**}	0,146 ^{ns}
<i>Ma3</i>	1	52	114.138	8,290	0,726	0,185	0,457 ^{ns}	0,092 ^{ns}

Notas: *S* é o número de sítios polimórficos; θ_W é a taxa de mutação por geração (WATTERSON, 1975); π corresponde à diversidade nucleotídica, que é a proporção da soma das diferenças entre pares de indivíduos e o total de comparações de pares (TAJIMA, 1983); *D* (TAJIMA, 1989) e R_2 (RAMOS-ONSINS and ROZAS, 2002) correspondem aos testes de neutralidade; ^{ns}: Não significativo;

** : P-valor < 0,01; $\theta_{W(\text{Região})} = \frac{S}{\sum_{i=1}^{n-1} \frac{1}{i}}$; $\theta_{W(\text{Sítio})} = \frac{\theta_{W(\text{Região})}}{\text{Tamanho}}$; valores de $\theta_{W(\text{Sítio})}$ estão na escala de 10^{-4}

4.4 Estrutura Populacional e Fluxo Gênico

As estimativas de mutação (θ_{LM}) e fluxo gênico (\mathcal{M}) envolvendo as sete populações de sorgo são apresentadas na Figura 4.3. Uma inspeção desses resultados mostra que, em geral, as regiões genômicas contendo os locos *Dw1*, *Dw2*, *Dw4*, *Mal* e *Ma3* sugeriram elevadas taxas de fluxo gênico para múltiplos pares de populações, algo que não foi observado com a mesma magnitude na análise com todas as regiões simultaneamente. Embora os intervalos de credibilidade tenham sido, em geral, menores e mais precisos na análise total, as estimativas de cada região mostraram evidências importantes para a história de melhoramento do sorgo. Tais evidências não se resumem apenas às elevadas taxas de fluxo de gênico, mas também aos padrões específicos obtidos para cada região, sugerindo histórias evolutivas distintas.

Um exemplo claro da importância da análise individual ocorreu para as populações da ETI e de UGA. O número médio de indivíduos que migraram da ETI para a UGA por geração, e que forneceram contribuição

genética-evolutiva, foi bastante elevado considerando as regiões *Dw1* ($\mathcal{M}_{16} = 66,55$) e *Dw4* ($\mathcal{M}_{16} = 91,90$). No entanto, tal número de migrantes foi menor e próximo para as regiões *Dw2* ($\mathcal{M}_{16} = 11,03$) e *Ma1* ($\mathcal{M}_{16} = 11,40$) e também para o *Total* ($\mathcal{M}_{16} = 7,11$), embora ainda tenha sido elevado. Para a região *Ma3*, o número médio de indivíduos que migraram da ETI para a UGA por geração foi bastante inferior ($\mathcal{M}_{16} = 0,55$).

O número médio de migrantes por geração que forneceu contribuição genética-evolutiva variou de 0,55 (*Ma3*, ETI–UGA) a 99,10 (*Dw2*, NIG–EUA). As maiores taxas médias de fluxo gênico foram observadas para as regiões *Dw4* (52,07), *Ma1* (48,97) e *Ma3* (47,35), enquanto que as menores taxas foram verificadas para *Dw1* (44,56), *Total* (43,33) e *Dw2* (41,78). O maior fluxo gênico observado foi da NIG para os EUA na região *Dw2* (99,10), com a primeira tendo a menor quantidade de polimorfismo (5 SNPs) entre todas as populações. Algo similar também foi observado para as populações da AFS e de UGA na região *Ma3*. No primeiro caso, elevado fluxo gênico foi verificado da AFS para as populações da ETI (65,94), da NIG (88,33) e do SUD (79,23), utilizando um total de 18 SNPs. No segundo caso, elevado fluxo gênico foi observado de UGA para a IND (95,26), o SUD (93,89) e os EUA (98,19), considerando 15 SNPs.

Todos os possíveis pares envolvendo as sete populações apresentaram taxas de fluxo gênico assimétricas, seguindo o modelo proposto. Isso significa que, por exemplo, o número médio de indivíduos que migraram da ETI para o SUD ($\mathcal{M}_{15} = 46,61$) por geração não foi o mesmo que do SUD para a ETI ($\mathcal{M}_{51} = 66,67$), com base nos intervalos de credibilidade não sobrepostos na análise total. Do total de 21 combinações envolvendo pares de populações, apenas 11 mostraram assimetria a partir da análise com todas as regiões. Tais combinações assimétricas são: (i) ETI-SUD ($\mathcal{M}_{15} = 46,61$ e $\mathcal{M}_{51} = 66,66$), (ii) ETI-EUA ($\mathcal{M}_{17} = 21,95$ e $\mathcal{M}_{71} = 41,33$), (iii) IND-AFS ($\mathcal{M}_{24} = 32,16$ e $\mathcal{M}_{42} = 58,82$), (iv) NIG-AFS ($\mathcal{M}_{34} = 67,83$ e $\mathcal{M}_{43} = 46,68$), (v) NIG-SUD ($\mathcal{M}_{35} = 73,51$ e $\mathcal{M}_{53} = 26,73$), (vi) NIG-EUA ($\mathcal{M}_{37} = 88,72$ e $\mathcal{M}_{73} = 58,68$), (vii) AFS-SUD ($\mathcal{M}_{45} = 93,64$ e $\mathcal{M}_{54} = 37,46$), (viii) AFS-UGA ($\mathcal{M}_{46} = 14,59$ e $\mathcal{M}_{64} = 47,04$), (ix) AFS-EUA ($\mathcal{M}_{47} = 49,87$ e $\mathcal{M}_{74} = 35,77$), (x) SUD-EUA ($\mathcal{M}_{57} = 42,42$ e $\mathcal{M}_{75} = 33,78$) e (xi) UGA-EUA ($\mathcal{M}_{67} = 57,09$ e $\mathcal{M}_{76} = 89,08$). As demais combinações mostraram assimetria com base nos padrões específicos de cada região, mostrando, novamente, a importância da análise individual no entendimento do processo evolutivo.

4.5 Recombinação e Desequilíbrio de Ligação

As taxas de mutação (por sítio) obtidas a partir do modelo de coalescência com recombinação ($\theta_{L\rho}$) mostraram estimativas superiores às taxas de mutação (por sítio) obtidas a partir da abordagem clássica (θ_w) (Tabela 4.2). Embora $\theta_{L\rho}$ e θ_w tenham apresentado uma elevada correlação de Spearman ($r = 0,779$), os resultados diferentes mostraram que a máxima verossimilhança foi importante para obter estimativas mais precisas da variabilidade genética. Isso porque θ_w , por pouco explorar a investigação genealógica a partir dos dados, não capitaliza toda a informação disponível nos mesmos. Além disso, maior precisão deve ter sido obtida pela estimação simultânea de $\theta_{L\rho}$ e da taxa de recombinação dependente do polimorfismo ($\frac{r_L}{\mu_L}$), mostrando um cenário evolutivo possivelmente mais real.

Ao comparar as taxas de mutação (por sítio) dos modelos de fluxo gênico (θ_{LM}) e recombinação ($\theta_{L\rho}$), observou-se grande diferença em suas estimativas (Figura 4.3 e Tabela 4.2). Sempre que θ_{LM} mostrou estimativa inferior a 2×10^{-4} para certa região e população, $\theta_{L\rho}$ mostrou estimativa superior, atingindo um valor dez vezes maior ou mais para alguns casos (*Dw1*: NIG; *Dw2*: NIG; *Dw4*: NIG e AFS; *Ma1*: UGA). Em contrapartida,

sempre que θ_{LM} foi superior a 2×10^{-4} , a estimativa de $\theta_{L\rho}$ foi inferior, com θ_{LM} atingindo um valor dez vezes maior ou bem acima para alguns casos (*Dw1*: AFS, SUD e EUA; *Dw2*: AFS e UGA; *Dw4*: ETI, IND, UGA e EUA; *Ma1*: IND, AFS e EUA; *Ma3*: ETI, IND e EUA). A correlação de Spearman entre θ_{LM} e $\theta_{L\rho}$ foi igual à 0,484, mostrando uma estimativa inferior quando $\theta_{L\rho}$ foi correlacionado com θ_W da abordagem clássica ($r = 0,779$).

A região genômica que mostrou a maior taxa de mutação média (escala de 10^{-4}) foi a *Dw1* (2,27), seguida das regiões *Dw4* (1,36), *Dw2* (1,30), *Ma3* (1,17) e *Ma1* (0,73). As maiores variabilidades genéticas, em ordem decrescente e de acordo com a região e a população, foram: (i) 3,414 (*Dw1*, AFS); (ii) 2,780 (*Dw1*, EUA); (iii) 2,722 (*Dw1*, UGA); (iv) 2,659 (*Dw2*, EUA); (v) 2,189 (*Dw1*, NIG); (vi) 2,182 (*Ma3*, EUA); e (vii) 2,054 (*Dw4*, EUA).

As taxas de recombinação constante (LAMARC) que são dependentes da mutação ($\frac{r_L}{\mu_L}$) apresentaram variações entre as regiões e populações (Tabela 4.2). A maior taxa média de $\frac{r_L}{\mu_L}$ (escala de 10^{-3}) por Kb foi para a região *Dw2* (1,578), seguida das regiões *Ma3* (1,209), *Dw1* (0,347), *Ma1* (0,213) e *Dw4* (0,129). As médias elevadas de $\frac{r_L}{\mu_L}$ das regiões *Dw2* e *Ma3* ocorreram principalmente devido às altas estimativas das populações da NIG (10,159) e da IND (6,789), respectivamente. As taxas de recombinação independentes da mutação (ρ_L) mostraram menores variações, com o maior valor médio (escala de 10^{-3}) por Kb apresentado pela região *Ma3* (0,246), seguida das regiões *Dw2* (0,121), *Dw1* (0,115), *Dw4* (0,057) e *Ma1* (0,030). As taxas mais elevadas de ρ_L ocorreram para a região *Ma3*, por meio das populações da IND (0,882), da ETI (0,329) e o Subpainele (0,346). A taxa de ρ_L da NIG na região *Dw2* (0,224), embora elevada, mostrou ser bastante inferior ao valor de ρ_L da IND na região *Ma3* (0,882). Ambas tiveram valores parecidos de recombinação quando houve dependência da mutação, mas apresentaram valores bem distintos na sua independência.

A Tabela 4.3 apresenta as taxas de recombinação mínima ($\rho_{r(\text{Min})}$), média ($\rho_{r(\text{Med})}$) e máxima ($\rho_{r(\text{Max})}$), estimadas via coalescência (rhomap), obtidas para cada região genômica e população. Ao contrário das taxas de recombinação mostradas na Tabela 4.2, que foram obtidas a partir do modelo de recombinação constante (KUHNER ET AL., 2000), tais recombinações foram obtidas a partir do modelo de recombinação variável ao longo da região (AUTON and McVEAN, 2007). De forma geral, os padrões obtidos indicaram processos de recombinação específicos para cada região, e também para cada população dentro de certa região. Por exemplo, as menores e maiores taxas de recombinação detectadas para as regiões foram: (i) 0,053 (UGA) e 1,721 (ETI) para *Dw1*; (ii) 0,021 (UGA) e 76,320 (Subpainele) para *Dw2*; (iii) 0,010 (NIG) e 1,625 (ETI) para *Dw4*; (iv) 0,004 (AFS) e 1,934 (EUA) para *Ma1*; e (v) 0,058 (SUD) e 2,039 (IND) para *Ma3*. Considerando todas as populações, a maior taxa de recombinação média foi a da região *Dw2* (3,395), seguida das regiões *Ma3* (0,465), *Dw1* (0,403), *Dw4* (0,241) e *Ma1* (0,060). Além disso, algumas populações apresentaram recombinações máximas bastante diferentes da taxa de recombinação média para certa região, como foi o caso da NIG (6,729) e dos EUA (38,439) na região *Dw2*. As taxas de recombinação dos EUA e do Subpainele nesta região foram extremamente elevadas, estando muito além do valor médio observado.

Enquanto a Tabela 4.3 apresenta as taxas de recombinação mínima, média e máxima para todas as regiões genômicas e populações, o Gráfico 4.4 mostra as taxas de recombinação (ρ_r/Kb) entre todos os SNPs adjacentes para algumas regiões e populações. Tais regiões mostraram padrões bastante variáveis de recombinação, e portanto foram utilizadas como exemplo para discussão. Esses padrões podem ser claramente visualizados por meio da escala diferente do eixo y. O caso (A), correspondendo à população da ETI na região *Dw1*, e o caso (D),

Tabela 4.2: Taxas de mutação ($\theta_{L\rho}$) e recombinação ($\frac{r_L}{\mu_L}$ e ρ_L por Kb) estimadas via coalescência (LAMARC) para cada uma das cinco regiões genômicas e sete populações de sorgo, bem como para o subpainel de 298 indivíduos. A taxa de mutação estimada via abordagem clássica (θ_W) também é mostrada, e ambas θ_W e $\theta_{L\rho}$ (escala de 10^{-4}) correspondem ao número médio de mutações por sítio. $\frac{r_L}{\mu_L}$ (dependente da mutação) e ρ_L (independente da mutação) (escala de 10^{-3}) correspondem ao número médio de recombinações entre dois indivíduos quaisquer até a ocorrência de um evento de coalescência. Neste modelo (KUHNER ET AL., 1995), $\theta_{L\rho}$, $\frac{r_L}{\mu_L}$ e ρ_L são constantes ao longo das gerações. IC é o intervalo de confiança. $E[r_{LAMARC}^2]$ corresponde ao desequilíbrio de ligação de alta ordem calculado a partir do ρ_L de toda a região

Região	População	S^1	θ_W (Sítio)	$\theta_{L\rho}$ (Sítio)	IC $\theta_{L\rho}$ (95%)	$\frac{r_L}{\mu_L}$ /Kb	IC $\frac{r_L}{\mu_L}$ (95%)	ρ_L /Kb	$E[r_{LAMARC}^2]$
<i>Dw1</i>	ETI	47	0,924	1,275	0,771 – 2,153	0,237	0,119 – 0,451	0,042	0,996
	IND	50	1,083	1,635	0,943 – 2,427	0,194	0,093 – 0,368	0,044	0,996
	NIG	58	1,288	2,189	1,339 – 4,191	0,237	0,114 – 0,479	0,072	0,994
	AFS	60	1,565	3,414	2,477 – 5,466	0,237	0,119 – 0,451	0,042	0,996
	SUD	53	1,021	1,851	1,181 – 2,441	0,449	0,279 – 0,605	0,115	0,989
	UGA	47	1,196	2,722	1,412 – 5,336	0,306	0,128 – 0,491	0,116	0,989
	EUA	86	1,243	2,780	2,283 – 3,414	0,547	0,425 – 0,694	0,211	0,982
	Subpainel	89	1,148	3,572	3,025 – 4,105	0,565	0,455 – 0,726	0,281	0,976
<i>Dw2</i>	ETI	31	0,538	0,892	0,455 – 1,214	0,628	0,314 – 0,912	0,092	0,992
	IND	31	0,593	1,153	0,637 – 1,736	0,301	0,116 – 0,541	0,057	0,995
	NIG	05	0,098	0,134	0,097 – 0,170	10,159	5,202 – 14,19	0,224	0,982
	AFS	36	0,829	1,402	0,880 – 2,871	0,187	0,061 – 0,347	0,043	0,997
	SUD	68	1,156	1,718	1,238 – 2,477	0,217	0,148 – 0,365	0,062	0,995
	UGA	42	0,943	1,129	0,656 – 3,600	0,108	0,029 – 0,383	0,020	0,998
	EUA	81	1,033	2,659	2,125 – 2,981	0,394	0,339 – 0,524	0,172	0,986
	Subpainel	85	0,968	2,871	2,428 – 3,211	0,627	0,541 – 0,768	0,296	0,976
<i>Dw4</i>	ETI	46	0,850	1,282	0,756 – 2,029	0,130	0,022 – 0,146	0,040	0,998
	IND	42	0,856	1,343	0,653 – 2,469	0,040	0,010 – 0,093	0,013	0,999
	NIG	38	0,793	0,892	0,489 – 1,298	0,041	0,000 – 0,067	0,009	0,999
	AFS	33	0,809	0,896	0,468 – 5,525	0,006	0,000 – 0,068	0,001	0,999
	SUD	43	0,779	1,339	0,793 – 2,025	0,095	0,067 – 0,234	0,031	0,998
	UGA	38	0,909	1,693	0,781 – 2,867	0,096	0,038 – 0,179	0,040	0,998
	EUA	49	0,666	2,054	1,389 – 2,338	0,253	0,214 – 0,403	0,127	0,993
	Subpainel	54	0,655	2,112	1,721 – 2,329	0,374	0,301 – 0,530	0,192	0,989
<i>Ma1</i>	ETI	28	0,454	0,549	0,234 – 0,963	0,010	0,000 – 0,032	0,002	0,999
	IND	30	0,536	0,896	0,521 – 1,401	0,063	0,032 – 0,125	0,022	0,999
	NIG	25	0,458	0,555	0,331 – 1,054	0,005	0,000 – 0,025	0,001	0,999
	AFS	25	0,538	0,490	0,232 – 1,070	0,000	0,000 – 0,004	0,000	1,000
	SUD	26	0,413	0,607	0,336 – 0,836	0,036	0,007 – 0,081	0,008	0,999
	UGA	24	0,504	1,291	0,422 – 2,153	0,007	0,000 – 0,028	0,004	0,999
	EUA	39	0,465	0,687	0,620 – 0,765	0,283	0,244 – 0,349	0,075	0,997
	Subpainel	39	0,415	0,247	0,229 – 0,268	1,302	1,163 – 1,515	0,124	0,995
<i>Ma3</i>	ETI	35	0,745	1,144	0,943 – 1,326	1,311	1,139 – 1,621	0,329	0,983
	IND	34	0,798	0,592	0,529 – 0,642	6,789	6,215 – 7,717	0,882	0,957
	NIG	36	0,865	1,704	1,062 – 2,738	0,142	0,076 – 0,285	0,053	0,997
	AFS	18	0,508	0,957	0,451 – 1,986	0,072	0,000 – 0,190	0,015	0,999
	SUD	28	0,584	1,153	0,738 – 1,690	0,186	0,094 – 0,358	0,047	0,997
	UGA	15	0,413	0,478	0,219 – 1,075	0,000	0,000 – 0,031	0,000	1,000
	EUA	51	0,798	2,182	1,923 – 2,679	0,619	0,493 – 0,717	0,297	0,985
	Subpainel	52	0,726	2,848	2,556 – 3,422	0,554	0,429 – 0,622	0,346	0,983

¹Número de sítios polimórficos

Tabela 4.3: Taxas de recombinação mínima ($\rho_{r(\text{Min})}$), média ($\rho_{r(\text{Med})}$) e máxima ($\rho_{r(\text{Max})}$) estimadas via coalescência (rhomap) para cada uma das cinco regiões genômicas e sete populações de sorgo, bem como para o subpainel de 298 indivíduos. Neste modelo de coalescência (AUTON and McVEAN, 2007), tais taxas correspondem ao número esperado de recombinações que acontece para determinado conjunto (par) de SNPs em uma distância de 1 Kb. Essas taxas são consideradas variáveis ao longo da região de interesse

Região	População	$\rho_{r(\text{Min})}$	$\rho_{r(\text{Med})}$	$\rho_{r(\text{Max})}$
<i>Dw1</i>	ETI	0,344	0,680	1,721
	IND	0,073	0,234	1,169
	NIG	0,110	0,285	1,076
	AFS	0,115	0,364	1,252
	SUD	0,110	0,406	1,435
	UGA	0,053	0,223	0,698
	EUA	0,145	0,432	1,320
	Subpainel	0,186	0,599	1,712
<i>Dw2</i>	ETI	0,329	0,768	2,029
	IND	0,035	0,329	1,548
	NIG	0,302	2,420	6,729
	AFS	0,183	0,491	1,583
	SUD	0,285	0,888	1,954
	UGA	0,021	0,273	1,301
	EUA	0,210	7,125	38,439
	Subpainel	0,346	14,868	76,320
<i>Dw4</i>	ETI	0,060	0,288	1,625
	IND	0,036	0,061	0,181
	NIG	0,010	0,033	0,621
	AFS	0,100	0,278	0,905
	SUD	0,117	0,334	1,146
	UGA	0,079	0,381	1,131
	EUA	0,089	0,291	1,554
	Subpainel	0,148	0,265	0,660
<i>Ma1</i>	ETI	0,009	0,011	0,021
	IND	0,018	0,162	1,357
	NIG	0,009	0,022	0,134
	AFS	0,004	0,006	0,018
	SUD	0,013	0,016	0,025
	UGA	0,013	0,017	0,050
	EUA	0,039	0,159	1,934
	Subpainel	0,034	0,085	1,053
<i>Ma3</i>	ETI	0,279	0,619	1,132
	IND	0,558	0,838	2,039
	NIG	0,152	0,366	1,499
	AFS	0,171	0,340	0,551
	SUD	0,058	0,194	0,898
	UGA	0,599	0,695	0,938
	EUA	0,198	0,340	0,756
	Subpainel	0,181	0,327	1,241

correspondendo à população da IND na região *Ma3*, parecem se assemelhar na variação de ρ_r , que permanece entre 0,500 e 2,000 recombinações médias por Kb. O caso (B), que corresponde à população dos EUA na região *Dw2*, mostra um valor de 38,439 recombinações médias por Kb, uma estimativa bastante elevada e diferente dos padrões das demais regiões e populações. Por fim, o caso (C) mostra os baixos padrões de recombinação da população do SUD na região *Ma1*, não ultrapassando o valor de 0,025 recombinações médias por Kb.

O DL calculado via medida tradicional (r^2) e coalescência ($E[r_{\text{rhomap}}^2]$) é mostrado na Figura 4.5. É possível observar que os padrões de ambos foram muito próximos em alguns casos. Foi o que aconteceu, por exemplo, para as regiões e populações *Dw2* (IND), *Dw2* (UGA) e *Ma1* (IND, NIG, AFS, SUD e UGA), onde, em geral, forte DL foi observado em grandes blocos de associação. Contudo, é possível perceber que o r^2 não parece fornecer uma definição clara da extensão de possíveis blocos para casos em que os padrões são completamente elevados na região. Um exemplo claro acontece com a população da NIG na região *Dw4*. Embora ambas as medidas evidenciam que há forte DL nesta região, há certa descontinuidade nos padrões mostrados por r^2 , algo que não acontece com o $E[r_{\text{rhomap}}^2]$. Possivelmente, esse comportamento está associado ao viés comumente verificado com o uso do r^2 , de modo que o DL via coalescência ($E[r_{\text{rhomap}}^2]$) pode ser fundamental para uma melhor definição das associações.

O DL de alta ordem ($E[r_{\text{LAMARC}}^2]$, Tabela 4.2) sugeriu estimativas extremamente elevadas (0,957 – 1,000) para todas as regiões genômicas e populações. Tais estimativas foram obtidas a partir das baixas taxas de recombinação (ρ_L) sugeridas pelo LAMARC considerando o polimorfismo de toda a região genômica. Embora houve casos em que o r^2 e o $E[r_{\text{rhomap}}^2]$ foram bastante elevados e mostraram estar próximos de $E[r_{\text{LAMARC}}^2]$, este DL não mostrou estimativa representativa das associações para todas as regiões e populações. Tal evidência se torna ainda mais clara ao comparar as estimativas de 0,957 – 0,999 do $E[r_{\text{LAMARC}}^2]$ com o DL médio a partir das estimativas do r^2 (Figura 4.5), que variaram da seguinte forma, de acordo com a região genômica: (i) 0,127 – 0,351 para região *Dw1*; (ii) 0,106 – 0,486 para a região *Dw2*; (iii) 0,212 – 0,567 para a região *Dw4*; (iv) 0,259 – 0,713 para a região *Ma1*; e (iv) 0,059 – 0,320 para a região *Ma3*.

Os padrões observados na Figura 4.5 claramente evidenciaram uma estrutura específica do DL para as populações ao longo das regiões. Essa especificidade pode ser visualizada de duas diferentes maneiras: (i) DL entre regiões genômicas; e (ii) DL entre populações para uma dada região. No primeiro caso, é possível perceber que o DL variou consideravelmente de acordo com a região do genoma. As regiões que pareceram ter DL mais extensivo foram *Ma1*, *Dw4* e *Dw2*, em ordem decrescente de associação. As regiões *Dw1* e *Ma3* também apresentaram DL relativamente extensivo, só que em menor proporção e para uma ou outra população.

No segundo caso, o DL mostrou-se de forma bastante específica de acordo com a população. Com exceção da região *Ma1*, em que praticamente todas as populações apresentaram forte bloco de associação, todas as demais tiveram populações com padrões específicos de DL. Por exemplo, as populações da IND e de UGA mostraram maior DL em comparação com as demais para ambas as regiões *Dw1* e *Dw2*. Isso também ocorreu com as populações da ETI, IND e NIG na região *Dw4*, e do SUD na região *Ma3*.

Um resultado muito curioso foi verificado para o Subpainel de indivíduos. Para todas as regiões, o DL presente no Subpainel foi inferior ou próximo àquele verificado para as demais populações. Isso mostra que a mistura populacional de indivíduos oriundos de diferentes localidades não gerou acréscimos nos padrões do DL, algo que pode eventualmente acontecer.

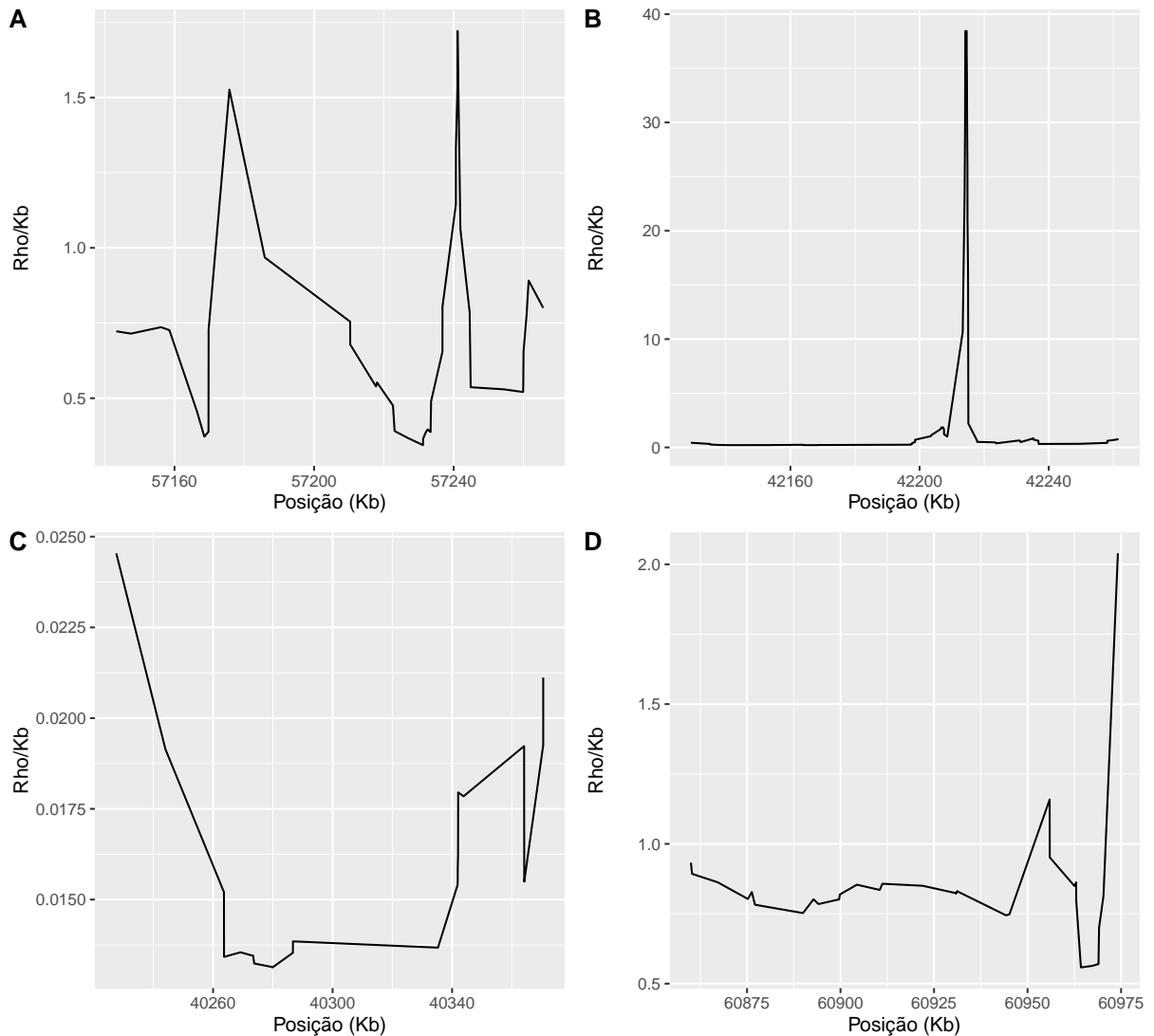


Figura 4.4: Taxas de recombinação (ρ_r/Kb) estimadas via coalescência (*rhomap*) entre todos os pares de SNPs adjacentes para algumas das regiões e populações mostradas na Tabela 4.3. No presente gráfico, todas as taxas de recombinação entre SNPs adjacentes são apresentadas, e não apenas as taxas mínima, média e máxima, como foi feito na referida Tabela. (A) Região *Dw1*, População da ETI; (B) Região *Dw2*, População dos EUA; (C) Região *Mal*, População do SUD; e (D) Região *Ma3*, População da IND. Note a variação de escala mostrada no eixo y, indicando taxas de recombinação bastante distintas de acordo com a região genômica e a população. No caso (B), uma taxa média de 38,439 recombinções por Kb é mostrada entre dois SNPs adjacentes para a população dos EUA na região *Dw2*, sugerindo um possível ponto quente de recombinação (*hotspot*)

As correlações de Spearman entre r^2 e $E[r_{\text{rhomap}}^2]$ são mostradas na Tabela 4.4. É possível observar que a associação entre ambas as medidas é muito próxima de zero para todas as regiões genômicas e populações. Embora há vários casos em que r^2 e $E[r_{\text{rhomap}}^2]$ mostraram padrões de associação distintos (Figura 4.5), os resultados das correlações pouco refletem aquilo que, em geral, foi verificado nos *heatmaps* de DL. Isto ocorre porque as medidas de correlação são uma estatística de resumo.

Tabela 4.4: Correlações de Spearman entre o desequilíbrio de ligação via medida tradicional (r^2) e coalescência ($E[r_{\text{rhomap}}^2]$) para cinco regiões genômicas e sete populações de sorgo, bem como para o subpainel de 298 indivíduos

População	Regiões Genômicas				
	<i>Dw1</i>	<i>Dw2</i>	<i>Dw4</i>	<i>Ma1</i>	<i>Ma3</i>
ETI	-0,097	-0,147	0,024	-0,020	0,064
IND	-0,117	0,199	0,052	-0,055	-0,057
NIG	0,015	0,103	0,028	0,038	-0,052
AFS	-0,006	0,054	0,016	-0,077	-0,023
SUD	-0,063	-0,110	-0,122	-0,117	0,156
UGA	-0,023	-0,186	-0,067	-0,034	0,094
EUA	-0,027	0,001	-0,021	0,092	0,036
Subpainel	-0,024	-0,044	-0,013	-0,006	0,071

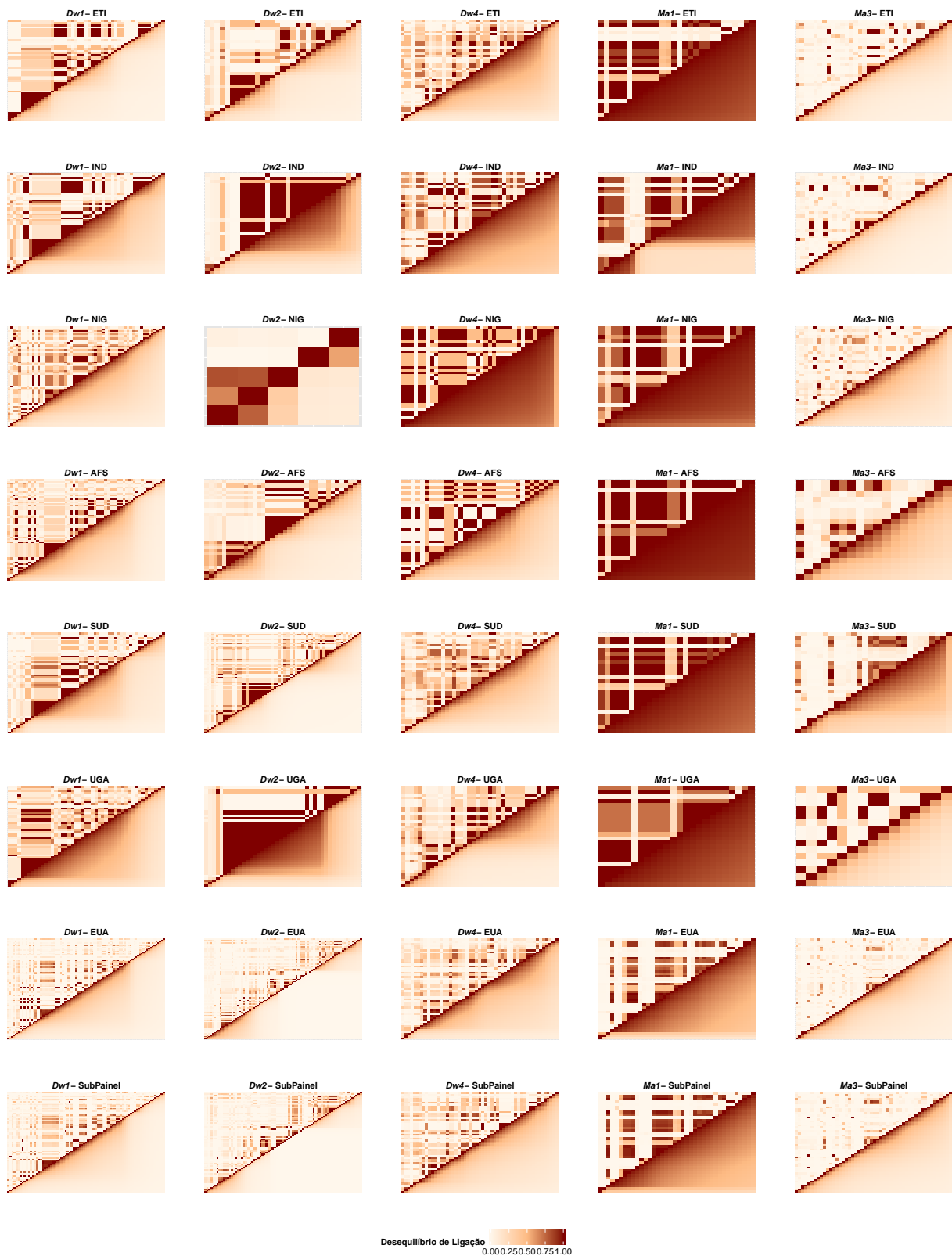


Figura 4.5: Gráfico do desequilíbrio de ligação via medida tradicional (r^2) e coalescência ($E[r_{\text{rhomap}}^2]$) para cada região genômica e população de sorgo. Os padrões do r^2 e $E[r_{\text{rhomap}}^2]$ são mostrados acima e abaixo da diagonal principal dos gráficos de *heatmap*, respectivamente. O $E[r_{\text{rhomap}}^2]$ foi calculado com base nas estimativas de recombinação variável obtidas via *rhomap* (ρ_r), considerando uma população em equilíbrio entre recombinação e deriva genética

5 DISCUSSÃO

No presente estudo, os princípios da teoria da coalescência foram utilizados para investigar a estrutura populacional e o desequilíbrio de ligação de um painel mundial de acessos de sorgo. Para tanto, análises de mutação, migração com fluxo gênico e recombinação foram realizadas para cinco regiões genômicas e sete populações que constituem este painel. Tais investigações foram realizadas com o propósito de fornecer informações relevantes para estudos futuros de mapeamento associativo em sorgo, bem como para um melhor entendimento do próprio painel. Além disso, como objetivo secundário, havia a expectativa de avaliar os modelos e pacotes disponíveis para análises baseadas em coalescência.

Uma das principais pressuposições assumida por grande parte dos modelos de coalescência é que as populações se comportam segundo o modelo proposto por FISHER (1930) e WRIGHT (1931). Assume-se que uma amostra recente de indivíduos foi aleatoriamente tomada a partir de uma grande população panmítica, que é baseada em cruzamentos aleatórios entre os indivíduos. No entanto, o painel aqui utilizado foi construído de maneira a agrupar indivíduos historicamente importantes para o melhoramento do sorgo, não sendo necessariamente uma amostra aleatória de uma população panmítica. Assim, as estimativas aqui obtidas podem ter apresentado certo viés e variâncias elevadas (HAMBLIN *ET AL.*, 2004), além do fato de que a taxa de autofecundação em populações cultivadas de sorgo pode ser superior a 90% (HAMBLIN *ET AL.*, 2005). Isso significa que os cruzamentos não acontecem, em sua grande maioria, de forma cruzada e aleatória, tendo implicações na forma pela qual a teoria da coalescência é proposta (FU, 1997; NORDBORG and DONNELLY, 1997; HAMBLIN *ET AL.*, 2005). Portanto, estudos futuros de coalescência em sorgo deveriam levar essa realidade em consideração.

O número reduzido de indivíduos para certas populações pode ter contribuído para a obtenção de menor polimorfismo para as regiões genômicas de interesse. Consequentemente, as estimativas de mutação, fluxo gênico e recombinação podem ter apresentado certo viés em relação ao processo evolutivo da espécie, embora o interesse tenha permanecido especificamente no painel aqui utilizado. Normalmente, o aumento da amostra resulta na elevação logarítmica da acurácia a partir da estimação por máxima verossimilhança (FELSENSTEIN, 2006). No entanto, esse aumento deve ocorrer de forma a elevar o polimorfismo, pois uma quantidade inferior a 20 locos pode resultar em estimativas pouco confiáveis dos parâmetros, principalmente quando o grau de diferenciação entre populações é pequeno (WILSON and RANNALA, 2003). No presente estudo, poucas foram as populações que apresentaram quantidades inferiores a 20 SNPs, mas o polimorfismo foi considerado baixo por se tratar de variações de nucleotídeos únicos. Embora os SNPs sejam, em geral, abundantes ao longo do genoma, normalmente são menos informativos quando comparados com outros marcadores moleculares, tais como os microssatélites.

A utilização da técnica de GBS possibilitou a obtenção de dados de sequências para a realização deste trabalho. Muitos dados perdidos foram inicialmente obtidos, e o método de imputação via NPUTE foi utilizado para a recuperação dessas informações. Embora grande parte do genoma dos indivíduos deva estar em homozigose, o NPUTE desconsiderou a possibilidade de haver locos em heterozigose, o que poder ser algo pouco real. Informações de heterozigotos certamente seriam importantes para se ter maior variação genética, influenciando diretamente na estimação dos parâmetros populacionais. No entanto, acredita-se que, em geral, a técnica de GBS, aliada à imputação via NPUTE, forneceu genótipos bastante razoáveis.

Desvios do modelo padrão baseado em neutralismo apenas foram detectados para a região *Mal*, com

base na estatística D (TAJIMA, 1989). Embora esta região tenha apresentado as menores taxas médias de recombinação populacional via LAMARC ($\rho_L = 0,030 \times 10^{-3}/\text{Kb}$) e rhomap ($0,060 \rho_r/\text{Kb}$), sua reduzida variabilidade pode ter sido consequência de um processo de seleção. Tal evidência também foi sugerida por KLEIN ET AL. (2008) ao detectarem menor variação na região que contém os locos *Mal* e *Dw2*, sugerindo, especificamente, a ocorrência de seleção direcional em decorrência das decisões do melhoramento no processo de conversão de linhagens. Contudo, o valor positivo de D aqui detectado, indicando a presença de ramos internos mais longos na genealogia, pode sugerir a ocorrência de seleção balanceadora ou até mesmo uma recente mistura entre diferentes populações (HEIN ET AL., 2004).

No entanto, outros processos evolutivos podem ter contribuído para a significância da estatística D (WAKELEY, 2009; KORNELIUSSEN ET AL., 2013) na região *Mal*, tais como fluxo gênico e endogamia, mostrando que D , e outros testes clássicos de neutralidade, devem pouco esclarecer a história evolutiva ocorrida (FU and LI, 1999; WAKELEY, 2009; KUHNER, 2009). Mais do que isso, é difícil saber de que maneira D é afetado na presença de vários processos evolutivos (HEIN ET AL., 2004; KUHNER, 2009). Portanto, deve-se considerar D apenas como evidência inicial de que processos mais complexos podem ter ocorrido ao longo das gerações, lembrando que valores não significativos de D podem desconsiderar evidências evolutivas importantes (MORRELL ET AL., 2003) devido ao uso restrito da variação.

A estimação de fluxo gênico usando o método proposto por BEERLI and FELSENSTEIN (2001) mostrou evidências importantes da história de melhoramento do sorgo. Embora este método seja computacionalmente intensivo (PALCZEWSKI and BEERLI, 2013) por ser baseado em abordagem completa, as regiões aqui investigadas não apresentaram grandes quantidades de SNPs a ponto de torná-lo inviável. Assim, o fluxo gênico pôde ser estimado a partir de diversas simulações, repetições e diferentes estratégias de busca no contexto bayesiano, na expectativa de que muitas histórias pudessem ser percorridas e supostamente consideradas.

Em geral, elevadas taxas de fluxo gênico foram obtidas para todas as regiões genômicas e todas elas simultaneamente. Isso significa que intenso intercâmbio de acessos deve ter historicamente ocorrido entre as sete populações aqui consideradas, muitas das quais estão situadas em regiões geográficas próximas. No caso dos EUA, que está localizado mais distantemente, o programa de conversão de linhagens deve ter favorecido a entrada de acessos provenientes da África, sendo, em geral, de maior magnitude em comparação à entrada de acessos dos EUA na África. Assim, acredita-se que as taxas de fluxo gênico aqui detectadas retratam eventos migratórios mais recentes, embora parte deles possa ter ocorrido em gerações mais antigas do período de domesticação, envolvendo populações africanas mais próximas. Eventos migratórios mais antigos podem representar a grande motivação do uso da teoria da coalescência em comparação aos métodos clássicos para inferir estrutura populacional. Além disso, o elevado fluxo gênico aqui detectado também pode indicar que as populações foram recentemente formadas a partir de uma mesma origem, a qual deve ter ocorrido na Etiópia e no Sudão.

Embora o fluxo gênico tenha sido considerado como elevado, é preciso ter cuidado com a sua magnitude. Originalmente, o método proposto por BEERLI and FELSENSTEIN (2001) utiliza o polimorfismo disponível na amostra para estimar, simultaneamente, as taxas de mutação (θ_{LM}) e fluxo gênico (\mathcal{M}). \mathcal{M} corresponde à razão $\frac{m}{\mu}$ (BEERLI and FELSENSTEIN, 1999), mostrando que esta estimativa é dependente da mutação. O produto de \mathcal{M} pelo valor de θ da população recipiente resulta no número médio de migrantes por geração (neste caso, $2N_e m$), sendo este independente da mutação (BEERLI, 1998). Esta quantidade é preferida porque altos valores de \mathcal{M} podem indicar

tanto elevado fluxo gênico quanto baixa mutação, dificultando uma real interpretação (BEERLI, 1998). Contudo, no cenário aqui investigado, as taxas de mutação apresentaram, em geral, grandes intervalos de credibilidade. Isso significa que $2N_e m$ poderia ser sub ou superestimado, apresentando intervalos de credibilidade bastante elevados.

Grandes intervalos de credibilidade foram encontrados para as estimativas de fluxo gênico M , sugerindo um cenário bastante complexo. Isso pode ter acontecido pelo fato de que poucos locos estiveram disponíveis para poucas regiões genômicas, reduzindo a variação genética entre populações (BEERLI, 1998; WILSON and RANNALA, 2003; BEERLI, 2006). Com elevadas taxas de fluxo gênico, grande parte das mutações estará presente em todas as populações, de modo que poucos locos poderão ser insuficientes para se ter maior variação e confiança no processo de estimação (BEERLI, 2006). Além disso, muitos eventos terão que ser investigados pelo algoritmo de MCMC caminhando-se do presente em direção ao passado, algo que se torna ainda mais complexo à medida que aumenta o número de populações, de indivíduos e, assim, a quantidade de parâmetros (PALCZEWSKI and BEERLI, 2013). Portanto, a estimação confiável de fluxo gênico entre populações que estiveram em constante e intenso contato é um grande desafio (FAUBET *ET AL.*, 2007; PALCZEWSKI and BEERLI, 2013).

Outro ponto importante é que a análise de fluxo gênico do presente estudo foi realizada no contexto bayesiano. Embora tal abordagem tem sido recomendada para cenários demográficos complexos (BEAUMONT and RANNALA, 2004; BEERLI, 2006), a estimação dos parâmetros pode ser sensivelmente dependente das informações estabelecidas a priori (STEPHENS, 2007). Mesmo que a distribuição a posteriori não apresente similaridade com a distribuição a priori, o que é indicativo de que variação suficiente está disponível nos dados (BEERLI, 2006), a atribuição de informações iniciais é fundamental para a estimação confiável dos parâmetros. No presente estudo, distribuições típicas foram utilizadas como priori para θ (linear) e M (logarítmica) (KUHNER and SMITH, 2007). Porém, o cenário aqui investigado é bastante complexo, sendo recomendado, em trabalhos futuros, a utilização de outras distribuições a priori consideradas importantes, tais como exponencial e uniforme (BEERLI, 2006).

A análise de fluxo gênico para cada região individualmente mostrou evidências importantes da história evolutiva do sorgo. Por um lado, o uso simultâneo de múltiplas regiões não ligadas pode fornecer maior variação genética (HEY and NIELSEN, 2004; BEERLI, 2006) e possibilitar a investigação de processos que afetam todo o genoma (SOUSA and HEY, 2013). Por outro, tais regiões podem evidenciar taxas de fluxo gênico variáveis entre as populações, mostrando histórias evolutivas específicas. Na presença de seleção, por exemplo, certas regiões podem experimentar completa ausência de fluxo gênico entre populações (SOUSA *ET AL.*, 2013; SOUSA and HEY, 2013). No presente estudo, fluxo gênico assimétrico foi observado entre todos os pares de populações somente quando a análise individual das regiões foi considerada. Portanto, a análise de fluxo gênico para cada região genômica foi muito importante para detectar um cenário evolutivo mais complexo e, possivelmente, mais real.

O método proposto por BEERLI and FELSENSTEIN (2001) também foi utilizado em trabalhos anteriores envolvendo espécies domesticadas. MORRELL *ET AL.* (2003) investigaram o fluxo gênico entre três populações selvagens de cevada (*Hordeum vulgare* ssp. *spontaneum*) utilizando dados de sequência de nove locos genômicos. Em média, um indivíduo migrante por geração foi observado entre todas as populações a partir da análise contendo todos os locos. Porém, taxas de fluxo gênico ($M = 4N_e m$) elevadas (11,664, 17,844 e 21,272) foram detectadas na análise para cada loco separadamente, obtendo-se, muitas vezes, grandes intervalos de confiança. Embora os autores tenham concentrado suas conclusões na análise total, resultados específicos foram obtidos para cada loco, algo que também foi observado no presente estudo.

Também para cevada, HÜBNER ET AL. (2012) estimaram a taxa de fluxo gênico entre três populações selvagens e entre populações selvagem e cultivada utilizando marcadores SSRs. Os autores encontraram elevado fluxo gênico entre as três populações selvagens (2,87–8,45) e da população cultivada para a selvagem (3,14), como consequência de cruzamentos ocasionais e da dispersão de sementes realizada pelo homem e por outros animais.

As análises de coalescência mostraram que as taxas de recombinação variaram de acordo com a região genômica e a população em sorgo. Isso mostra que o valor médio de recombinação (1,4 ρ /Kb, *rhomap*) sugerido por MORRIS ET AL. (2013) é insuficiente para caracterizar os padrões de recombinação no genoma desta espécie. Os autores argumentam que diferentes taxas foram observadas para regiões centroméricas e teloméricas. No entanto, os resultados aqui obtidos via *rhomap* sugerem que as taxas de recombinação podem ser bastante variáveis para regiões contendo locos importantes para o melhoramento do sorgo, tais como *Dw2* (3,395 ρ_r /Kb), *Ma1* (0,060 ρ_r /Kb) e *Ma3* (0,465 ρ_r /Kb). Embora os dois primeiros estejam fisicamente ligados no cromossomo 6 (LIN ET AL., 1995; KLEIN ET AL., 2008), seus padrões de recombinação sugerem histórias distintas para cada população e o subpainel utilizado.

Em geral, as taxas médias de recombinação constantes e independentes da mutação (ρ_L /Kb), estimadas via LAMARC no presente estudo, mostraram estimativas superiores daquelas obtidas por HAMBLIN ET AL. (2005). Esses autores investigaram os padrões de recombinação em sorgo com base nos métodos aproximados de HUDSON (2001) e LI and STEPHENS (2003). Para tanto, seis regiões genômicas não ligadas foram analisadas para um painel constituído de 32 acessos (cultivados e selvagens) de diferentes localidades do mundo. Uma única e constante taxa de recombinação foi obtida para cada uma dessas seis regiões, assim como foi realizado no presente estudo utilizando o método de KUHNER ET AL. (2000). As taxas de recombinação obtidas por HAMBLIN ET AL. (2005) variaram de 2,073 – 15,660 $\times 10^{-6}$ /Kb e 0,997 – 9,957 $\times 10^{-6}$ /Kb utilizando os métodos de HUDSON (2001) e LI and STEPHENS (2003), respectivamente.

HAMBLIN ET AL. (2005) mostraram que a recombinação foi inferior para acessos cultivados em comparação aos acessos selvagens, considerando cinco das seis regiões estudadas. No entanto, uma comparação entre populações cultivadas também pode ser desejável. O nosso estudo mostrou que as taxas de recombinação entre populações cultivadas podem ser consideravelmente diferentes para certa região. Tal diferença foi aqui baseada no método de AUTON and McVEAN (2007), que forneceu taxas de recombinação variáveis ao longo de certa região.

O caso mais extremo dessa diferença ocorreu para a população dos EUA (38,439 ρ /Kb) e o Subpainel (76,320 ρ /Kb) na região *Dw2*. Grande parte dos eventos de recombinação aqui detectados deve ter ocorrido recentemente no contexto de coalescência, já que não há indivíduos selvagens de períodos anteriores à domesticação. Assim, acredita-se que as taxas elevadas de tais populações correspondam a pontos quentes de recombinação (*hospots*), com sítios adjacentes mostrando, em média, taxas próximas àquelas verificadas para as demais regiões. Taxa de recombinação elevada também foi observada para a população da NIG (6,729 ρ /Kb) na região *Dw2*, embora tal valor possa ter sido subestimado pela baixa quantidade de SNPs.

Estimativas de recombinação via coalescência também foram obtidas para outras espécies autógamas. Em *Arabidopsis thaliana*, NORDBORG ET AL. (2005) e KIM ET AL. (2007) encontraram taxas médias de 0,200 e 0,800 ρ /Kb para todo o genoma, respectivamente. KIM ET AL. (2007) observaram taxas mais elevadas (1,500 ρ /Kb) a curtas distâncias (0 – 5 Kb), obtendo-se valores próximos à média (0,800 ρ /Kb) a partir de 5 Kb. Os autores sugeriram a ocorrência de conversão gênica para explicar tais padrões, que evidenciam certa discrepância entre DL

de curtas e longas distâncias. Em sorgo, poucas evidências de recombinação foram observadas a curtas distâncias (CANIATO ET AL., 2014), sugerindo que a conversão gênica não deve explicar discrepâncias de associação para este espécie (HAMBLIN ET AL., 2005).

Em cevada, MORRELL ET AL. (2003) observaram taxas de recombinação variando de 0,094–0,511/intersítio para nove genes de importância. Esses autores utilizaram o mesmo método de KUHNER ET AL. (2000) investigado no presente estudo, com o qual uma única e constante taxa de recombinação é estimada para certa região. MORRELL ET AL. (2003) obtiveram as taxas de recombinação dependentes da mutação ($\frac{r}{\mu}$) para cada um dos nove genes, e os valores observados foram bastante elevados. As taxas médias de recombinação dependentes da mutação aqui detectadas variaram de acordo com a região, sendo muito inferiores às taxas obtidas por MORRELL ET AL. (2003). A região que mostrou o maior valor médio de $\frac{r}{\mu}$ /Kb foi a *Dw2* ($1,578 \times 10^{-3}$), seguida das regiões *Ma3* ($1,209 \times 10^{-3}$), *Dw1* ($0,347 \times 10^{-3}$), *Ma1* ($0,213 \times 10^{-3}$) e *Dw4* ($0,129 \times 10^{-3}$). Claramente, os genes estudados por MORRELL ET AL. (2003) são bem menores que as regiões aqui investigadas, contendo, possivelmente, menor polimorfismo. Assim, como esses autores se basearam apenas nas taxas de recombinação dependentes da mutação, é possível que altas estimativas de recombinação tenham sido obtidas apenas em função de menores taxas de mutação.

As análises de fluxo gênico e recombinação não foram realizadas simultaneamente no presente estudo. Embora seja possível investigar tais processos em um único procedimento de análise no LAMARC (MORRELL ET AL., 2003), as estimações não são realizadas com base em um modelo conjunto. Isso porque fluxo gênico e recombinação foram propostos por modelos de coalescência independentes, que consideram apenas a taxa de mutação como processo comum. Portanto, é esperado que as estimativas obtidas em ambos os casos apresentem certo viés, algo que pode ser maior ou menor à medida que outros processos evolutivos se tornam importantes para explicar a variação dos dados.

Os padrões do DL via r^2 e $E[r_{\text{rhomap}}^2]$ mostraram resultados muito próximos para algumas regiões e populações. No entanto, a medida tradicional r^2 mostrou certa descontinuidade de possíveis blocos de associação, principalmente quando houve forte DL entre os locos. Tal comportamento pode sugerir, muitas vezes, interpretações pouco precisas da estrutura do DL dentro de certa região. HAMBLIN ET AL. (2005) também observaram padrões pouco conclusivos do DL em sorgo com base em tal medida. Esses resultados reforçam a ideia de que o r^2 é fortemente influenciado pelo tamanho da amostra e/ou pela baixa frequência alélica (FLINT-GARCIA ET AL., 2003), fornecendo estimativas variáveis e pouco precisas de associação (PRITCHARD and PRZEWSKI, 2001; NORDBORG and TAVARÉ, 2002). Acredita-se que o r^2 pode ser útil apenas para indicar que certo DL está presente em determinada região, fornecendo padrões pouco refinados de associação (EVANS and CARDON, 2005).

Em contrapartida, o DL obtido via coalescência sugeriu padrões que parecem fazer sentido para uma espécie autógama. Para muitos casos, grandes blocos de haplótipos, bem definidos e estruturados, foram observados, mostrando um extenso DL entre locos fisicamente ligados no genoma. Tais blocos não foram observados para todas as populações de certa região, mostrando que os padrões do DL foram específicos para cada população. Isso sugere que um único *pool* gênico pode ser insuficiente para a realização de mapeamento associativo em sorgo, algo que também foi sugerido em estudo anterior com humanos (EVANS and CARDON, 2005).

No entanto, a abordagem do $E[r^2]$ aqui utilizada (equação 2.30) foi proposta considerando que a amostra populacional é muito grande. Claramente, algumas das populações investigadas não apresentaram grande quantidade de indivíduos, possivelmente por questões de disponibilidade ou recursos de captação. Isso ocorreu

principalmente para as populações da AFS e de UGA. Então, embora as menores amostras estejam próximas às quantidades que são normalmente utilizadas em estudos de coalescência, o $E[r^2]$ pode ter sido inflacionado por questões de amostragem. Tal ponto é muito importante, e deve ser considerado em estudos futuros.

Além disso, embora o DL via coalescência tenha fornecido padrões importantes, apenas a recombinação foi considerada em seu cálculo. Em geral, todas as regiões tiveram elevado fluxo gênico entre as populações. O número médio de indivíduos que entraram e saíram de certa população, considerando o cenário complexo aqui investigado, deve ter influenciado nos padrões históricos do DL obtidos. Da mesma forma, o processo de seleção balanceadora, sugerido para a região *Mal*, também pode ter influenciado nos padrões elevados do DL. Portanto, é desejável que a interpretação genealógica do DL (McVEAN, 2002) seja estendida para contemplar vários processos evolutivos simultaneamente, tais como fluxo gênico, recombinação e seleção. Para uma espécie autógama tal como o sorgo, desvios do equilíbrio panmítico são esperados (HAMBLIN ET AL., 2005), e isso também deveria ser levado em consideração. Espera-se que uma estimativa mais precisa do DL via coalescência possa fornecer informações mais refinadas para estudos de mapeamento associativo em sorgo.

No momento em que grandes quantidades de dados estão sendo geradas a partir do genoma de várias espécies, métodos de coalescência aproximados podem ser utilizados para investigar diferentes cenários evolutivos. Embora o uso eficiente de tais métodos, como o ABC (*Approximate Bayesian Computation*) e o PAC (*Product of Approximate Conditionals*), seja algo ainda discutível, acredita-se que abordagens aproximadas estarão cada vez mais presentes no contexto de coalescência para múltiplos locos (McVEAN and CARDIN, 2005). Contudo, o grande desafio continuará sendo o desenvolvimento de abordagens completas mais eficientes (SETHURAMAN and HEY, 2016) e que consideram vários processos simultaneamente (SOUSA and HEY, 2013), na busca de um entendimento mais profundo dos processos evolutivos.

6 CONSIDERAÇÕES FINAIS

No presente estudo, a estrutura populacional e o desequilíbrio de ligação foram investigados com base na teoria da coalescência. Para tanto, análises de mutação, migração com fluxo gênico e recombinação foram realizadas para cinco regiões genômicas e sete populações de sorgo, e resultados muito interessantes foram encontrados.

A experiência de ter tido que compreender os princípios da teoria da coalescência e utilizar pacotes especificamente desenvolvidos para tal abordagem foi incrível. As análises de fluxo gênico e recombinação constante ao longo da região foram realizadas com o pacote LAMARC. Trata-se de um software muito bem documentado e com informações suficientes para o uso adequado. Contudo, detalhes importantes para a correta interpretação dos resultados não são muito bem esclarecidos ou são apresentados de forma isolada em um ou outro arquivo, dificultando o entendimento. Além disso, por utilizar abordagem de verossimilhança completa, tal pacote demanda enorme tempo computacional, podendo se tornar inviável para casos em que muitos dados encontram-se disponíveis. O programa `rhomap`, disponível no pacote `LDhat`, foi utilizado para análises de recombinação variável ao longo da região. Trata-se de um pacote com documentação necessária para o uso correto, porém com poucas explicações para iniciantes. Foi especificamente desenvolvido para o uso da teoria da coalescência a partir de muitos dados, podendo ser uma alternativa interessante para estimar recombinação em cenários mais atuais.

Em geral, elevado fluxo gênico foi observado entre todas as populações de sorgo a partir das análises de cada região genômica e todas elas simultaneamente. Além disso, os resultados sugeriram que a história evolutiva de regiões específicas do genoma desta espécie pode ser distinta, até mesmo quando processos que influenciam todo o genoma estão sendo estudados. Ainda, o fluxo gênico mostrou ser assimétrico entre todos os possíveis pares de populações, sugerindo que a forma pela qual as populações se relacionaram e continuam se interagindo evolutivamente não é igual. Este resultado evidencia que os métodos clássicos para investigar a estrutura populacional podem ser insatisfatórios para explicar o processo histórico-evolutivo.

Utilizando o método proposto por KUHNER *ET AL.* (2000), que considera taxas de recombinação constantes para certa região (LAMARC), baixas estimativas foram observadas para todas as regiões genômicas e populações. No entanto, utilizando o método proposto por AUTON and McVEAN (2007), que considera taxas de recombinação variáveis (`rhomap`), altas e baixas estimativas foram observadas de acordo com a região genômica e a população. A magnitude de recombinação para certas populações foi bastante elevada, sendo um resultado inesperado em sorgo.

Os padrões do DL via medida tradicional (r^2) e coalescência ($E[r_{\text{rhomap}}^2]$) mostraram resultados próximos em alguns casos. No entanto, o r^2 apresentou padrões descontínuos em várias ocasiões, dificultando o entendimento e a caracterização dos blocos de associação. Em contrapartida, o DL via coalescência ($E[r_{\text{rhomap}}^2]$) forneceu resultados que pareceram ter sido mais consistentes, sendo uma estratégia eventualmente importante para um refinamento dos padrões de associação. O DL de alta ordem via coalescência ($E[r_{\text{LAMARC}}^2]$) forneceu estimativas bastante elevadas para todas as regiões e populações, um resultado muito diferente dos padrões observados por r^2 e $E[r_{\text{rhomap}}^2]$.

Acredita-se que os resultados aqui encontrados possam fornecer informações relevantes para estudos futuros de mapeamento associativo em sorgo. Embora o contato entre as sete populações aqui estabelecidas tenha

vido historicamente intenso, os padrões do DL mostraram resultados específicos de acordo com a região genômica e a população. Tais resultados sugerem que o mapeamento genético a partir de um único *pool* gênico pode ser insuficiente para a detecção de associações causais importantes para características quantitativas em sorgo.

REFERÊNCIAS

- ABDURAKHMONOV, I. Y. and A. ABDUKARIMOV, 2008 Application of association mapping to understanding the genetic diversity of plant germplasm resources. *International journal of plant genomics* **2008**: 1–18, doi:10.1155/2008/574927.
- AGRAMA, H. and M. TUINSTRA, 2004 Phylogenetic diversity and relationships among sorghum accessions using SSRs and RAPDs. *African Journal of Biotechnology* **2**: 334–340.
- ALDRICH, P. and J. DOEBLEY, 1992 Restriction fragment variation in the nuclear and chloroplast genomes of cultivated and wild *Sorghum bicolor*. *Theoretical and Applied Genetics* **85**: 293–302.
- ALDRICH, P., J. DOEBLEY, K. SCHERTZ, and A. STEC, 1992 Patterns of allozyme variation in cultivated and wild *Sorghum bicolor*. *Theoretical and Applied Genetics* **85**: 451–460.
- ALLENDORE, F. W. and G. H. LUIKART, 2007 Conservation and the genetics of populations. Wiley. com.
- ARDLIE, K. G., L. KRUGLYAK, and M. SEIELSTAD, 2002 Patterns of linkage disequilibrium in the human genome. *Nature Reviews Genetics* **3**: 299–309.
- AUTON, A. and G. McVEAN, 2007 Recombination rate estimation in the presence of hotspots. *Genome research* **17**: 1219–1227.
- BAHLO, M. and R. GRIFFITHS, 2000 Inference from gene trees in a subdivided population. *Theoretical population biology* **57**: 79–95.
- BEAUMONT, M. A. and B. RANNALA, 2004 The Bayesian revolution in genetics. *Nature Reviews Genetics* **5**: 251–261.
- BECQUET, C. and M. PRZEWORSKI, 2007 A new approach to estimate parameters of speciation models with application to apes. *Genome research* **17**: 1505–1519.
- BEERLI, P., 1998 Estimation of migration rates and population sizes in geographically structured populations. *NATO ASI SERIES A LIFE SCIENCES* **306**: 39–54.
- BEERLI, P., 2006 Comparison of bayesian and maximum-likelihood inference of population genetic parameters. *Bioinformatics* **22**: 341–345.
- BEERLI, P. and J. FELSENSTEIN, 1999 Maximum-likelihood estimation of migration rates and effective population numbers in two populations using a coalescent approach. *Genetics* **152**: 763–773.
- BEERLI, P. and J. FELSENSTEIN, 2001 Maximum likelihood estimation of a migration matrix and effective population sizes in n subpopulations by using a coalescent approach. *Proceedings of the National Academy of Sciences* **98**: 4563–4568.
- BEERLI, P. and M. PALCZEWSKI, 2010 Unified framework to evaluate panmixia and migration direction among multiple sampling locations. *Genetics* **185**: 313–326.

- BEISSINGER, T. M., C. N. HIRSCH, R. S. SEKHON, J. M. FOERSTER, J. M. JOHNSON, G. MUTTONI, B. VAILLANCOURT, C. R. BUELL, S. M. KAEPLER, and N. DE LEON, 2013 Marker density and read depth for genotyping populations using genotyping-by-sequencing. *Genetics* **193**: 1073–1081.
- BHOSALE, S. U., B. STICH, H. F. W. RATTUNDE, E. WELTZIEN, B. I. HAUSSMANN, C. T. HASH, P. RAMU, H. E. CUEVAS, A. H. PATERSON, A. E. MELCHINGER, *ET AL.*, 2012 Association analysis of photoperiodic flowering time genes in west and central african sorghum [*Sorghum bicolor* (L.) Moench]. *BMC plant biology* **12**: 32.
- BOUCHET, S., D. POT, M. DEU, J.-F. RAMI, C. BILLOT, X. PERRIER, R. RIVALLAN, L. GARDES, L. XIA, P. WENZL, *ET AL.*, 2012 Genetic structure, linkage disequilibrium and signature of selection in sorghum: lessons from physically anchored dart markers. *PLoS One* **7**: e33470.
- BROWN, P. J., W. L. ROONEY, C. FRANKS, and S. KRESOVICH, 2008 Efficient mapping of plant height quantitative trait loci in a sorghum association population with introgressed dwarfing genes. *Genetics* **180**: 629–637.
- CANIATO, F. F., M. T. HAMBLIN, C. T. GUIMARAES, Z. ZHANG, R. E. SCHAFFERT, L. V. KOCHIAN, and J. V. MAGALHAES, 2014 Association mapping provides insights into the origin and the fine structure of the sorghum aluminum tolerance locus, *altsb*. *PloS One* **9**: e87438.
- CASA, A., S. MITCHELL, M. HAMBLIN, H. SUN, J. BOWERS, A. PATERSON, C. AQUADRO, and S. KRESOVICH, 2005 Diversity and selection in sorghum: simultaneous analyses using simple sequence repeats. *Theoretical and Applied Genetics* **111**: 23–30.
- CASA, A. M., S. E. MITCHELL, J. D. JENSEN, M. T. HAMBLIN, A. H. PATERSON, C. F. AQUADRO, and S. KRESOVICH, 2006 Evidence for a selective sweep on chromosome 1 of cultivated sorghum. *Crop science* **46**: S–27.
- CASA, A. M., G. PRESSOIR, P. J. BROWN, S. E. MITCHELL, W. L. ROONEY, M. R. TUINSTRAN, C. D. FRANKS, and S. KRESOVICH, 2008 Community resources and strategies for association mapping in sorghum. *Crop science* **48**: 30–40.
- CHANTEREAU, J., G. TROUCHE, J. RAMI, M. DEU, C. BARRO, and L. GRIVET, 2001 RFLP mapping of QTLs for photoperiod response in tropical sorghum. *Euphytica* **120**: 183–194.
- CHILDS, K. L., F. R. MILLER, M.-M. CORDONNIER-PRAATT, L. H. PRATT, P. W. MORGAN, and J. E. MULLET, 1997 The sorghum photoperiod sensitivity gene, *Ma3*, encodes a phytochrome B. *Plant Physiology* **113**: 611–619.
- CRASTA, O., W. XU, D. ROSENOW, J. MULLET, and H. NGUYEN, 1999 Mapping of post-flowering drought resistance traits in grain sorghum: association between QTLs influencing premature senescence and maturity. *Molecular and General Genetics MGG* **262**: 579–588.
- DAVEY, J. W., P. A. HOHENLOHE, P. D. ETTER, J. Q. BOONE, J. M. CATCHEN, and M. L. BLAXTER, 2011 Genome-wide genetic marker discovery and genotyping using next-generation sequencing. *Nature Reviews Genetics* **12**: 499–510.
- DE ALENCAR FIGUEIREDO, L., C. CALATAYUD, C. DUPUIITS, C. BILLOT, J.-F. RAMI, D. BRUNEL, X. PERRIER, B. COURTOIS, M. DEU, and J.-C. GLASZMANN, 2008 Phylogeographic evidence of crop neodiversity in sorghum. *Genetics* **179**: 997–1008.

- DE WET, J., 1978 Systematics and evolution of Sorghum sect. Sorghum (Gramineae). *American Journal of Botany* pp. 477–484.
- DE WET, J. M. and J. HARLAN, 1971 The origin and domestication of *Sorghum bicolor*. *Economic Botany* **25**: 128–135.
- DEU, M., F. RATTUNDE, and J. CHANTEREAU, 2006 A global view of genetic diversity in cultivated sorghums using a core collection. *Genome* **49**: 168–180.
- DEU, M., F. SAGNARD, J. CHANTEREAU, C. CALATAYUD, D. HÉRAULT, C. MARIAC, J.-L. PHAM, Y. VIGOUROUX, I. KAPRAN, P. TRAORE, ET AL., 2008 Niger-wide assessment of in situ sorghum genetic diversity with microsatellite markers. *Theoretical and Applied Genetics* **116**: 903–913.
- DOGGETT, H. ET AL., 1988 *Sorghum*. Number 2. ed., Longman Scientific and Technical.
- ELSHIRE, R. J., J. C. GLAUBITZ, Q. SUN, J. A. POLAND, K. KAWAMOTO, E. S. BUCKLER, and S. E. MITCHELL, 2011 A robust, simple genotyping-by-sequencing (GBS) approach for high diversity species. *PloS One* **6**: e19379.
- EVANS, D. M. and L. R. CARDON, 2005 A comparison of linkage disequilibrium patterns and estimated population recombination rates across multiple populations. *The American Journal of Human Genetics* **76**: 681–687.
- FAUBET, P., R. S. WAPLES, and O. E. GAGGIOTTI, 2007 Evaluating the performance of a multilocus bayesian method for the estimation of migration rates. *Molecular Ecology* **16**: 1149–1166.
- FEARNHEAD, P. and P. DONNELLY, 2001 Estimating recombination rates from population genetic data. *Genetics* **159**: 1299–1318.
- FEARNHEAD, P. and P. DONNELLY, 2002 Approximate likelihood methods for estimating local recombination rates. *Journal of the Royal Statistical Society: Series B (Statistical Methodology)* **64**: 657–680.
- FELSENSTEIN, J., 1981 Evolutionary trees from DNA sequences: a maximum likelihood approach. *Journal of molecular evolution* **17**: 368–376.
- FELSENSTEIN, J., 1988 Phylogenies from molecular sequences: inference and reliability. *Annual review of genetics* **22**: 521–565.
- FELSENSTEIN, J., 1992 Estimating effective population size from samples of sequences: inefficiency of pairwise and segregating sites as compared to phylogenetic estimates. *Genetical research* **59**: 139–147.
- FELSENSTEIN, J., 2006 Accuracy of coalescent likelihood estimates: do we need more sites, more sequences, or more loci? *Molecular biology and evolution* **23**: 691–700.
- FELSENSTEIN, J. and G. A. CHURCHILL, 1996 A hidden Markov model approach to variation among sites in rate of evolution. *Molecular Biology and Evolution* **13**: 93–104.
- FELSENSTEIN, J., M. K. KUHNER, J. YAMATO, and P. BEERLI, 1999 Likelihoods on coalescents: a Monte Carlo sampling approach to inferring parameters from population samples of molecular data. *Lecture Notes-Monograph Series* pp. 163–185.

- FELTUS, F., G. HART, K. SCHERTZ, A. CASA, S. KRESOVICH, S. ABRAHAM, P. KLEIN, P. BROWN, and A. PATERSON, 2006 Alignment of genetic maps and QTLs between inter-and intra-specific sorghum populations. *Theoretical and applied genetics* **112**: 1295–1305.
- FISHER, R., 1930 *The genetical theory of natural selection*.
- FISHER, R. A., 1935 The logic of inductive inference. *Journal of the Royal Statistical Society* pp. 39–82.
- FISHER, R. A. *ET AL.*, 1922 On the dominance ratio. *Proceedings of the royal society of Edinburgh* **42**: 321–341.
- FLINT-GARCIA, S. A., J. M. THORNSBERRY, and E. S. IV BUCKLER, 2003 Structure of linkage disequilibrium in plants*. *Annual Review of Plant Biology* **54**: 357–374.
- FOLKERTSMA, R. T., H. F. W. RATTUNDE, S. CHANDRA, G. S. RAJU, and C. T. HASH, 2005 The pattern of genetic diversity of guinea-race *Sorghum bicolor* (L.) Moench landraces as revealed with SSR markers. *Theoretical and Applied Genetics* **111**: 399–409.
- FRISSE, L., R. HUDSON, A. BARTOSZEWICZ, J. WALL, J. DONFACK, and A. DI RIENZO, 2001 Gene conversion and different population histories may explain the contrast between polymorphism and linkage disequilibrium levels. *The American Journal of Human Genetics* **69**: 831–843.
- FU, Y.-X., 1997 Coalescent theory for a partially selfing population. *Genetics* **146**: 1489–1499.
- FU, Y.-X., 1998 Probability of a segregating pattern in a sample of DNA sequences. *Theoretical population biology* **54**: 1–10.
- FU, Y.-X. and W.-H. LI, 1993 Statistical tests of neutrality of mutations. *Genetics* **133**: 693–709.
- FU, Y.-X. and W.-H. LI, 1999 Coalescing into the 21st century: an overview and prospects of coalescent theory. *Theoretical Population Biology* **56**: 1–10.
- FUTUYMA, D. J., 2005 *Evolution*. Sunderland, ma. USA: Sinauer Associates. Gouy MAM, Guindon S, Gascuel O (2010). SeaView version **4**: 221–224.
- GARBER, E. D. *ET AL.*, 1950 Cytotaxonomic studies in the genus *Sorghum*. *University of California Publications in Botany* **23**: 283–362.
- GEYER, C. J., 1991 Markov chain monte carlo maximum likelihood. Interface Foundation of North America. Retrieved from the University of Minnesota Digital Conservancy, <http://purl.umn.edu/58440>.
- GLAUBITZ, J. C., T. M. CASSTEVENS, F. LU, J. HARRIMAN, R. J. ELSHIRE, Q. SUN, and E. S. BUCKLER, 2014 TASSEL-GBS: a high capacity genotyping by sequencing analysis pipeline. *PLoS One* **9**: E90346.
- GREEN, P. J., 1995 Reversible jump Markov chain Monte Carlo computation and Bayesian model determination. *Biometrika* **82**: 711–732.
- GRIFFITHS, R. and S. TAVARÉ, 1994a Ancestral inference in population genetics. *Statistical Science* pp. 307–319.

- GRIFFITHS, R. and S. TAVARÉ, 1994b Simulating probability distributions in the coalescent. *Theoretical Population Biology* **46**: 131–159.
- GRIFFITHS, R. and S. TAVARÉ, 1995 Unrooted genealogical tree probabilities in the infinitely-many-sites model. *Mathematical biosciences* **127**: 77–98.
- GRIFFITHS, R. C., 1981 Neutral two-locus multiple allele models with recombination. *Theoretical Population Biology* **19**: 169–186.
- GRIFFITHS, R. C. and P. MARJORAM, 1996 Ancestral inference from samples of DNA sequences with recombination. *Journal of Computational Biology* **3**: 479–502.
- GRIFFITHS, R. C. and S. TAVARÉ, 1994c Sampling theory for neutral alleles in a varying environment. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences* **344**: 403–410.
- GUPTA, P. K., S. RUSTGI, and P. L. KULWAL, 2005 Linkage disequilibrium and association studies in higher plants: present status and future prospects. *Plant molecular biology* **57**: 461–485.
- HAHN, M. W., 2007 Detecting natural selection on cis-regulatory DNA. *Genetica* **129**: 7–18.
- HAHN, M. W., 2008 Toward a selection theory of molecular evolution. *Evolution* **62**: 255–265.
- HAMBLIN, M. T., A. M. CASA, H. SUN, S. C. MURRAY, A. H. PATERSON, C. F. AQUADRO, and S. KRESOVICH, 2006 Challenges of detecting directional selection after a bottleneck: lessons from *Sorghum bicolor*. *Genetics* **173**: 953–964.
- HAMBLIN, M. T., M. G. S. FERNANDEZ, A. M. CASA, S. E. MITCHELL, A. H. PATERSON, and S. KRESOVICH, 2005 Equilibrium processes cannot explain high levels of short-and medium-range linkage disequilibrium in the domesticated grass *Sorghum bicolor*. *Genetics* **171**: 1247–1256.
- HAMBLIN, M. T., S. E. MITCHELL, G. M. WHITE, J. GALLEGO, R. KUKATLA, R. A. WING, A. H. PATERSON, and S. KRESOVICH, 2004 Comparative population genetics of the panicoid grasses: sequence polymorphism, linkage disequilibrium and selection in a diverse sample of *Sorghum bicolor*. *Genetics* **167**: 471–483.
- HAMBLIN, M. T., M. G. SALAS FERNANDEZ, M. R. TUINSTRAN, W. L. ROONEY, and S. KRESOVICH, 2007 Sequence variation at candidate loci in the starch metabolism pathway in sorghum: prospects for linkage disequilibrium mapping. *Crop Science* **47**: S–125.
- HAMILTON, M., 2009 Population genetics. Wiley. com.
- HARLAN, J. and J. DE WET, 1972a A simplified classification of cultivated sorghum. *Crop Science* **12**: 172–176.
- HARLAN, J. and J. DE WET, 1972b A simplified classification of cultivated sorghum. *Crop Science* **12**: 172–176.
- HARRIS, K., P. SUBUDHI, A. BORRELL, D. JORDAN, D. ROSENOW, H. NGUYEN, P. KLEIN, R. KLEIN, and J. MULLET, 2007 Sorghum stay-green QTL individually reduce post-flowering drought-induced leaf senescence. *Journal of experimental botany* **58**: 327–338.

- HART, G., K. SCHERTZ, Y. PENG, and N. SYED, 2001 Genetic mapping of *Sorghum bicolor* (L.) Moench QTLs that control variation in tillering and other morphological characters. *Theoretical and Applied Genetics* **103**: 1232–1242.
- HARTL, D. L. and A. G. CLARK, 2007 *Principles of population genetics*, volume 116. Sinauer associates Sunderland.
- HASTINGS, W. K., 1970 Monte Carlo sampling methods using Markov chains and their applications. *Biometrika* **57**: 97–109.
- HEDRICK, P. W., 2010 *Genetics of populations*. Jones & Bartlett Publishers.
- HEIN, J., M. SCHIERUP, and C. WIUF, 2004 *Gene genealogies, variation and evolution: a primer in coalescent theory*. Oxford university press.
- HEY, J., 2010 Isolation with migration models for more than two populations. *Molecular biology and evolution* **27**: 905–920.
- HEY, J. and R. NIELSEN, 2004 Multilocus methods for estimating population sizes, migration rates and divergence time, with applications to the divergence of *Drosophila pseudoobscura* and *D. persimilis*. *Genetics* **167**: 747–760.
- HEY, J. and R. NIELSEN, 2007 Integration within the Felsenstein equation for improved Markov chain Monte Carlo methods in population genetics. *Proceedings of the National Academy of Sciences* **104**: 2785–2790.
- HILL, W. and A. ROBERTSON, 1968 Linkage disequilibrium in finite populations. *Theoretical and Applied Genetics* **38**: 226–231.
- HOLSINGER, K. E. and B. S. WEIR, 2009 Genetics in geographically structured populations: defining, estimating and interpreting F_{ST} . *Nature Reviews Genetics* **10**: 639–650.
- HÜBNER, S., T. GÜNTHER, A. FLAVELL, E. FRIDMAN, A. GRANER, A. KOROL, and K. J. SCHMID, 2012 Islands and streams: clusters and gene flow in wild barley populations from the levant. *Molecular ecology* **21**: 1115–1129.
- HUDSON, R. R., 1983 Properties of a neutral allele model with intragenic recombination. *Theoretical population biology* **23**: 183–201.
- HUDSON, R. R., 1985 The sampling distribution of linkage disequilibrium under an infinite allele model without selection. *Genetics* **109**: 611–631.
- HUDSON, R. R., 1987 Estimating the recombination parameter of a finite population model without selection. *Genetical research* **50**: 245–250.
- HUDSON, R. R., 1990 Gene genealogies and the coalescent process. *Oxford surveys in evolutionary biology* **7**: 44.
- HUDSON, R. R., 2001 Two-locus sampling distributions and their application. *Genetics* **159**: 1805–1817.
- JORDAN, D., E. S. MACE, R. HENZELL, P. KLEIN, and R. KLEIN, 2010 Molecular mapping and candidate gene identification of the Rf2 gene for pollen fertility restoration in sorghum [*Sorghum bicolor* (L.) Moench]. *Theoretical and applied genetics* **120**: 1279–1287.

- KAPLAN, N. and R. R. HUDSON, 1985 The use of sample genealogies for studying a selectively neutral m-loci model with recombination. *Theoretical population biology* **28**: 382–396.
- KARLIN, S. and J. MCGREGOR, 1967 The number of mutant forms maintained in a population. In Proceedings of the Fifth Berkeley Symposium on mathematics, Statistics and probability, volume 4, pp. 415–438.
- KARPER, R., 1932 A dominant mutation of frequent recurrence in sorghum. *American Naturalist* pp. 511–529.
- KEBEDE, H., P. SUBUDHI, D. ROSENOW, and H. NGUYEN, 2001 Quantitative trait loci influencing drought tolerance in grain sorghum (*Sorghum bicolor* L. Moench). *Theoretical and Applied Genetics* **103**: 266–276.
- KIM, S., V. PLAGNOL, T. T. HU, C. TOOMAJIAN, R. M. CLARK, S. OSSOWSKI, J. R. ECKER, D. WEIGEL, and M. NORDBORG, 2007 Recombination and linkage disequilibrium in *Arabidopsis thaliana*. *Nature genetics* **39**: 1151–1155.
- KIM, Y., S. FENG, and Z.-B. ZENG, 2008 Measuring and partitioning the high-order linkage disequilibrium by multiple order Markov chains. *Genetic epidemiology* **32**: 301–312.
- KIMMEL, G., R. M. KARP, M. I. JORDAN, and E. HALPERIN, 2008 Association mapping and significance estimation via the coalescent. *The American Journal of Human Genetics* **83**: 675–683.
- KIMURA, M., 1955 Solution of a process of random genetic drift with a continuous model. *Proceedings of the National Academy of Sciences of the United States of America* **41**: 144.
- KIMURA, M., 1969 The number of heterozygous nucleotide sites maintained in a finite population due to steady flux of mutations. *Genetics* **61**: 893.
- KINGMAN, J., 1982a Exchangeability and the evolution of large populations.
- KINGMAN, J. F., 1982b The coalescent. *Stochastic processes and their applications* **13**: 235–248.
- KINGMAN, J. F., 1982c On the genealogy of large populations. *Journal of Applied Probability* pp. 27–43.
- KINGMAN, J. F., 2000 Origins of the coalescent: 1974–1982. *Genetics* **156**: 1461–1463.
- KLEIN, R., P. KLEIN, A. CHHABRA, J. DONG, S. PAMMI, K. CHILDS, J. MULLET, W. ROONEY, and K. SCHERTZ, 2001 Molecular mapping of the Rf1 gene for pollen fertility restoration in sorghum (*Sorghum bicolor* L.). *Theoretical and applied genetics* **102**: 1206–1212.
- KLEIN, R. R., J. E. MULLET, D. R. JORDAN, F. R. MILLER, W. L. ROONEY, M. A. MENZ, C. D. FRANKS, and P. E. KLEIN, 2008 The effect of tropical sorghum conversion and inbred development on genome diversity as revealed by high-resolution genotyping. *Crop science* **48**: S–12.
- KORNELIUSSEN, T. S., I. MOLTKE, A. ALBRECHTSEN, and R. NIELSEN, 2013 Calculation of Tajimas D and other neutrality test statistics from low depth next-generation sequencing data. *BMC bioinformatics* **14**: 289.
- KREITMAN, M., 2000 Methods to detect selection in populations with applications to the human. *Annual review of genomics and human genetics* **1**: 539–559.

- KUHNER, M. K., 2006 LAMARC 2.0: maximum likelihood and Bayesian estimation of population parameters. *Bioinformatics* **22**: 768–770.
- KUHNER, M. K., 2009 Coalescent genealogy samplers: windows into population history. *Trends in Ecology & Evolution* **24**: 86–93.
- KUHNER, M. K. and L. P. SMITH, 2007 Comparing likelihood and bayesian coalescent estimation of population parameters. *Genetics* **175**: 155–165.
- KUHNER, M. K., J. YAMATO, and J. FELSENSTEIN, 1995 Estimating effective population size and mutation rate from sequence data using Metropolis-Hastings sampling. *Genetics* **140**: 1421–1430.
- KUHNER, M. K., J. YAMATO, and J. FELSENSTEIN, 1998 Maximum likelihood estimation of population growth rates based on the coalescent. *Genetics* **149**: 429–434.
- KUHNER, M. K., J. YAMATO, and J. FELSENSTEIN, 2000 Maximum likelihood estimation of recombination rates from population data. *Genetics* **156**: 1393–1401.
- LEWONTIN, R., 1964 The interaction of selection and linkage. i. general considerations; heterotic models. *Genetics* **49**: 49.
- LEWONTIN, R. and K.-I. KOJIMA, 1960 The evolutionary dynamics of complex polymorphisms. *Evolution* pp. 458–472.
- LI, H., B. HANDSAKER, A. WYSOKER, T. FENNEL, J. RUAN, N. HOMER, G. MARTH, G. ABECASIS, R. DURBIN, ET AL., 2009 The sequence alignment/map format and SAMtools. *Bioinformatics* **25**: 2078–2079.
- LI, N. and M. STEPHENS, 2003 Modeling linkage disequilibrium and identifying recombination hotspots using single-nucleotide polymorphism data. *Genetics* **165**: 2213–2233.
- LIN, Y.-R., K. F. SCHERTZ, and A. H. PATERSON, 1995 Comparative analysis of QTLs affecting plant height and maturity across the poaceae, in reference to an interspecific sorghum population. *Genetics* **141**: 391.
- LOHSE, K., R. HARRISON, and N. H. BARTON, 2011 A general method for calculating likelihoods under the coalescent process. *Genetics* **189**: 977–987.
- MACE, E. and D. JORDAN, 2010 Location of major effect genes in sorghum (*Sorghum bicolor* (L.) Moench). *Theoretical and applied genetics* **121**: 1339–1356.
- MAGALHAES, J. V., D. F. GARVIN, Y. WANG, M. E. SORRELLS, P. E. KLEIN, R. E. SCHAFFERT, L. LI, and L. V. KOCHIAN, 2004 Comparative mapping of a major aluminum tolerance gene in sorghum and other species in the poaceae. *Genetics* **167**: 1905–1914.
- MAILUND, T., A. E. HALAGER, M. WESTERGAARD, J. Y. DUTHEIL, K. MUNCH, L. N. ANDERSEN, G. LUNTER, K. PRÜFER, A. SCALLY, A. HOBOLTH, ET AL., 2012 A new isolation with migration model along complete genomes infers very different divergence processes among closely related great ape species. *PLoS Genet* **8**: e1003125.

- MARCHINI, J. and B. HOWIE, 2010 Genotype imputation for genome-wide association studies. *Nature Reviews Genetics* **11**: 499–511.
- MARJORAM, P. and P. JOYCE, 2011 Practical implications of coalescent theory. *Problem Solving Handbook in Computational Biology and Bioinformatics* **1**: 63.
- MARJORAM, P. and S. TAVARÉ, 2006 Modern computational approaches for analysing molecular genetic variation data. *Nature Reviews Genetics* **7**: 759–770.
- McVEAN, G., 2007a Linkage disequilibrium, recombination and selection. *Handbook of Statistical Genetics*, Third Edition pp. 909–944.
- McVEAN, G., 2007b The structure of linkage disequilibrium around a selective sweep. *Genetics* **175**: 1395–1406.
- McVEAN, G., P. AWADALLA, and P. FEARNHEAD, 2002 A coalescent-based method for detecting and estimating recombination from gene sequences. *Genetics* **160**: 1231–1241.
- McVEAN, G. A., 2002 A genealogical interpretation of linkage disequilibrium. *Genetics* **162**: 987–991.
- McVEAN, G. A. and N. J. CARDIN, 2005 Approximating the coalescent with recombination. *Philosophical Transactions of the Royal Society B: Biological Sciences* **360**: 1387–1393.
- McVEAN, G. A., S. R. MYERS, S. HUNT, P. DELOUKAS, D. R. BENTLEY, and P. DONNELLY, 2004 The fine-scale structure of recombination rate variation in the human genome. *Science* **304**: 581–584.
- METROPOLIS, N., A. W. ROSENBLUTH, M. N. ROSENBLUTH, A. H. TELLER, and E. TELLER, 1953 Equation of state calculations by fast computing machines. *The journal of chemical physics* **21**: 1087.
- MORDEN, C., J. DOEBLEY, and K. SCHERTZ, 1990 Allozyme variation among the spontaneous species of *Sorghum* section *Sorghum* (Poaceae). *Theoretical and Applied Genetics* **80**: 296–304.
- MORDEN, C. W., J. F. DOEBLEY, and K. F. SCHERTZ, 1989 Allozyme variation in old world races of *Sorghum bicolor* (Poaceae). *American journal of botany* pp. 247–255.
- MORRELL, P. L., K. E. LUNDY, and M. T. CLEGG, 2003 Distinct geographic patterns of genetic diversity are maintained in wild barley (*Hordeum vulgare* ssp. *spontaneum*) despite migration. *Proceedings of the National Academy of Sciences* **100**: 10812–10817.
- MORRIS, G. P., P. RAMU, S. P. DESHPANDE, C. T. HASH, T. SHAH, H. D. UPADHYAYA, O. RIERA-LIZARAZU, P. J. BROWN, C. B. ACHARYA, S. E. MITCHELL, ET AL., 2013 Population genomic and genome-wide association studies of agroclimatic traits in sorghum. *Proceedings of the National Academy of Sciences* **110**: 453–458.
- MULTANI, D. S., S. P. BRIGGS, M. A. CHAMBERLIN, J. J. BLAKESLEE, A. S. MURPHY, and G. S. JOHAL, 2003 Loss of an MDR transporter in compact stalks of maize br2 and sorghum dw3 mutants. *Science* **302**: 81–84.
- MURPHY, R. L., R. R. KLEIN, D. T. MORISHIGE, J. A. BRADY, W. L. ROONEY, F. R. MILLER, D. V. DUGAS, P. E. KLEIN, and J. E. MULLET, 2011 Coincident light and clock regulation of pseudoresponse regulator protein 37 (PRR37) controls photoperiodic flowering in sorghum. *Proceedings of the National Academy of Sciences* **108**: 16469–16474.

- MURPHY, R. L., D. T. MORISHIGE, J. A. BRADY, W. L. ROONEY, S. YANG, P. E. KLEIN, and J. E. MULLET, 2014 Ghd7 (Ma6) represses sorghum flowering in long days: Alleles enhance biomass accumulation and grain production. *The Plant Genome* **7**.
- MUTEGI, E., F. SAGNARD, K. SEMAGN, M. DEU, M. MURAYA, B. KANYENJI, S. DE VILLIERS, D. KIAMBI, L. HERSELMAN, and M. LABUSCHAGNE, 2011 Genetic structure and relationships within and between cultivated and wild sorghum (*Sorghum bicolor* (L.) Moench) in kenya as revealed by microsatellite markers. *Theoretical and applied genetics* **122**: 989–1004.
- MYUNG, I. J., 2003 Tutorial on maximum likelihood estimation. *Journal of mathematical Psychology* **47**: 90–100.
- NATH, H. and R. GRIFFITHS, 1993 The coalescent in two colonies with symmetric migration. *Journal of mathematical biology* **31**: 841–851.
- NATH, H. and R. GRIFFITHS, 1996 Estimation in an island model using simulation. *Theoretical population biology* **50**: 227–253.
- NIELSEN, R., 2001 Statistical tests of selective neutrality in the age of genomics. *Heredity* **86**: 641–647.
- NIELSEN, R. and M. SLATKIN, 2013 *An introduction to population genetics: theory and applications*. Sinauer Associates Sunderland, MA.
- NIELSEN, R. and J. WAKELEY, 2001 Distinguishing migration from isolation: a Markov chain Monte Carlo approach. *Genetics* **158**: 885–896.
- NORDBORG, M., 2007 Coalescent theory. *Hand of statistical genetics*.
- NORDBORG, M. and P. DONNELLY, 1997 The coalescent process with selfing. *Genetics* **146**: 1185–1195.
- NORDBORG, M., T. T. HU, Y. ISHINO, J. JHAVERI, C. TOOMAJIAN, H. ZHENG, E. BAKKER, P. CALABRESE, J. GLADSTONE, R. GOYAL, ET AL., 2005 The pattern of polymorphism in *Arabidopsis thaliana*. *PLoS biology* **3**: e196.
- NORDBORG, M. and S. TAVARÉ, 2002 Linkage disequilibrium: what history has to tell us. *TRENDS in Genetics* **18**: 83–90.
- NOTOHARA, M., 1990 The coalescent and the genealogical process in geographically structured population. *Journal of mathematical biology* **29**: 59–75.
- ODONG, T., J. VAN HEERWAARDEN, J. JANSEN, T. J. VAN HINTUM, and F. VAN EEUWIJK, 2011 Determination of genetic structure of germplasm collections: are traditional hierarchical clustering methods appropriate for molecular marker data? *Theoretical and Applied Genetics* **123**: 195–205.
- OHTA, T. and M. KIMURA, 1969a Linkage disequilibrium at steady state determined by random genetic drift and recurrent mutation. *Genetics* **63**: 229.
- OHTA, T. and M. KIMURA, 1969b Linkage disequilibrium due to random genetic drift. *Genetical Research* **13**: 47–55.
- OHTA, T. and M. KIMURA, 1971 Linkage disequilibrium between two segregating nucleotide sites under the steady flux of mutations in a finite population. *Genetics* **68**: 571.

- ORAGUZIE, N. C., E. H. RIKKERINK, S. E. GARDINER, H. D. SILVA, ET AL., 2007 Association mapping in plants. Springer-Verlag GmbH.
- PALCZEWSKI, M. and P. BEERLI, 2013 A continuous method for gene flow. *Genetics* **194**: 687–696.
- PARADIS, E., 2010 pegas: an R package for population genetics with an integrated–modular approach. *Bioinformatics* **26**: 419–420.
- PATERSON, A. H., J. E. BOWERS, R. BRUGGMANN, I. DUBCHAK, J. GRIMWOOD, H. GUNDLACH, G. HABERER, U. HELLSTEN, T. MITROS, A. POLIAKOV, ET AL., 2009 The *Sorghum bicolor* genome and the diversification of grasses. *Nature* **457**: 551–556.
- PEARSE, D. E. and K. A. CRANDALL, 2004 Beyond F_{ST} : analysis of population genetic data for conservation. *Conservation Genetics* **5**: 585–602.
- PEREIRA, M. and M. LEE, 1995 Identification of genomic regions affecting plant height in sorghum and maize. *Theoretical and Applied Genetics* **90**: 380–388.
- POLAND, J. A., P. J. BROWN, M. E. SORRELLS, and J.-L. JANNINK, 2012 Development of high-density genetic maps for barley and wheat using a novel two-enzyme genotyping-by-sequencing approach. *PloS one* **7**: e32253.
- POLAND, J. A. and T. W. RIFE, 2012 Genotyping-by-sequencing for plant breeding and genetics. *The Plant Genome* **5**: 92–102.
- PRITCHARD, J. K. and M. PRZEWORSKI, 2001 Linkage disequilibrium in humans: models and data. *The American Journal of Human Genetics* **69**: 1–14.
- PRITCHARD, J. K., M. STEPHENS, N. A. ROSENBERG, and P. DONNELLY, 2000 Association mapping in structured populations. *The American Journal of Human Genetics* **67**: 170–181.
- PTAK, S. E., K. VOELPEL, and M. PRZEWORSKI, 2004 Insights into recombination from patterns of linkage disequilibrium in humans. *Genetics* **167**: 387–397.
- QUINBY, J., 1966 Fourth maturity gene locus in sorghum. *Crop Science* **6**: 516–518.
- QUINBY, J., 1975 The genetics of sorghum improvement. *Journal of Heredity* **66**: 56–62.
- QUINBY, J., J. HESKETH, and R. VOIGT, 1973 Influence of temperature and photoperiod on floral initiation and leaf number in sorghum. *Crop Science* **13**: 243–246.
- QUINBY, J. and R. KARPER, 1945 Inheritance of three genes that influence time of floral initiation and maturity date in milo. *Journal of the American Society of Agronomy*.
- QUINBY, J. and R. KARPER, 1954 Inheritance of height in sorghum.
- QUINBY, J. R., 1974 Sorghum improvement and the genetics of growth.
- R CORE TEAM, 2015 R: A Language and Environment for Statistical Computing.

- RAFALSKI, A. and M. MORGANTE, 2004 Corn and humans: recombination and linkage disequilibrium in two genomes of similar size. *TRENDS in Genetics* **20**: 103–111.
- RAMÍREZ-SORIANO, A., S. E. RAMOS-ONSINS, J. ROZAS, F. CALAFELL, and A. NAVARRO, 2008 Statistical power analysis of neutrality tests under demographic expansions, contractions and bottlenecks with recombination. *Genetics* **179**: 555–567.
- RAMOS-ONSINS, S. E. and J. ROZAS, 2002 Statistical properties of new neutrality tests against population growth. *Molecular biology and evolution* **19**: 2092–2100.
- ROBERTS, A., L. McMILLAN, W. WANG, J. PARKER, I. RUSYN, and D. THREADGILL, 2007 Inferring missing genotypes in large SNP panels using fast nearest-neighbor searches over sliding windows. *Bioinformatics* **23**: i401–i407.
- ROONEY, W. and C. W. SMITH, 2000 Techniques for developing new cultivars. *Sorghum: Origin, history, technology and production*. John Wiley & Sons, New York pp. 329–347.
- ROONEY, W. L. and S. AYDIN, 1999 Genetic control of a photoperiod-sensitive response in *Sorghum bicolor* (L.) Moench. *Crop science* **39**: 397–400.
- ROSA, J. R. B. F., 2011 Análise do desequilíbrio de ligação e da estrutura populacional do germoplasma brasileiro de cana-de-açúcar. Ph.D. thesis, Escola Superior de Agricultura Luiz de Queiroz, Universidade de São Paulo.
- ROSENBERG, N. A. and M. NORDBORG, 2002 Genealogical trees, coalescent theory and the analysis of genetic polymorphisms. *Nature Reviews Genetics* **3**: 380–390.
- ROSENOW, D., J. DAHLBERG, J. STEPHENS, F. MILLER, D. BARNES, G. PETERSON, J. JOHNSON, and K. SCHERTZ, 1997 Registration of 63 converted sorghum germplasm lines from the sorghum conversion program. *Crop Sci* **37**: 1399–1400.
- ROSENOW, D., J. QUISENBERRY, C. WENDT, and L. CLARK, 1983 Drought tolerant sorghum and cotton germplasm. *Agricultural Water Management* **7**: 207–222.
- SAGNARD, F., M. DEU, D. DEMBÉLÉ, R. LEBLOIS, L. TOURÉ, M. DIAKITÉ, C. CALATAYUD, M. VAKSMANN, S. BOUCHET, Y. MALLÉ, ET AL., 2011 Genetic diversity, structure, gene flow and evolutionary relationships within the *Sorghum bicolor* wild–weedy–crop complex in a western african region. *Theoretical and applied genetics* **123**: 1231–1246.
- SANSALONI, C., 2012 Desenvolvimento e aplicações de DArT (Diversity Arrays Technology) e genotipagem por sequenciamento (Genotyping-By-Sequencing) para análise genética em Eucalyptus, volume 1. Tese (Doutorado). Universidade de Brasília.
- SANTOS, F., C. CASELA, and J. WAQUIL, 2005 Melhoramento de sorgo, volume 2. UFV: Viçosa.
- SETHURAMAN, A. and J. HEY, 2016 IMa2p–parallel MCMC and inference of ancient demography under the isolation with migration (IM) model. *Molecular ecology resources* **16**: 206–215.
- SIMONSEN, K. L., G. A. CHURCHILL, and C. F. AQUADRO, 1995 Properties of statistical tests of neutrality for DNA polymorphism data. *Genetics* **141**: 413–429.

- SINGH, H. P. and H. LOHITHASWA, 2006 Sorghum. In *Cereals and Millets*, pp. 257–302, Springer.
- SLATKIN, M., 2008 Linkage disequilibrium—understanding the evolutionary past and mapping the medical future. *Nature Reviews Genetics* **9**: 477–485.
- SLATKIN, M., L. EXCOFFIER, ET AL., 1996 Testing for linkage disequilibrium in genotypic data using the expectation-maximization algorithm. *Heredity* **76**: 377–383.
- SMITH, C. W. and R. A. FREDERIKSEN, 2000 *Sorghum: Origin, history, technology, and production*, volume 2. John Wiley & Sons.
- SNOWDEN, J., 1936 *Cultivated Races of Sorghum..* London; Trustees of Bentham-Moxon Fund.
- SOUSA, V. and J. HEY, 2013 Understanding the origin of species with genome-scale data: modelling gene flow. *Nature Reviews Genetics* **14**: 404–414.
- SOUSA, V. C., M. CARNEIRO, N. FERRAND, and J. HEY, 2013 Identifying loci under selection against gene flow in isolation-with-migration models. *Genetics* **194**: 211–233.
- STACKLIES, W., H. REDESTIG, M. SCHOLZ, D. WALTHER, and J. SELBIG, 2007 *pcaMethods*—a bioconductor package providing pca methods for incomplete data. *Bioinformatics* **23**: 1164–1167.
- STEPHENS, J., F. MILLER, and D. ROSENOW, 1967 Conversion of alien sorghums to early combine genotypes. *Crop Science* **7**: 396–396.
- STEPHENS, M., 2007 Inference under the coalescent. *Handbook of statistical genetics*.
- STUMPE, M. P. and G. A. McVEAN, 2003 Estimating recombination rates from population-genetic data. *Nature Reviews Genetics* **4**: 959–968.
- SUBUDHI, P., D. ROSENOW, and H. NGUYEN, 2000 Quantitative trait loci for the stay green trait in sorghum (*Sorghum bicolor* L. Moench): consistency across genetic backgrounds and environments. *Theoretical and Applied Genetics* **101**: 733–741.
- SUKUMARAN, S., W. XIANG, S. R. BEAN, J. F. PEDERSEN, S. KRESOVICH, M. R. TUINSTRAN, T. T. TESSO, M. T. HAMBLIN, and J. YU, 2012 Association mapping for grain quality in a diverse sorghum collection. *The Plant Genome* **5**: 126–135.
- SVED, J., 1968 The stability of linked systems of loci with a small population size. *Genetics* **59**: 543.
- SVED, J., 1971 Linkage disequilibrium and homozygosity of chromosome segments in finite populations. *Theoretical population biology* **2**: 125–141.
- TAJIMA, F., 1983 Evolutionary relationship of DNA sequences in finite populations. *Genetics* **105**: 437–460.
- TAJIMA, F., 1989 Statistical method for testing the neutral mutation hypothesis by DNA polymorphism. *Genetics* **123**: 585–595.
- TAKAHATA, N., 1988 The coalescent in two partially isolated diffusion populations. *Genetical research* **52**: 213–222.

- TAKAHATA, N. and M. SLATKIN, 1990 Genealogy of neutral genes in two partially isolated populations. *Theoretical Population Biology* **38**: 331–350.
- TAO, Y., A. HARDY, J. DRENTH, R. HENZELL, B. FRANZMANN, D. JORDAN, D. BUTLER, and C. McINTYRE, 2003 Identifications of two different mechanisms for sorghum midge resistance through QTL mapping. *Theoretical and Applied Genetics* **107**: 116–122.
- TAO, Y., R. HENZELL, D. JORDAN, D. BUTLER, A. KELLY, and C. McINTYRE, 2000 Identification of genomic regions associated with stay green in sorghum by testing RILs in multiple environments. *Theoretical and Applied Genetics* **100**: 1225–1232.
- TEMPLETON, A. R., 2006 *Population genetics and microevolutionary theory*. John Wiley & Sons.
- TUINSTRAN, M., E. GROTE, P. GOLDSBROUGH, and G. EJETA, 1996 Identification of quantitative trait loci associated with pre-flowering drought tolerance in sorghum. *Crop Science* **36**: 1337–1344.
- VARSHNEY, R. K., S. N. NAYAK, G. D. MAY, and S. A. JACKSON, 2009 Next-generation sequencing technologies and their implications for crop genetics and breeding. *Trends in biotechnology* **27**: 522–530.
- WAKELEY, J., 1997 Using the variance of pairwise differences to estimate the recombination rate. *Genetical research* **69**: 45–48.
- WAKELEY, J., 2009 *Coalescent theory: an introduction*, volume 1. Roberts & Company Publishers.
- WALL, J. D., 2000 A comparison of estimators of the population recombination rate. *Molecular Biology and Evolution* **17**: 156–163.
- WALL, J. D., 2004 Estimating recombination rates using three-site likelihoods. *Genetics* **167**: 1461–1473.
- WATTERSON, G., 1975 On the number of segregating sites in genetical models without recombination. *Theoretical population biology* **7**: 256–276.
- WEIR, B. and W. HILL, 1986 Nonuniform recombination within the human beta-globin gene cluster. *American journal of human genetics* **38**: 776.
- WEIR, B. S. ET AL., 1990 *Genetic data analysis. Methods for discrete population genetic data.*. Sinauer Associates, Inc. Publishers.
- WILSON, G. A. and B. RANNALA, 2003 Bayesian inference of recent migration rates using multilocus genotypes. *Genetics* **163**: 1177–1191.
- WIUF, C. and J. HEIN, 1999 Recombination as a point process along sequences. *Theoretical population biology* **55**: 248–259.
- WRIGHT, S., 1931 Evolution in mendelian populations. *Genetics* **16**: 97–159.
- XU, W., P. K. SUBUDHI, O. R. CRASTA, D. T. ROSENOW, J. E. MULLET, and H. T. NGUYEN, 2000 Molecular mapping of QTLs conferring stay-green in grain sorghum (*Sorghum bicolor* L. Moench). *Genome* **43**: 461–469.

- YU, J. and E. S. BUCKLER, 2006 Genetic association mapping and genome organization of maize. *Current Opinion in Biotechnology* **17**: 155–160.
- ZHU, C., M. GORE, E. S. BUCKLER, and J. YU, 2008 Status and prospects of association mapping in plants. *The plant genome* **1**: 5–20.
- ZÖLLNER, S. and J. K. PRITCHARD, 2005 Coalescent-based association mapping and fine mapping of complex trait loci. *Genetics* **169**: 1071–1092.
- ZÖLLNER, S. and A. VON HAESLER, 2000 A coalescent approach to study linkage disequilibrium between single-nucleotide polymorphisms. *The American Journal of Human Genetics* **66**: 615–628.