



Universidade Federal de Pernambuco
Centro de Ciências Biológicas
Departamento de Genética
Programa de Pós-Graduação em Genética

Bruna Piereck Moura

Análise comparativa de transposons classe II (*CACTA* e *Mutator*) nos genomas de *Vigna unguiculada*, *Phaseolus vulgaris* e *Medicago truncatula*

Recife - PE

2015

BRUNA PIERECK MOURA

Análise comparativa de transposons classe II (*CACTA* e *Mutator*) nos genomas de *Vigna unguiculada*, *Phaseolus vulgaris* e *Medicago truncatula*

Dissertação apresentada ao Programa de Pós-Graduação em Genética da Universidade Federal de Pernambuco como parte dos requisitos exigidos para obtenção do título de Mestre em Genética.

Orientadora: Ana Christina Brasileiro Vidal

Recife - PE

2015

Catálogo na fonte
Elaine Barroso
CRB 1728

Moura, Bruna Piereck

Análise comparativa de transposons classe II (CACTA e Mutator) nos genomas de *Vigna unguiculada*, *Phaseolus vulgaris* e *Medicago truncatula* / Bruna Piereck Moura- Recife: O Autor, 2015.

96 folhas: il., fig., tab.

Orientadora: Ana Christina Brasileiro Vidal

Coorientadora: Ana Maria Benko Iseppon

Dissertação (mestrado) – Universidade Federal de Pernambuco.

Centro de Biociências. Genética, 2015.

Inclui referências

1. Genética vegetal 2. Feijão 3. Bioinformática I. Pontes Filho, Nicodemos Teles de (orientador) II. Título

581.35

CDD (22.ed.)

UFPE/CB-2017-275

BRUNA PIERECK MOURA

Análise comparativa de transposons classe II (*CACTA* e *Mutator*) nos genomas de *Vigna unguiculata*, *Phaseolus vulgaris* e *Medicago truncatula*

Dissertação apresentada ao Programa de Pós-Graduação em Genética, Área de Concentração Genética, da Universidade Federal de Pernambuco, como requisito parcial para obtenção do título de Mestre em Genética.

Aprovado em: 06/03/2015

BANCA EXAMINADORA

Dra. Ana Christina Brasileiro Vidal
Universidade Federal de Pernambuco

Dr. Tercilio Calsa junior
Universidade Federal de Pernambuco

Dr. Éderson Akio Kido
Universidade Federal de Pernambuco

Dra. Katia Castanho Scortecci
Universidade Federal do Rio Grande do Norte

AGRADECIMENTOS

Eu, Bruna, gostaria de começar agradecendo imensamente as minhas orientadoras uma dupla de ANAs, Ana Christina Brasileiro Vidal e Ana Maria Benko Iseppon que fizeram o possível o meu sonho de entrar no mestrado juntamente com o PPGG e o CNPq. Ao entrar no LGBV recebi como objeto de estudo os elementos transponíveis, sem nunca ter ouvido falar deles até então, acabei por me apaixonar por eles com toda sua complexidade, me vi fascinada e até um pouco assustada. No entanto agora, como um vício, apesar das dificuldades e estresses, não pretendo largá-los e me vejo embebida na vontade de sempre aprofundar-me mais no conhecimento sobre eles. Graças a essas duas mães científicas que me adotaram e com ajuda delas foi possível caminhar nesse mundo chamado Pós-Graduação e aprender muito, especialmente que sempre temos muito mais para aprender, pois a vida é mutável, assim como os transposons e conhecê-la é um projeto sem fim.

Devo também a elas a oportunidade de conhecer e trabalhar com pessoas maravilhosas que me ajudam no dia-a-dia, tirando dúvidas e construindo novas ideias e novos caminhos: Lidiane Amorim, que me convidou a voltar ao laboratório e que foi minha tutora junto com Luis Carlos Belarmino, que me iniciou na bioinformática; João Pacifico e Marx Oliveira, que estavam sempre ao lado tirando as dúvidas das mais básicas às mais insolúveis. Flávia Araújo, Sheyla Lima e Mitalle Matos, que junto comigo sentavam para discutir os projetos e conhecer melhor as sequências e ferramentas de estudo de cada uma. Obrigada também a Carol Wanderlei e Hévila Mendes, que também me socorreram nas correrias da vida acadêmica; e a Bianca Ximenes e ao Joselito Nascimento, bem como Raul

Maia e Sérgio Paiva, que em meio a tudo arrumaram tempo para solucionar questões de uma usuária básica da bioinformática.

Eu não posso deixar de incluir aqueles que me ajudaram a trilhar o caminho até aqui. Desde os meus professores no ensino médio que me inspiraram e apresentaram à minha paixão, a Genética, até a professora Bereneuza, que me acolheu no meu primeiro estágio e Adauto Gomes Neto que como meu tutor na GENETEC me acompanhava diariamente. A Celuza que me apresentou a Valesca Padonfi, uma profissional que admiro muito e que me ensinou muito enquanto tive o prazer de trabalhar com ela. Sem esquecer-se de agradecer à Vanessa Valetin por todos os e-mails trocados, avisos dados e risadas compartilhadas, e aos professores que estiveram conosco durante as disciplinas guiando o caminho para aprofundar o conhecimento sobre genética.

Imprescindíveis na minha vida, agradeço aos amigos feitos nessa jornada e aos que eu trouxe comigo, mais que isso tenho que agradecer à minha família pelo apoio que me foi dado, e pela paciência de aguentar um humor gerido pelo projeto de Mestrado, instável, incrível e alucinante. Viviane e Rafael, se não fossem vocês eu não teria como aguentar a carga emocional de um Mestrado e meu pai que a seu modo discreto e calado torceu por mim. E um agradecimento aos céus, ao meu Avô amado que estudou comigo, que me acompanhou nas entrevistas de estágio, que vibrava comigo como se cada conquista fosse dele próprio, que ouvia cada explicação pacientemente mesmo sem entender nada apenas para satisfazer o meu desejo de compartilhar com ele cada nova descoberta.

Obrigada a todos que fizeram desse projeto possível em todas as suas variáveis.

*“If you’re too big to follow the rivers,
How you ever gonna find the sea”*

- Emeli Sande-
(River ♪)

RESUMO

Os elementos transponíveis (TEs) são sequências móveis de DNA, presentes em quase todos os organismos já estudados. Constituem parte abundante dos genomas, influenciando em sua estrutura e atividade de diversas maneiras. Dessa forma, conhecer a composição e a distribuição de TEs em leguminosas como *Vigna unguiculata*, *Phaseolus vulgaris* e *Medicago truncatula* é parte importante do conhecimento da influência destes nos genomas vegetais. Usando ferramentas bioinformáticas baseadas em homologia, esse trabalho identificou um total de 5,64 Mb de transposons e 48,42 Mb de retrotransposons, distribuídos em todos os cromossomos de *P. vulgaris*; 14,04 Mb de transposons e 51,20 de retrotransposons em *M. truncatula*, e 2,38 Mb e 9,34 Mb de transposons e retrotransposons, respectivamente, em *V. unguiculata*. Todos apresentaram *CACTA* e *Mutator* entre as três superfamílias de transposons mais abundantes no genoma. Além disso, a identificação de transcritos em cultivares contrastantes quanto à tolerância à seca em *V. unguiculata* mostrou que a cultivar tolerante apresentou maior número de transcritos de *CACTA* e *Mutator* que a cultivar susceptível. Interessantemente apesar de possuir o menor genoma, *M. truncatula* foi a leguminosa com maior quantidade de TEs, possivelmente por sua maior tolerância à atividade destes elementos ou menor controle epigenético. Por outro lado, a despeito dos riscos inerentes à sua atividade, a maior abundância de transcritos na cultivar tolerante à seca de *V. unguiculata* indica que a atividade dos TEs pode conferir uma vantagem evolutiva.

Palavras-chave: Elementos Transponíveis. Feijão-caupi. Feijão comum. Transposase. Bioinformática.

ABSTRACT

Transposable elements (TEs) are mobile DNA sequences, present in almost all organisms studied to date. They constitute the most part of genomes, acting in their structure and activity at different levels. Therefore, the knowledge of TE composition and distribution in legumes like *Vigna unguiculata*, *Phaseolus vulgaris* and *Medicago truncatula* are important for a better knowledge about their influence on plant genomes. Bioinformatics based on homology search was used in this work and identified, in all chromosomes, about 5,64 Mb of transposons and 48,42 Mb of retrotransposons in *P. vulgaris*; 14,04 Mb of transposons and 51,20 of retrotransposons in *M. truncatula*; and 2,38 Mb and 9,34 Mb of transposon and retrotransposons, respectively, in *V. unguiculata*. All species had *CACTA* and *Mutator* among the three most abundant superfamilies of transposons. Additionally, the identification of transcripts at the drought tolerant *V. unguiculata* cultivar showed a higher number of TEs (*CACTA* and *Mutator*) than the susceptible cultivar. Interestingly, the smallest legume genome, *M. truncatula*, was the one richer in TE sequences, maybe because of its higher tolerance to TE activity or lower epigenetic control. On the other hand, despite all risks related to TE activity, the greater abundance of transcripts in the *V. unguiculata* drought tolerant cultivar indicates that the activity of TEs may confer an evolutionary advantage.

Key words: Transposable Elements. Common bean. Cowpea. Transposase. Bioinformatics.

LISTA DE ILUSTRAÇÕES

Ilustrações	Página
Figura 1 - Estrutura esquemática de um retrotransposon (<i>LTR - Long Terminal Repeat</i>)	24
Figura 2 - Estrutura esquemática de um transposon (<i>Mutator</i>)	25
Figura 3 - Representação esquemática das TIRs (Repetições Terminais Invertidas, do inglês: <i>Terminal Inverted Repeat</i>)	26
Figura 4 - Esquema de geração das duplicações de sítio alvo (TSD, do inglês: <i>Target Site Duplication</i>)	27
Figura 5 - Estrutura esquemática de um elemento CACTA autônomo	35
Figura 6 - Estrutura do elemento Spm e dos transcritos codificados por Spm. Fonte: Masson <i>et al.</i> (1991)	36
Figura 7 - Figura esquemática de um elemento MuK (<i>Mutator killer</i>)	39
Figura 8 - Fluxograma de etapas referentes à metodologia	53
Figura 9 - Representação gráfica da quantidade em número absoluto e em Megabases (Mb) das superfamílias de elementos transponíveis por espécie	56
Figura 10 - Domínios gênicos mais abundantes encontrados em associação com transposons <i>Mutator</i> na cultivar tolerante 'Pingo de Ouro' de <i>Vigna unguiculata</i>	60
Figura 11 - Domínios gênicos mais abundantes encontrados em associação com transposons <i>Mutator</i> na cultivar tolerante 'Pingo de Ouro' de <i>Vigna unguiculata</i>	61
Figura 12 - Domínios gênicos mais abundantes encontrados em associação com transposons <i>CACTA</i> na cultivar sensível 'Santo Inácio' de <i>Vigna unguiculata</i>	61
Figura 13 - Domínios gênicos mais abundantes encontrados em associação com transposons <i>CACTA</i> na cultivar sensível 'Santo Inácio' de <i>Vigna unguiculata</i>	62
Figura 14 - Alinhamento do domínio da transposase de <i>Mutator</i> das cultivares de <i>Vigna unguiculata</i> 'Pingo de ouro Ouro' (PO) e 'Santo Inácio' (SI)	63

Figura 15 - Disposição cromossômica das transposases (TNPs) e quinases de <i>V. unguiculata</i> em cromossomo de <i>P. vulgaris</i>	64
Quadro 1 - Quantidade de sequências de <i>Mutator</i> identificadas	58
Quadro 2 - Distribuição dos domínios por frame de leitura.	59

LISTA DE TABELAS

Tabelas e Anexos	Página
Tabela 1 - Espécies da família Fabaceae com genoma completo sequenciado	48
Tabela 2 - Quantidades em números absolutos de elementos transponíveis por cromossomo para <i>Phaseolus vulgaris</i> (Pv) e <i>Medicago truncatula</i> (Mt)	57
Tabela 3 - Comparação da composição de elementos transponíveis (TEs) entre angiospermas.	68
Anexo 1 - Lista das sondas utilizadas neste trabalho.	82
Anexo 2 - Porcentagem (%) e quantidades em Megabases de seqüências para cada superfamília calculada para o tamanho real dos genomas de <i>Medicago truncatula</i> , <i>Phaseolus vulgaris</i> e <i>Vigna unguiculata</i>	89
Anexo 3 - Alinhamento das transposases de <i>Mutator</i> das cultivares contrastantes para tolerância a seca e salinidade de <i>Vigna unguiculata</i>	90
Anexo 4 - Lista com todos os domínios de genes relacionados a seqüências de <i>Mutator</i> identificados em <i>Vigna unguiculata</i>	91

LISTA DE ABREVIATURAS, SIGLAS E SÍMBOLOS

Item	Definição
Ac	Loco Ativador; <i>Locus Activation</i>
BLAST	Ferramenta de Busca de Alinhamento Local Básico; <i>Basic Local Alignment Search Tool</i>
BLASTn	Blast de nucleotídeos; <i>nucleotide BLAST</i>
BLASTx	Busca de proteínas em banco de dados usando nucleotídeos traduzidos; <i>Search protein database using a translated nucleotide BLAST</i>
CBW	Condicional de Baum-Welch
CIG	<i>Center for Integrative Genomics</i>
db	Banco de dados; <i>Data base</i>
DDBJ	Banco de DNA do Japão; <i>DNA Bank of Japan</i>
DDE	Motivo de aminoácidos conservado, dois resíduos D (Aspartato) e E (Glutamato)
DMS	Cirurgia de Modelo dinâmico; <i>Dynamic Model Surgery</i>
Ds	<i>Locus de Dissociação; Locus of Dissociation</i>
Em/Spm	<i>Enhancer e Supressor-mutator</i>
EMBL-EBI	Laboratório de Biologia Molecular Europe – Instituto Europeu de Bioinformática; <i>European Molecular Biology Laboratory – European Bioinformatics Institute</i>
EN	Endonuclease
EST	Etiqueta de Sequências expressas; <i>Expressed Sequence Tag</i>
GAG	Proteína do capsídeo viral, também presente em Retrotransposons
HMM	Modelo oculto de Markov; <i>Hidden Markov Model</i>
ICR	Região complementar interna; <i>Internal Complementary Region</i>
INT	Integrase
JGI	<i>Joint Genome Institute (United States Department of Energy)</i>
kda	Quilo Dalton (do inglês kilodalton), medida de peso das proteínas
KEGG	<i>Kyoto encyclopedia of Genes and Genomes</i>
KOG	Grupos de Ortólogos Eucariotos; <i>EuKaryotic Orthologous Groups</i>
LINEs	Longos elementos intercalados; <i>Long Interspersed Elements</i>
LTRs	Repetições terminais longas; <i>Long Terminal Repeats</i>
mipsREdat	<i>MIPS Repeat Element tadabase</i>

miRNA	Micro RNA
MITEs	Miniatura invertida-repetida de Elemento Transponível; <i>Miniature Inverted-repeat Transposable Element</i>
MuK	<i>Mutator-Killer</i>
MULE	<i>Mutator-like</i>
NCBI	Centro Nacional de Informação Biotecnológica; <i>National Center for Biotechnology Information</i>
nt	Nucleotídeo
ORF	Quadro aberto de leitura; <i>Open Reading Frame</i>
PA	Proteína aspártica
pb	Pares de Bases
PFAM	<i>Protein Families Database</i>
PGI	<i>Pigeonpea genomics initiative</i>
PLE	<i>Penelop-like Elements</i>
PO	Pingo de Ouro – Cultivar tolerante
RH	RNAse H
RM	RepeatMasker
RNapol	RNA polymerase
RNA _t	RNA transportador
RT	Transcriptase reversa; <i>Reverse Transcriptase</i>
RU	<i>Repbase update</i>
SI	Santo Inácio, cultivar susceptível
SINE	Pequenos elementos intercalados; <i>Short Interspersed Elements</i>
siRNA	Pequeno RNA de interferência
sRNA	Pequeno RNA; <i>Small RNA</i>
STR	Regiões subterminais; <i>Subterminal regions</i>
TEs	Elementos transponíveis; <i>Transposable elements</i>
TIR	Repetições terminais invertidas; <i>Terminal Inverted Repeats</i>
TNP/Tnp	Transposase
Tport	Transportador
TSD	Duplicação de sítio alvo; <i>Target Site Duplication</i>
UniProt	<i>Universal Protein Resource</i>
UTR	Região não traduzida; <i>Untranslated region</i>
V(D)J	V – Variável, D – Diversificado e J – junto, combinado; V – <i>Variable</i> , D – <i>Diversity</i> , J – <i>Joining</i>

LISTA DE WEBSITES CITADOS

Sigla	Banco de Dados	Site	Acesso
CGKB	<i>Crop Genebank Knowledge Base</i>	http://cropgenebank.sgrp.cgiar.org/	-
CGP	<i>The Compositae Genome Project</i>	http://compgenomics.ucdavis.edu/	-
DDBJ	<i>DNA data Bank of Japan</i>	http://www.ddbj.nig.ac.jp	-
EMBL-EBI	<i>European Molecular Biology Laboratory – European Bioinformatics Institute</i>	http://www.ebi.ac.uk	-
IPGI	<i>The International Peanut Genome Initiative</i>	http://www.peanutbase.org/home	
KEEG	<i>Kyoto Encyclopaedia of Genes and Genomes</i>	www.genome.jp/kegg/	-
KOG	<i>EuKaryotic Orthologous Groups</i>	http://bioinfo03.ibi.unicamp.br/vigna/	-
NCBI	<i>National Center for Biotechnology Information</i>	http://www.ncbi.nlm.nih.gov/	-
NordEST	<i>NordEST – Vigna Project</i>	http://bioinfo03.ibi.unicamp.br/vigna/	Restrito
mipsREdat	<i>MIPS Repeat Element tadabase</i>	http://www.transplantdb.eu/node/2249#)	-
PANTHER	<i>PANTHER</i>	www.pantherdb	-
PeanutBase	<i>PeanutBase</i>	http://peanutbase.org/	Restrito*
PFAM	<i>Protein Families Database</i>	pfam.sanger.ac.uk/	-
PGI	<i>Pigenonpea genomics initiative</i>	http://www.nrcpb.org/content/pigeon-pea-genome	-
Phytozome	<i>Phytozome</i>	http://www.phytozome.net/	-
RefSeq	<i>Reference Sequence Database</i>	http://www.ncbi.nlm.nih.gov/RefSeq/	
RU	<i>Repbase Update</i>	http://www.girinst.org	Restrito**
TAIR	<i>The Arabidopsis Information Resource</i>	www.arabidopsis.org/	-
UniProt	<i>Universal Protein Resource</i>	www.uniprot.org/	-

Ps.: Na última coluna (Acesso) estão sinalizados com “ - ” os bancos de dados públicos e por tanto de acesso livre, e com “Restrito” bancos de dados de acesso restrito a determinado grupo de pesquisa e seus colaboradores. * O banco de dados PeanutBase permite o acesso e análise online, porém o Download e a análise completa do genoma disponível apenas são liberados após aceite do termo de não publicação (informação completa no site). ** Acesso permitido após cadastro com e-mail institucional.

SUMÁRIO

1 INTRODUÇÃO	16
2 REVISÃO DA LITERATURA	18
2.1 As culturas de <i>Vigna unguiculata</i> , <i>Phaseolus vulgaris</i> e <i>Medicago truncatula</i>	18
2.2 Elementos transponíveis: descoberta, conceito, atuação e abrangência	20
2.3 Classificação dos elementos transponíveis	23
2.3.1 retrotransposon	28
2.3.1 dna-transposon	30
2.4 Sequências de estudo	34
2.4.1 CACTA	37
2.4.2 Mutator	40
2.5 Elementos transponíveis e a regulação da expressão	42
2.6 Bioinformática	42
2.6.1 Identificação e caracterização dos elementos transponíveis	42
2.6.2 Bancos de dados	46
3 OBJETIVOS	49
4 MATERIAL E MÉTODOS	50
4.1 Mineração de sequências candidatas	50
4.2 Recuperação e tradução das sequências	50
4.3 Caracterização	51
4.4 Inferência, genoma sequenciado X genoma real	51
4.5 Ancoragem	52
5 RESULTADOS	54
6 DISCUSSÃO	65
7 CONCLUSÃO	74
REFERÊNCIAS	76
ANEXOS 1	82
ANEXO 2	89
ANEXO 3	90
ANEXO 4	91

1 INTRODUÇÃO

Descritos pela primeira vez em milho por Bárbara McClintock em 1950, os elementos transponíveis (TEs) são sequências genômicas com capacidade de se movimentar através do genoma. Hoje já se sabe que os TEs figuram entre os componentes mais abundantes em genomas eucariotos e que estão dispersos em todos os reinos, atuando como força evolutiva e propulsora da diversidade. Sua variedade os distinguiu inicialmente em duas grandes classes de acordo com seu mecanismo básico de ação: a classe I (retrotransposons), que é caracterizada por se mover via sequências de RNA, e a classe II (DNA-transposons ou transposons), que baseia a sua dinâmica na movimentação direta das sequências de DNA. Em decorrência de sua mobilidade, os elementos transponíveis são capazes de carregar fragmentos de diversas sequências e, dependendo da superfamília a que pertencem, seus impactos vão além das mutações acarretadas por sua mobilidade, podendo interferir na expressão de outros genes.

Com estruturas distintas, as superfamílias de transposons *CACTA* e *Mutator* se diferenciam principalmente pelo tamanho pequeno das repetições terminais invertidas (TIRs), variando de poucas dezenas de pares de bases em *CACTA* a centenas de pares de bases em *Mutator*, e pelo número de transposases essenciais à mobilidade de cada um. Apesar das diferenças estruturais, estas superfamílias destacam-se por sua característica comum: preferência de inserção em regiões ricas em genes, podendo exercer maior influência sobre estas regiões. Sua atividade, no entanto, nem sempre é benéfica, o que leva à necessidade de mecanismos de controle. Neste contexto, o controle

epigenético ocorre especialmente por metilação, o que justifica o aumento da sua atividade quando o genoma encontra-se hipometilado.

Tendo em vista tantos aspectos relacionados a esses elementos e sua relação com os genomas, a identificação, caracterização e quantificação desses elementos consistem nos primeiros passos para a compreensão de seus mecanismos e funções, além de serem essenciais à compreensão das dimensões de sua influência.

V. unguiculata (L.) Walp, mais conhecida como feijão-caupi, é uma das principais fontes de proteína vegetal e carboidratos para famílias de baixa renda em diversas regiões tropicais. Embora seu genoma completo ainda não esteja disponível, bancos de dados de transcriptômica enriquecem o conhecimento sobre esta Fabaceae. Complementarmente *Phaseolus vulgaris* L. e *Medicago truncatula* Gaertn representam importantes leguminosas, onde a primeira apresenta grande importância socioeconômica e a segunda é reconhecida como leguminosa modelo. Ambas possuem genoma sequenciado disponível permitindo análises comparativas e maiores inferências. Desta forma, o presente trabalho visou conhecer a quantidade e a distribuição dos elementos *CACTA* e *Mutator* nessas leguminosas, avaliando genes que são possivelmente influenciados pelos mesmos e associando-os ao mapeamento genético e citogenético do feijão-caupi. Tais estudos auxiliam no avanço dos conhecimentos sobre o genoma das leguminosas estudadas e na compreensão tanto do papel estrutural quanto funcional desses elementos. Além disso, TEs são excelentes candidatos para uso em projetos de mapeamento genético e cromossômico, tendo em vista sua abundância, ubiquidade genômica e rápidas taxas evolutivas.

2 REVISÃO DA LITERATURA

2.1 As culturas de *Vigna unguiculata*, *Phaseolus vulgaris* e *Medicago truncatula*

Nativo da África (Freire *et al.*, 1999) o feijão-caupi [*V. unguiculata* (L.) Walp] é uma Fabaceae cultivada em 97 países entre todos os continentes (Freire-Filho, 2011). A sua produção mundial soma mais de 5,4 milhões de toneladas. No Brasil a estimativa é que entre 2005 e 2009 a produção do feijão-caupi alcançou uma produção média de ~513,62 toneladas (15,48% da produção de feijão). Presente na mesa de 26,39 milhões de pessoas no Norte e Nordeste, com um consumo médio calculado em 20 kg/pessoa/ano, é uma das principais fontes de proteínas, carboidratos, vitaminas e minerais, especialmente entre comunidades rurais e população de baixa renda o feijão-caupi é bastante cultivado nas regiões Norte e Nordeste do Brasil por pequenos produtores entre os quais agricultores familiares, embora encontre-se em processo de expansão para as regiões Centro-oeste e Sudeste (Freire-Filho *et al.*, 2011),

O feijão-caupi se destaca por seu vigor e poder adaptativo a diversas condições edafoclimáticas como tolerância a estiagem prolongada e crescimento em solos de baixa fertilidade, o que torna essa cultura interessante fonte de genes para o melhoramento genético. Contudo, apesar de sua diversidade genética capaz de tolerar uma variedade de estresses, sua produção ainda é menor do que o potencial observado para essa leguminosa (Rodrigues *et al.*, 2005; Freire-Filho *et al.*, 2011). Assim projetos vêm sendo desenvolvidos objetivando o melhoramento do feijão-caupi, com ênfase para tolerância à seca, salinidade e a múltiplas viroses, bem como a doenças fúngicas e bacterianas além de

características anatômicas, que melhorem a produção como porte compacto e ereto, melhorando a produção e facilitando a colheita (Pimentel *et al.*, 2002; Wang *et al.*, 2003; Freire-Filho *et al.*, 2011).

O feijão comum (*P. vulgaris* L.), apesar de cultivado em todos os países de clima tropical e subtropical, tem no Brasil seu centro de maior produção e consumo mundial. Cultivado especialmente por agricultores familiares, seu plantio ocorre também em regiões semiáridas onde a perda da produção é elevada. Extremamente sensível às variações ambientais, sua produção nacional 2014/15 chegou a 3.338,4 toneladas (CONAB, jan. de 2015). No entanto, a sua produção nacional anual é considerada baixa em relação ao seu potencial produtivo. Fonte de proteína vegetal e nutrientes, é utilizado muitas vezes para substituir carnes entre populações de baixa renda, formando junto com o arroz a base alimentar de várias famílias (Dos-Santos *et al.*, 2011; Gomes *et al.*, 2012; Brito *et al.*, 2013; Mingotte *et al.*, 2013).

Parente da alfafa, *Medicago truncatula* Gaertn., por sua vez, destaca-se como leguminosa modelo, de pequeno porte, genoma diploide com tamanho de 5×10^8 pb, passível de autopolinização natural e desenvolvimento rápido. Importante para o entendimento dos mecanismos de fixação de nitrogênio na interação planta-rizóbio, esta leguminosa permite uma melhor compreensão sobre estrutura do genoma, função de genes e proteínas da família Fabaceae, possuindo diversas variações fenotípicas para características importantes como hábito de crescimento, tempo de floração, especificidade simbiótica e resistência a estresses (Cook, 1999; Kurdyukov *et al.*, 2013).

2.2 Elementos transponíveis: descoberta, conceito, atuação e abrangência

Os elementos transponíveis foram descobertos em 1950 pela pesquisadora Bárbara McClintock (ganhadora do Nobel de 1983 por tal achado) alterando o padrão de cores dos grãos de milho. Ela os chamou de *loci* mutáveis ou genes instáveis em sua primeira descrição no comunicado em abril de 1950 e neste fez um alerta sobre sua descoberta: “As observações acumuladas sobre esses numerosos *loci* mutáveis são tão extensivas, que nenhum breve relato daria suficiente informação para preparar o leitor para um julgamento independente sobre a natureza desse fenômeno” (McClintock, 1950; nobelprize.org, web 18 de nov de 2014).

Alterações cromossômicas no braço curto do cromossomo nove de milho em taxas muito acima do comum foram associadas a um *locus* muito instável, hoje conhecido amplamente como ‘elemento transponível’ (TE). Ao avaliar essas alterações, McClintock suspeitou de um mecanismo básico associado a mudanças de estado na cromatina que acarretariam em tais mutações. Ao reconhecer o primeiro *locus* responsável pela quebra cromossômica ela o chamou de *locus* de Dissociação (*Ds*). Ao tentar localizar a posição exata de *Ds* no cromossomo o resultado foi inesperado, mostrando esse *locus* em locais diferentes, inclusive em outros cromossomos além daquele onde foi identificado inicialmente. Além disso, as novas posições também estavam associadas a posições de quebras cromossômicas, indicando que a natureza das mutações estava relacionada à transposição de *Ds* (McClintock, 1950; McClintock, 1951).

De acordo com a pesquisadora, as inserções de *Ds* sintetizavam todas as outras ocorrências provocadas, como: (1) fusão de cromátides; (2) deleção de

segmentos da cromátide; (3) translocação com ponto de quebra no *locus Ds*; (4) duplicações de segmentos com ponto de quebra no *locus Ds*; (5) transposição de *Ds* de uma posição a outra; (6) perda de atividade detectável de *Ds*; (7) mudanças no próprio *locus Ds*, chamado “mudança de estado”. Este último levando a alterações quanto à sua frequência e tempo de desenvolvimento em que seria expresso em gerações futuras. No entanto, um segundo *locus*, chamado Ativador (*Ac*), mostrou-se fator dominante para ocorrência das anomalias cromossômicas. Sua atividade se mostrou independente de qualquer fator externo e na sua ausência nenhuma atividade em *Ds* era observada, nem qualquer quebra cromossômica; comprovando, conseqüentemente, a dependência de um dos *loci*. Observado isso, McClintock classificou *Ac* e *Ds* em *loci* autônomo e não autônomo, respectivamente. Em seus relatos, McClintock suspeita que tais elementos eram compostos de matéria similar ou igual; mas, além da certeza de que eles se moviam através do genoma, deixa diversos questionamentos sobre seus mecanismos e atividade (McClintock, 1950; McClintock, 1951).

Apesar de ignorado por algum tempo, seu trabalho quebrou dogmas e abriu as portas para uma nova visão sobre a estrutura do genoma, sua organização, controle de sua expressão, e sua capacidade de gerar alterações, que podiam ser importantes evolutivamente. Sem saber a total abrangência e influência de sua descoberta no genoma, Bárbara McClintock trouxe ao conhecimento científico uma fonte de variabilidade genética, arma secreta da adaptabilidade e evolução dos organismos.

Devido à sua natureza móvel, a interferência desses elementos pode ter resultados algumas vezes deletérios. A especificidade de inserção dos TEs ainda

não é bem compreendida. Inicialmente, acreditava-se na aleatoriedade de sua inserção; porém, hoje já se sabe que famílias específicas podem ter alvos bem estabelecidos ou regiões preferenciais de inserção (Li *et al.*, 2009; Liu *et al.*, 2009). A princípio justificado a partir da similaridade de sequência, análises recentes demonstram que isso não é sempre verdade, dessa forma é prudente afirmar que o mecanismo responsável pela escolha do local de inserção ainda não é bem compreendido. Uma nova possibilidade levantada é a de que diferentes tipos de transposases e integrases tenham preferências por diferentes tipos de cromatina (Bennetzen e Wang, 2014).

Além dos TEs atuarem na inativação de genes e geração de novos genes, vários pseudogenes e promotores de genes são derivados de transposons ou tiveram origem em duplicações induzidas por eles, podendo apresentar seguimentos de TEs inseridos em sua composição (Bennetzen, 2000; Hanada *et al.*, 2009). Adicionalmente, os TEs afetam a regulação transcricional e pós-transcricional de outros genes, alterando padrões epigenéticos e gerando pequenos RNAs (sRNAs) (Bennetzen *et al.*, 2014). Concomitantemente, estudos revelam a indução da expressão de TEs em resposta a condições de estresses biótico e abiótico em plantas, sendo observados em maior quantidade em plantas tolerantes que em plantas sensíveis a tais estresses. Há indícios de que eles possam promover superexpressão em genes próximos, embora os mecanismos não sejam conhecidos (Hanada *et al.*, 2009; Chénais *et al.*, 2012; Makarevitch *et al.*, 2014). Curiosamente, existem também evidências circunstanciais que indicam a transferência horizontal de TEs entre espécies, embora o mecanismo de transferência ainda permaneça obscuro (Bennetzen, 2000; Hanada *et al.*, 2009; Liu *et al.*, 2009; Bennetzen e Wang, 2014).

Os elementos transponíveis compõem um grupo significativo de sequências, que foram inicialmente designados como “DNA-lixo” e referenciados apenas para justificar a correlação negativa entre o tamanho dos genomas e a complexidade das espécies (Biémont e Vieira, 2006; Gemayel *et al.*, 2012). Atualmente, é reconhecido que os TEs possuem um papel crucial na evolução dos genomas, apresentando-se como parte constituinte em todos os eucariotos já estudados, com exceção apenas do protozoário *Plasmodium falciparum* (Chénais *et al.*, 2012). Além de abranger uma grande variedade de espécies, os TEs compõem uma grande proporção dos genomas analisados. Em plantas podem variar entre cerca de 10% do genoma de *Arabidopsis thaliana* L. (The *Arabidopsis* Initiative, 2000) a 85% do genoma de milho (Chénais *et al.*, 2012). Já em humanos a estimativa é de que aproximadamente 45% do genoma seja composto por TEs, enquanto que em ratos esse valor está ao redor de 40%, na mosca da fruta em torno de 15% e no mosquito (*Anopheles sp.*) 25% (Feschotte, 2004; Chénais *et al.*, 2012).

2.3 Classificação dos elementos transponíveis

Como consequência de sua abrangência e diversidade, surgiu a necessidade de organizar os TEs em um sistema de classificação complexo. De modo semelhante ao sistema de classificação dos seres vivos, com base em seu mecanismo básico de transposição e sua composição enzimática, os TEs foram classificados em duas grandes classes, com níveis hierárquicos menores: subclasse, ordem, superfamília e família. Esta forma de organizá-los foi proposta por Wicker *et al.* (2007) e buscou unificar a sua classificação, como descrito a seguir:

(1) Classe I - Retrotransposons (Figura 1): Se movem a partir de um RNA gerado a partir de sua sequência, que é traduzido para cDNA por uma transcriptase reversa, e transportado por uma transposase, que reconhece o sítio de inserção e promove a inserção e fixação da cópia do retrotransposon. Este mecanismo é conhecido como “copia-e-cola”. Dessa forma, os retrotransposons multiplicam seu número de cópias e aumentam o tamanho do genoma a cada evento de transposição (Wicker *et al.*, 2007; Joly-Lopez e Bureau, 2014).



Figura 1: Estrutura esquemática de um retrotransposon (*LTR - Long Terminal Repeat*). As setas pretas representam as longas repetições terminais e as caixas de cor cinza as ORFs (Quadro aberto de leitura; do inglês: *Open Reading Frame*) contidas no elemento. Dentre as ORFs podem ser observadas: GAG (proteína de capsídeo), PA (proteínase aspártica), INT (integrase), RT (transcriptase reversa) e RH (RNase H). Fonte: Adaptada de Wicker *et al.* (2007).

(2) Classe II - DNA transposons ou simplesmente transposons (Figura 2): São caracterizados por se moverem através do mecanismo conhecido como “corta-e-cola”. Podem possuir uma ou duas transposases (TNP) em sua estrutura, sendo esta responsável por cortar o DNA e o movê-lo diretamente a outra posição. Por esse motivo seu número de cópias tende a ser bem menor que o dos retrotransposons. Contudo, alguns elementos apresentam capacidade de multiplicar seu número no genoma, apresentando-se como exceção à regra (Wicker *et al.*, 2007; Gifford *et al.*, 2013; Joly-Lopez e Bureau, 2014). Dessa forma, os transposons podem ser subdivididos em duas subclasses: subclasse I, que se movimenta sem alterar o número de cópias, e subclasse II, que se move a partir das cópias geradas (Wicker *et al.*, 2007).



Figura 2: Estrutura esquemática de um transposon (*Mutator*). Os triângulos brancos representam as duplicações de sítio alvo (TSD, do inglês: *Target Site Duplication*); as setas pretas, as repetições terminais invertidas (TIRs, do inglês: *Terminal Inverted Repeats*), enquanto a caixa cinza refere-se à ORF (Quadro aberto de leitura, do inglês: *Open Reading Frame*), representada pela transposase (TNP). Fonte: Esquema do Autor.

As ordens permitem separar os elementos de acordo com as peculiaridades enzimáticas e organizacionais dos diferentes elementos; seguidas pelas superfamílias, que são baseadas na similaridade de sequências, incluindo aquelas das regiões não codificantes, e por último as famílias que são determinadas por sequências que apresentam alta similaridade nas regiões codificantes ou nas repetições terminais (pelo menos 80%, para seguimentos de no mínimo 80 pb). Em alguns casos, as famílias podem ser definidas também por elementos autônomos e não autônomos, que se diferenciam respectivamente pela presença ou não da transposase. Os não autônomos não possuem a transposase, mas apresentam toda estrutura necessária para serem por ela reconhecidos e possuem a capacidade de se mover dependendo da enzima de outros elementos (Wicker *et al.*, 2007).

Foi convencionado que toda classificação abaixo do nível hierárquico de ordem deve ser escrita em itálico. Wicker *et al.* (2007) também propuseram um sistema de nomenclatura para os elementos. Por exemplo, o elemento *Caspar* representante da superfamília *CACTA*, ordem TIR, subclasse I, classe II seria anotado → DTC_*Caspar*_AA123456-3, DTC representando a classificação geral (DNA-transposon + Ordem TIR + Superfamília CACTA) seguido do nome da família e número de acesso ou do BAC em que se encontra a sequência. Caso

alguma das classificações seja desconhecida, essa deverá ser substituída pela letra 'x' minúscula (Wicker *et al.*, 2007).

Para classificar um elemento encontrado, Wicker *et al.* (2007) sugerem uma ordem de anotação: (1) BLASTn contra um banco específico de TEs; (2) BLASTx contra um banco de dados de proteínas caracterizadas provavelmente possibilitaria identificação da classe, ordem e até superfamília; (3) análise minuciosa da estrutura quanto à presença ou ausência de repetições terminais invertidas (TIRs, do inglês: *-Terminal Inverted Repeats* – Figura 3), de repetições terminais diretas, de repetições subterminais diretas e inversas, sítio de ligação ao DNA, e duplicações de sítio alvo (TSDs, do inglês: *Target Site Duplication*, Figura 4) (Wicker *et al.*, 2007).

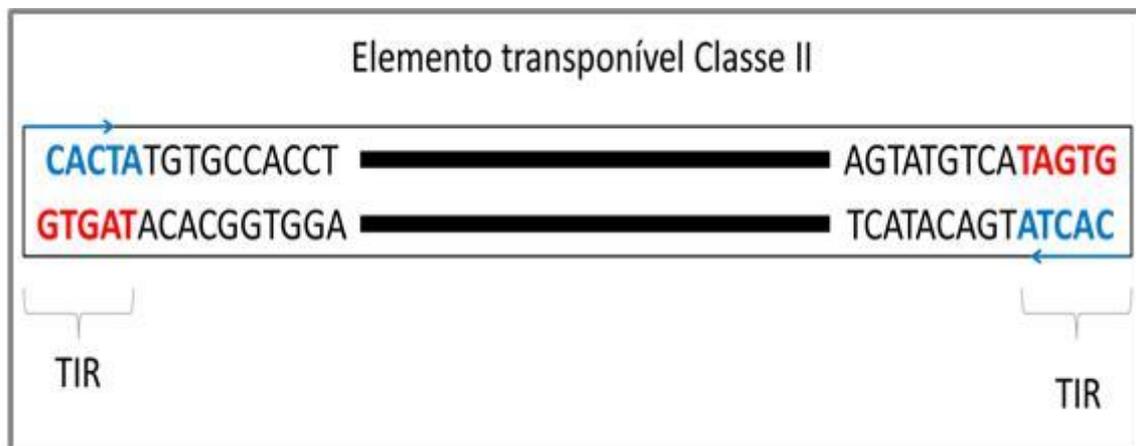


Figura 3: Representação esquemática das TIRs (Repetições Terminais Invertidas, do inglês: *Terminal Inverted Repeat*). Em azul e vermelho estão representadas as TIRs. As linhas pretas representam as regiões internas do elemento, que no exemplo em questão refere-se a um elemento da superfamília *CACTA*. Fonte: Esquema do Autor.

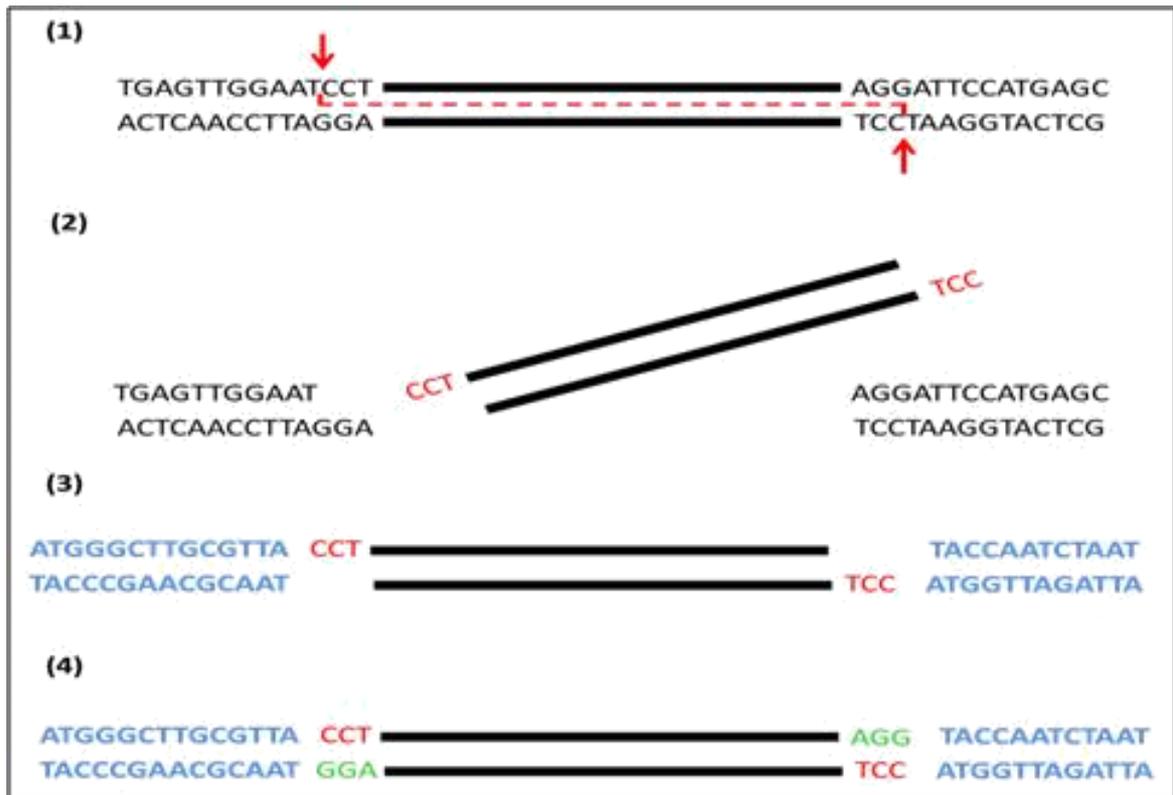


Figura 4: Esquema de geração das duplicações de sítio alvo (TSD, do inglês: *Target Site Duplication*). Em (1) representação de um elemento transponível (TE), setas vermelhas indicam local de corte enzimático. Em (2) a excisão do TE, nucleotídeos em preto representam a sequência genômica que permanece e em vermelho a região do TE que será duplicada. Em (3) o TE inserido em uma nova região genômica, representada em azul. Em (4) o TE inserido e os nucleotídeos em verde representam a região duplicada. Fonte: Esquema do Autor.

Seguindo esse sistema, uma superfamília conhecida como *MITE* (*Miniature Inverted Repeat*) não se encaixa claramente em nenhum dos grupos, apesar de serem geralmente classificados como elementos da classe II por serem flanqueados por TIRs. Os *MITEs* são elementos exclusivamente não autônomos que podem ser ativados por diferentes superfamílias de TEs, mas seu mecanismo básico de transposição é a partir da inserção de suas cópias, como ocorre em geral para elementos de classe I (Wicker *et al.*, 2007; Guernonprez *et al.*, 2012).

2.3.1 Retrotransposons

Os retrotransposons são divididos em cinco ordens: (1) LTRs (*Long Terminal Repeats*); (2) DIRS-LIKE; (3) PLEs (*Penelope-like*); (4) LINEs (*Long Interspersed Nuclear Elements*) e (5) SINEs (*Short Interspersed Nuclear Elements*), sendo as últimas quatro conhecidas também como não LTRs.

(1) LTRs: Menos abundantes em animais, são predominantes nas plantas, e encontrados em todos os reinos. Evolutivamente relacionados aos retrovírus, os LTRs codificam estruturas proteicas características destes, tais como o domínio GAG e POL. O primeiro codifica uma proteína estrutural *virus-like* e a segunda representa um complexo contendo: a proteína aspártica (PA), a Transcriptase Reversa (RT), a RNase H (RH) e o DDE integrase (INT). É possível encontrar ainda ORFs de função desconhecida e fragmentos capturados de outras sequências. Além das evidências de sua interferência na regulação de outras sequências, são elementos interessantes para estudos de filogenia (Wicker et al., 2007; Du *et al.*, 2010b).

(2) DIRS-LIKE: Não possuem integrase do tipo DDE nem geram TSD, apesar de possuírem repetições terminais invertidas. Estas são às vezes seguidas, na terminação 3', por uma região interna de complementariedade (ICR, *Internal Complementary Region*), importante para a replicação deste elemento. Além de possuírem os domínios GAG, RT/RH, elementos desta ordem se diferenciam dos demais por possuírem uma provável tirosina recombinase e um motivo C-terminal dentro do domínio RT/RH com provável função de adenina metiltransferase. Já foram identificados em vertebrados, nematoides, algas e fungos (Goodwin e Poulter, 2004; Wicker *et al.*, 2007).

- (3) PLEs: Geralmente apresentam organização parcialmente em tandem, formando pseudo LTRs, sendo, por isso, durante algum tempo inseridos na superfamília LTR. Acredita-se que essas formações possam ser importantes para a sua atividade, embora o início do elemento seja anterior, com aproximadamente 170 pb. Diferente do observado para LTRs e LINEs, estes elementos são observados frequentemente, mantendo íntrons que deixam sua RT truncada. Sua única ORF codifica uma RT/EN conectada por uma sequência *link*, com peculiaridade adicional de que sua RT é mais semelhante a uma telomerase, devido a seus motivos funcionais (Evgen'ev e Arkhipova, 2005; Wicker *et al.*, 2007).
- (4) LINEs: São elementos cujo comprimento é de várias quilobases e possuem caracteristicamente duas ORFs (ORF1 e ORF2). A primeira codifica o domínio GAG e a segunda RT/EM. Estas podem apresentar ainda o domínio *zinc-finger-like*, podendo ser ricas em cisteínas, que funcionam como domínios de ligação ao DNA. Geralmente são flanqueados por regiões não traduzidas (UTRs, *Untranslated Regions*) e apresentam na região 3' uma calda poli-A, ou apenas uma região rica em adenina (Schmidt, 1999; Wicker *et al.*, 2007).
- (5) SINEs: Com poucas centenas de pares de bases, são compostos por uma região 5' derivada de RNA transportador (RNAt) que apresenta dois motivos conservados, os quais possuem homologia com o promotor reconhecido pela RNAPol III. A região 3', por sua vez, possui sequências relacionadas aos LINEs, sendo situada em posição anterior a uma região de cauda poli-A ou rica em adenina (Schmidt, 1999; Wicker *et al.*, 2007).

2.3.2 DNA transposons

Os DNA transposons ou simplesmente transposons são divididos em duas subclasses (I e II). A subclasse I possui a ordem TIR com nove superfamílias e a ordem Crypton, com uma superfamília de mesmo nome. A subclasse II possui as ordens Helitron e Maverick, ambas com superfamília de mesmo nome.

Subclasse I

- Superfamílias da ordem TIR

- 1) *CACTA*: Assim como PIF, esta superfamília apresenta duas ORFs, que a partir de transcritos alternativos codificam duas TNPs essenciais para a atividade deste elemento. Apresentam TIRs entre 10 e 30 pb e TSD de 3 pb. Diferente da maioria dos elementos, apresentam STRs com repetições invertidas e diretas de 10 pb a 20 pb cada. Também conhecidos como elementos *Em/Spm*, por serem os dois primeiros elementos identificados nesta superfamília. Até o momento apenas foram identificados em plantas (Masson *et al.*, 1991; Wicker *et al.*, 2003; Wicker *et al.*, 2007).
- 2) *hAT*: Assim nomeada ao juntar as iniciais de três elementos bem caracterizados, *hobo*, *Ac-Ds* e *Tam3*. Esta superfamília se caracteriza por apresentar TIRs de até 10 pb, TSD de 8 pb e pelo menos uma ORF. São identificados quatro domínios: (1) BED *zinc-finger*, que é N-terminal; (2) um de ligação ao DNA que também participa da oligomerização; (3) domínio catalítico (TNP), e (4) um longo domínio de inserção com diversas α -hélices, seguindo a α -hélice um resíduo de glutamato que completa a tríade DDE do domínio catalítico. Já foram reportados em humanos, animais e plantas

sendo observados ativos recentemente nos dois últimos (Wicker *et al.*, 2007; Arensburger *et al.*, 2011).

- 3) *Merlin*: Detectados pela primeira vez em bactérias, mas também observados em humanos e animais. Seus elementos apresentam TIRs de 20 pb a 140 pb e codificam uma transposase maior que 10 kb, que gera TSDs de 8 pb e 9 pb (Feschotte, 2004; Wicker *et al.*, 2007).
- 4) *Mutator*: Com TIRs que podem ser bastante pequenas com poucas dezenas de nucleotídeos ou se estender por centenas de pares de bases. Possuem uma *família* conhecida como *PACK-MULEs*, caracterizada por carregar fragmentos de outros genes. Seu nome derivou da grande quantidade de mutações geradas. Elementos dessa superfamília já foram encontrados em plantas, fungos, protozoários e bactérias (Wicker *et al.*, 2007; Hanada *et al.*, 2009).
- 5) *P*: Elementos inicialmente encontrados em insetos, detectados posteriormente também em metazoários e fungos. Possuem TSDs de 8 pb e 150 pb terminais essenciais para sua atividade que incluem TIRs com aproximadamente 40 pb (Kaufman e Riot, 1992; Wicker *et al.*, 2007).
- 6) *PIF-Harbinger*: Ao contrário da maioria, esses elementos possuem duas ORFs, uma para a transposase e outra para proteína de ligação ao DNA. Sua TIR de pelo menos 14 pb possui o motivo G(N)₅GTT e sua TSD mostra-se com 3 pb. Encontrados em plantas, animais e fungos, esta superfamília apresenta também a capacidade de mobilizar membros da superfamília *MITE* (Grzebelus, *et al.*, 2007; Wicker *et al.*, 2007; Sinzelle *et al.*, 2008).

- 7) *PiggBack*: Também denominada IFP2 é encontrada principalmente em animais, especificamente em insetos. Esses elementos são identificados especialmente por se inserirem preferencialmente na vizinhança de sequências 'TTAA', gerando TSD com 4 pb (geralmente 'TTAA'), possuindo caracteristicamente dois ou três resíduos C/G na sua TIR de tamanho diminuto (~16 pb) e STR com repetições assimétricas diretas de 'TTAA'. Sua transposase única é rica em resíduos de cisteína (Penton *et al.*, 2002; Wicker *et al.*, 2007).
- 8) *Tc1-Mariner*: Encontrada em plantas, animais e humanos, essa superfamília é constituída por TIRs de 40 pb a 50 pb e transposase com motivo DDE/D base para classificação das famílias. Sua transposase apresenta-se com preferência de inserção na vizinhança de pares de base 'TA', gerando TSD de 2 pb (Wicker *et al.*, 2007; Gil *et al.*, 2013).
- 9) *Transib*: Apresentam TIRs com tamanho entre 9 pb e ~500 pb e geram TSD de 5 pb ricos em GC. Curiosamente a transposase desse elemento possui similaridade com proteínas envolvidas com o sistema de recombinação V(D)J, além de sua TIR ser aparentemente ancestral às sequências reconhecidas pela proteína responsável pela recombinação V(D)J. Identificados primeiramente em *Drosophila melanogaster*, já foram observados também em vários invertebrados e, apesar de serem detectados em fungos e plantas, há indícios de que sua ocorrência nesses organismos seja decorrente de transferência horizontal (Wicker *et al.*, 2007; Chen e Li, 2008).

- Superfamília *Crypton* da ordem *Crypton*:

Encontrada inicialmente em fungos patogênicos, já foram identificados em animais, algas e até humanos. Sua transposase se assemelha à tirosina recombinase observada em retrotransposons *DIRS*, porém não possuem RT, sugerindo transposição via DNA. Apesar de não possuírem TIRs, apresentam entre 4 e 6 pb de repetições diretas conservadas na extremidade e TSD de 2 pb (TA) (Wicker *et al.*, 2007; Kojima e Jurka, 2011).

Subclasse II

- Superfamília *Helitron* da ordem *Helitron*:

Transpõe-se via mecanismo do círculo-rolante, comum em bactérias e não geram TSD. São definidos pelos motivos 'TC' na extremidade 5' e 'CTRR' na extremidade 3' (no qual 'R' é uma purina) e possuem 16-20 pb palindrômicos característicos. São elementos autônomos, que produzem pelo menos um tipo-Y2 de tirosina recombinase, com domínio helicase e um iniciador de atividade replicativa. Podem ser observados carregando fragmentos e são encontrados em plantas, animais e fungos (Wicker *et al.*, 2007; Dong *et al.*, 2011).

- Superfamília *Maverick* da ordem *Maverick*:

Também conhecidos como *Polintons* são flanqueados por grandes TIRs e codificam até 11 proteínas, algumas com homologia ao DNA viral. Embora artigos recentes discutam a possibilidade de que estes pertençam à classe I pela quantidade de elementos virais observados, *Mavericks* codificam polimerase B sugerindo mecanismo típico de elemento de classe II. Já

foram encontrados em diversos eucariotos e até o momento não foram identificados em plantas (Wicker *et al.*, 2007; Krupovic *et al.*, 2014).

2.4 Sequências de estudo

2.4.1 CACTA

Os TEs da superfamília *CACTA* (classe II, ordem TIR), também conhecidos como *Em/Spm*, de “*Enhancer*” e “*Supressor-mutator*” (primeiros elementos desta superfamília identificados por Peterson e McClintock, independentemente), foram formalmente nomeados *CACTA* pela presença de cinco nucleotídeos ‘CACTA’ nas regiões TIR. É possível observar também a combinação ‘CACTG’ muitas vezes, sendo estes cinco nucleotídeos a única parte conservada da região terminal entre diferentes famílias (Li *et al.*, 2009).

Esta superfamília apresenta TIRs entre 8 e 30 pb, seguidas das repetições subterminais (STRs - *Subterminal Regions*) diretas e inversas mais conservadas que a TIR, embora ainda variáveis. Seu tamanho varia de 10 a 20 pb para cada repetição. A região TIR junto com a STR pode somar entre 200 e 500 nt, e são as regiões reconhecidas pelas duas transposases características destes elementos (Figura 5). Para *CACTA*, as transposases representam a região mais conservada, podendo ser às vezes a única região eficiente na identificação destes elementos, dado que os cinco nucleotídeos ‘CACTA’ possuem tamanho ineficiente na identificação por homologia. O reconhecimento das regiões TIR e STR pela TNP justifica que as regiões de inserção também apresentem similaridade com as mesmas, porém não explica sua preferência de inserção na região 3’UTR. No entanto, é importante ressaltar que novos trabalhos

demonstram que similaridade de sequência talvez não seja o principal mecanismo de seleção do novo local de inserção (Li *et al.*, 2009; Bennetzen e Wang, 2014).

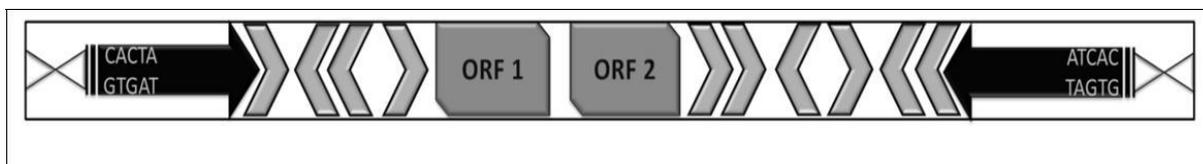


Figura 5: Estrutura esquemática de um elemento *CACTA* autônomo. Os triângulos brancos de cada lado representam as duplicações de sítio alvo (TSD); as setas pretas representam as repetições terminais invertidas (TIRs); as setas em cinza as repetições subterminais (STRs) diretas e inversas, e ao centro os dois quadros de cor cinza representam as duas ORFs, que a partir de transcritos alternativos formam TNP 1 e TNP 2, homólogas de TnpA e TnpD. Fonte: Esquema do Autor.

Estes elementos já foram identificados carregando fragmentos gênicos e até segmentos gênicos completos, que podem alcançar tamanhos maiores que 3 kb em espécies como o milho, *Ipomoea nil* (L.) Roth (conhecida como *Japanese morning glory*), soja e erva-bezerra (*Antirrhinum*). Em milho, foi constatada a presença de até cinco fragmentos de *loci* diferentes em um mesmo elemento, além de transcritos quiméricos e *splicing* alternativo. As análises fenéticas com sequências de milho mostraram que elementos que capturam fragmentos gênicos se agrupam separadamente daqueles que possuem transposase, indicando a possibilidade de um pré-requisito para atividade de captura (Li *et al.*, 2009; Lee *et al.*, 2012).

Essa superfamília apresenta a geração de TSD com 3 pb, geralmente ricas em A/T. Diferentemente de muitos transposons *CACTA* apresenta duas ORFs essenciais à transposição em sua estrutura. De acordo com Masson *et al.* (1991), a TnpA foi a primeira transposase de *CACTA* a ser descoberta em 1986. Apenas em 1991, Masson e colaboradores descobriram as TnpB, TnpC e TnpD (Figura 6) e, mais tarde, os autores provaram que apenas duas delas eram

essenciais a transposição, a TnpA e a TnpD. A TnpA codifica uma proteína de 68 KDa de ligação ao DNA que interage com as regiões STR. Já TnpD é composta por duas ORFs, que fusionadas formam uma proteína de 131 kDa. Dentre as duas ORFs que formam TnpD, uma delas é a mesma da TnpA (Figura 6) (Masson *et al.*, 1991; Wicker *et al.*, 2003; Ping-fang, 2006; Li *et al.*, 2009).

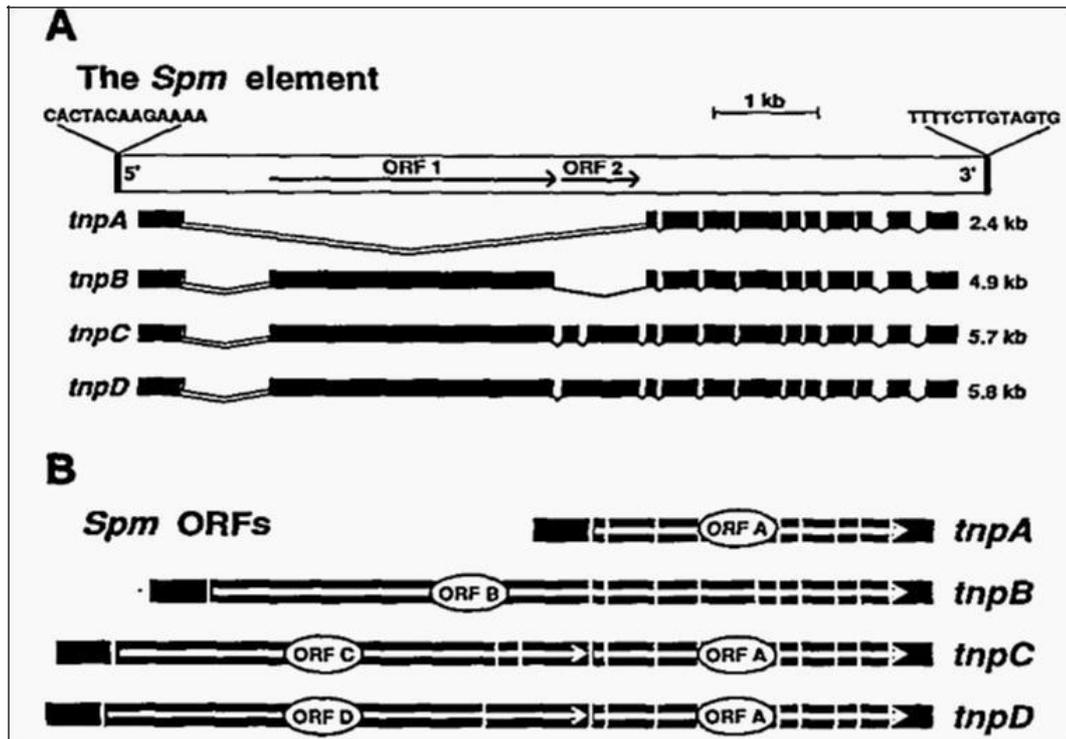


Figura 6: Estrutura do elemento *Spm* e dos transcritos codificados por *Spm*. Fonte: Masson *et al.* (1991).

- (A) O elemento *Spm* é representado pela grande caixa branca. Setas pretas dentro da caixa representam a localização e o tamanho das duas maiores ORFs. Logo abaixo a representação dos íntrons e éxons para cada uma das proteínas. Quadros preenchidos representam éxons dos transcritos alternativos designados TnpA, TnpB, TnpC e TnpD. O tamanho de cada transcrito sem a calda Poli-A está indicado à direita de cada um. Os dois primeiros éxons diferem em 126 pb. O primeiro éxon alternativo é indicado pela caixa em branco após o éxon1.
- (B) O diagrama mostra a maioria das ORFs de cada transcrito deste elemento. A quebra em cada diagrama corresponde a sítios de splicing. ORF-A codifica um polipeptídeo de 68 kDa e está presente como sequência codificante separada nos transcritos TnpA, TnpC e TnpD. ORF-B codifica um polipeptídeo de 171 kDa, fusão da ORF1 com a ORF-A; enquanto TnpC e TnpD são a fusão das ORF1 e ORF2 respectivamente.

Para compreender como estes transcritos influenciavam mudanças no genoma de plantas, Masson *et al.* (1991) fizeram testes transformando plantas de tabaco para que estas possuíssem cada uma das proteínas funcionais individualmente e depois duas a duas ou três a três, incluindo testes em que TnpC e TnpD não possuíam a extensão da ORF-A. Desta forma, conseguiram evidências suficientes para afirmar que tanto TnpA quanto TnpD são essenciais para a atividade, ao contrário dos demais. A TnpA possui ainda uma característica peculiar, capaz de desmetilar sua própria região promotora e conseqüentemente de participar ativamente da regulação desta superfamília. Sugere-se que ela consiga alterar o estado de metilação durante a geração da nova fita, sobrepondo-se à fita hemimetilada. As enzimas com atividade de transposase, chamadas TNP1 e TNP2, foram identificadas posteriormente e são homólogas a TnpA e TnpD, respectivamente (Ping-fang, 2006; Saze *et al.*, 2012).

A superfamília *CACTA* contém elementos autônomos e não autônomos, sendo encontrada até o momento apenas em plantas. Apresenta-se como uma das principais superfamílias da classe II por proporcionar grande interação com regiões gênicas. Além da obviedade de sua interferência, inserindo-se dentro de genes e provocando silenciamento, este elemento interfere direta e indiretamente no controle da expressão de genes próximos (Wicker *et al.*, 2003; Lee *et al.*, 2012; Takahashi *et al.*, 2012).

2.4.2 Mutator

A superfamília *Mutator* ou *Mutator-like (MULEs)* foi descoberta em 1978 por Robertson quando estudava uma linhagem de milho que produzia muitas mutações espontâneas e por isso recebeu esse nome. *Mutators* possuem TIRs

variando entre 50 pb e 600 pb, que são bastante divergentes entre as famílias e se comportam como promotores destes elementos, tornando possível a transcrição nas duas direções. Classificados em autônomos e não autônomos, os TEs dessa superfamília geram TSDs entre 9 e 11 pb. São caracterizados por apresentarem apenas uma transposase chamada MuDR, assim chamada em homenagem a Don Robertson, que descobriu os primeiros elementos desta superfamília. Esta transposase é formada por dois componentes chamados MudrA (120 kDa) e MudrB (23 kDa), o primeiro semelhante a TNP de bactérias, e o segundo, sem qualquer semelhança com banco de dados, acredita-se que participe do processo de inserção e fixação do inserto, e assim como a maioria, possui o motivo DDE (Jiang *et al.*, 2004; Diao e Lisch, 2006; Wicker *et al.*, 2007; Hanada *et al.*, 2009)

Uma mesma família de *Mutator* deve possuir a mesma TIR embora possa apresentar-se bastante variável quanto às sequências internas. Isso ocorre porque as sequências internas em sua maioria não estão relacionadas à TNP e sim a fragmentos genômicos. De forma geral, elementos da superfamília *Mutator* apresentam preferência de inserção em regiões de cópia única ou baixo número de cópias e ricas em 'CG', sendo por esse motivo mais observados em regiões gênicas, mais especificamente na proximidade 5' dos genes. Além disso, certas famílias podem ser encontradas inserindo-se predominantemente na vizinhança de famílias gênicas específicas (Jiang *et al.*, 2004; Diao e Lisch, 2006; Wicker *et al.*, 2007; Hanada *et al.*, 2009).

A família mais conhecida dentro da superfamília *Mutator* é a *Pack-MULEs* (Pack derivado do inglês *package* ou pacote), nomeada assim por sua característica de capturar fragmentos de diversos loci cromossômicos. *Pack-*

MULEs podem se fundir promovendo a formação de sequências quiméricas, algumas das quais são expressas com indícios de se apresentarem de forma funcional. Em arroz, os elementos encontrados com níveis mais altos de expressão foram aqueles com maior número de fragmentos capturados (Jiang *et al.*, 2004; Diao e Lisch, 2006; Wicker *et al.*, 2007; Hanada *et al.*, 2009).

O silenciamento da superfamília *Mutator* ocorre também a partir de uma família específica de elementos nomeada de *MuK* (*Mutator Killer*), atuante como *locus* dominante. Esta se caracteriza por apresentar deleções internas em relação a elementos *Mutator* e duplicações parciais invertidas da região proteica (*MudrA* – com ~1300 pb), além da TIR da extremidade 3' que segue a sequência de *MuDR* (Figura 7). Embora a exata sequência de *MuK* não seja conhecida, tem-se sugerido que, ao ser transcrito, *MuK* produz um grampo de aproximadamente 4 kb, que é processado em sRNAs e estes provocam a inativação da transposase *MuDR* e conseqüentemente o silenciamento dos elementos (Diao e Lisch, 2006; Skibbe *et al.*, 2012).

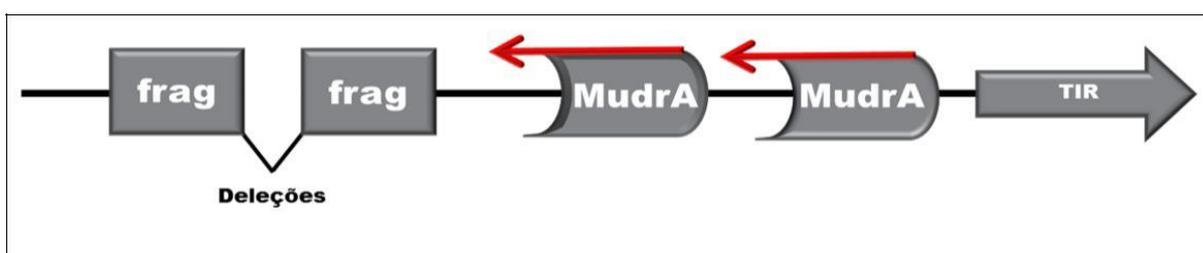


Figura 7: Figura esquemática de um elemento *MuK* (*Mutator Killer*). As linhas pretas representam porções estruturais do elemento *MuK*. A linha preta em "V" representa porção onde ocorre deleção em relação a um elemento *Mutator*. As caixas de cor cinza indicam fragmentos conservados de elementos seguidos de parte da transposase *MudrA* e sua duplicação. Setas vermelhas sinalizam sequências que foram invertidas e a seta cinza representa repetições terminais invertidas (TIRs).
Fonte: Adaptado de Skibbe *et al.* (2012).

A presença abundante de elementos *Mutator* em *Oryza sativa* L. e *Lotus japonicus* (Regel) K. Larsen indica a disseminação desse elemento antes da divisão entre mono e dicotiledôneas. A frequência de inserção de *Mutator* é geralmente mais baixa em cromossomos com grandes blocos de cromatina condensada e heterocromatina, aumentando em número e tamanho quanto mais se distanciam do centrômero *MULEs*, em arroz (Hanada *et al.*, 2009; Liu *et al.*, 2009).

2.5 Elementos transponíveis e a regulação da expressão

Assim como vários genes, os elementos transponíveis são ativados de forma diferenciada de acordo com o tecido e tempo de desenvolvimento, podendo ser induzidos por estresses ambientais (bióticos e abióticos). Além disso, TEs frequentemente carregam seus próprios promotores e possuem regimes regulatórios próprios. Sua inserção em regiões regulatórias pode acarretar em nova especificidade de tecido para o *locus* afetado ou alterar a expressão do *locus* afetado deixando-o sob o controle epigenético, que a princípio seria direcionada para o elemento transponível. No entanto, para evitar maiores danos, a maioria dos TEs está quase sempre inativa, seja por mutações que provoquem alterações na sua estrutura essencial, seja por mecanismos epigenéticos de silenciamento (Bennetzen, 2000; Bennetzen e Wang, 2014).

Entre os mecanismos epigenéticos se destacam a metilação do DNA e a ação de pequenos RNAs de interferência (siRNAs), embora sofram também influência de modificações das histonas, geralmente propondo-se a atuação de mais de um mecanismo em conjunto. Parte dessa manutenção é facilitada por

seu posicionamento em regiões não gênicas e, portanto, geralmente inativas. No entanto, diferente da maioria, alguns transposons classe II se inserem preferencialmente em regiões gênicas. Isto pode ser considerado um mecanismo apropriado que o ajuda a manter atividade suficiente e conseqüentemente o seu número de cópias, considerado baixo em relação a elementos de classe I (Bennetzen, 2000; Bennetzen e Wang, 2014).

Em plantas os siRNAs podem participar do controle transcricional (siRNA de 24 nt) e pós-transcricional (siRNA de 21 nt e 22 nt) dos TEs e de outros genes. Este processo é influenciado pelos transcritos quiméricos que podem levar ao silenciamento dos elementos transponíveis e de outros genes, a partir dos fragmentos carregados e transcritos. Estas quimeras podem ser traduzidas em proteínas que competem com as originais pelo substrato. Com conseqüências também sobre a metilação, os siRNAs são capazes de recrutar um tipo de metiltransferase responsável por rearranjar a metilação do DNA. No caso de células germinativas onde a metilação diminui drasticamente, o volume de siRNAs de 21 nt derivados de elementos transponíveis cresce para manter o controle do silenciamento e o mesmo padrão de metilação na linhagem seguinte. Já foram observados também processos de resposta positiva, como, por exemplo, quando a família *Evadé* de retrotransposon está com nível de expressão da proteína GAG mais alto, onde os siRNAs derivados do mesmo induzem a metilação *de novo* na região do domínio, que por sua vez induz a metilação do LTR e repressão gradual de sua expressão (Bennetzen, 2000; Li *et al.*, 2009; Cantu *et al.*, 2010; Saze *et al.*, 2012).

A metilação do DNA e das histonas compreende um dos principais, se não o principal, meio de silenciamento dos elementos transponíveis.

Conseqüentemente, com frequência a hipometilação acarreta o aumento da atividade de ambas as classes (I e II). Alterações na metilação de histonas específicas podem levar à ativação ou inativação de um grupo específico de TEs como, por exemplo, mutações na metiltransferase da histona H3K9, que pode elevar o nível de expressão dos *LTRs*, enquanto que danos na H3K4 desmetilase levam à transposição de não *LTRs*, exceto *LTRs* da família Tos17. Vem sendo constatado também que os TEs com tamanhos menores mostram uma tendência de se tornarem alvo de metilações assimétricas [metilações de regiões CHH (H = A, C ou T), geralmente *de novo*, promovidas pela metiltransferase CMT3]] ao contrário dos TEs de tamanhos maiores que apresentam metilação simétrica (regiões CG ou CHG, geralmente pela metiltransferase MET1) (Bennetzen, 2000; Zeh *et al.*, 2009; Cantu *et al.*, 2010; Saze *et al.*, 2012).

2.6 Bioinformática

2.6.1 Identificação e caracterização dos elementos transponíveis

O crescente número de sequências genômicas disponíveis tem se mostrado como uma potente fonte no estudo dos elementos transponíveis, reafirmando a diversidade e abrangência destes elementos, mostrando a necessidade de análises cada vez mais detalhadas sobre sua função e influência. As ferramentas bioinformáticas disponíveis até o momento permitem desde a identificação e caracterização de elementos conhecidos, como a descoberta de novas superfamílias, fornecendo dados para estudos de evolução gênica e permitindo inferências sobre suas funções. Contudo, a rápida evolução

tecnológica, a complexidade destas sequências e as dificuldades técnicas de acompanhar tais mudanças tornam o estudo desses elementos um trabalho de difícil padronização de sua anotação, gerando desafios no processo de identificação e caracterização (Lerat, 2010; Chénais *et al.*, 2012).

Dependendo do conhecimento sobre a sequência de TE que se pretende estudar, o pesquisador pode realizar buscas por homologia ou adicionar buscas baseadas em suas características estruturais. Para compreender e escolher melhor as opções devemos pensar na estrutura e natureza dos TEs (Quesneville *et al.*, 2005; Lerat, 2010).

Frequentemente, cópias de TEs pertencentes a uma mesma família apresentam estruturas que apesar de similares divergiram como consequência de sua natureza móvel e mutável. Como resultado de seus mecanismos intrínsecos, os TEs apresentam-se como mosaicos de sequências, que incluem desde suas estruturas básicas a sequências genômicas capturadas e transposons aninhados (TEs inseridos em TEs), acarretando em sequências de extrema complexidade e difíceis de serem identificadas em toda sua diversidade apenas por similaridade. Entretanto, a sua detecção pode ser melhorada utilizando grupos de referência que incorporem aspectos da sua diversidade estrutural, consensos de seus domínios e bancos de dados específicos de transposon, como o Repbase Update (RU; <http://www.girinst.org>). O RU é composto em sua maioria por consensos de famílias de TEs conhecidas, mas também por elementos completos depositados (Quesneville *et al.*, 2005; Buisine *et al.*, 2008; Lerat, 2010).

A busca baseada em homologia é ainda considerada a mais eficiente na identificação de TEs, sendo parte integrante da análise comparativa entre sequências de referência contra bancos de dados. As bibliotecas de sequências

de referência podem ser geradas pelo usuário, sendo gerada nesse caso uma biblioteca local, ou ainda bibliotecas gerais como no caso do banco RU. É importante ressaltar que essa técnica limita-se à detecção de elementos já conhecidos ou de elementos bastante similares. Adicionalmente, as ferramentas disponíveis tendem a identificá-los parcialmente ou em vários fragmentos como resultado das divergências acarretadas por sua natureza. Este problema se agrava quando as sondas envolvem cópias representativas em vez de consensos ou quando os consensos são gerados a partir de todas as cópias de uma superfamília, gerando resultados bastante distintos daqueles conseguidos com consensos de cópias pertencentes à mesma família ou compostos apenas por elementos autônomos (Tóth *et al.*, 2006; Buisine *et al.*, 2008; Lerat, 2010).

Alternativamente, visando melhores resultados na busca por homologia, Buisine *et al.* (2008) propuseram a criação de três grupos de consenso de referência: (1) com otimização do potencial codificante das ORFs contidas nos TEs; (2) que busca reunir o máximo possível da diversidade estrutural das superfamílias, e (3) que combine os dois critérios anteriores. Em seus testes foram encontrados, por exemplo, cerca de 12% a mais de sequências de transposon no genoma de *A. thaliana* (Tóth *et al.*, 2006; Buisine *et al.*, 2008; Lerat, 2010).

O RepeatMasker (RM) é uma das ferramentas mais utilizadas na identificação de elementos transponíveis por ser rápida e por permitir também a caracterização automatizada destes elementos. Inicialmente desenvolvido para mascarar sequências repetitivas facilitando a montagem e detecção de genes, o RM apresenta-se como ferramenta eficiente na identificação e caracterização das sequências repetitivas. Seu *software* inclui diversos algoritmos de alinhamento

como BLAST, Cross_match e HMM, (do inglês: *Hidden Markov Model*), utilizando RU como biblioteca padrão, porém, permite agregar outras bibliotecas de interesse no módulo UNIX (desktop). Com três arquivos básicos de saída e com opções de arquivos adicionais no módulo UNIX, o RM permite saber a descrição detalhada das repetições encontradas como tamanho, *e-value*, posição de início e fim das sequências, etc., além de listar a quantidade de elementos de cada superfamília após clusterização e um arquivo com a sequência mascarada. Outros programas já foram desenvolvidos com abordagens semelhantes como o CENSOR e o MARKERAID, outros como o PLOTREP e o GREEDIER para melhorar a seu desempenho, mas o RM se mantém entre os mais usados (Li *et al.*, 2008; Lerat, 2010; Tempel, 2012).

O perfil HMM é um algoritmo de alinhamento múltiplo e análise de resíduos conservados, promovendo análises robustas e gerando consensos que incorporam todas as probabilidades de alinhamentos. Embora largamente usado para modelagem de famílias de sequências de proteínas é raramente utilizado para famílias de sequências de DNA. De fato, os principais pacotes de *software* usados para identificação de perfis proteicos com resíduos conservados usando modelos HMM são explicitamente concentrados em modelagem de proteínas. A implementação de dois novos algoritmos por Edlefsen e Liu (2010) abrem novas possibilidades na aplicação do HMM para famílias de DNA. A condicional de Baum-Welch (CBW) é um procedimento alternativo que atualiza a parametrização para cada posição, portanto os alinhamentos não são atualizados em conjunto, garantindo a convergência ao melhor consenso. O algoritmo de Cirurgia de Modelo Dinâmico [do inglês *Dynamic Model Surgery (DMS)*], por sua vez, é responsável pela detecção de desvios. No entanto, a implementação destes

torna-se difícil uma vez que não são encontrados direcionamentos no artigo que descreve os benefícios dos mesmos (Edlefsen e Liu, 2010).

2.6.2 Bancos de dados

Atualmente os bancos de dados apresentam-se como importantes aliados das análises computacionais, armazenando e disponibilizando informações sobre todas as “ômicas” para diversas espécies. Entre os bancos de dados públicos destacam-se o NCBI (*National Center for Biotechnology Information*; <http://www.ncbi.nlm.nih.gov/>), o DDBJ (*DNA data Bank of Japan*; <http://www.ddbj.nig.ac.jp>) e o EMBL-EBI (*European Molecular Biology Laboratory – European Bioinformatics Institute*; <http://www.ebi.ac.uk>). Estes bancos de dados trabalham cooperativamente estabelecendo padrões que facilitam a submissão, a curagem e o intercâmbio de dados (Fernández-Suárez e Galperin, 2013; Nakamura *et al.*, 2013).

O Phytozome (<http://www.phytozome.net/>), por sua vez, é um projeto realizado em conjunto CIG (*Center for Integrative Genomics*) e o JGI (*United States Department of Energy – Joint Genome Institute*) que visa facilitar estudos genômicos e comparativos de plantas verdes, fornecendo acesso a 31 genomas sequenciados e anotados com o PFAM (*Protein Families Database*), KOG (*EuKaryotic Orthologous Groups*), KEGG (*Kyoto Encyclopedia of Genes and Genomes*) e PANTHER sempre que possível, usando também anotações disponíveis publicamente como RefSeq, UniProt (*Universal Protein Resource*), TAIR (*The Arabidopsis Information Resource*) e JGI (http://www.phytozome.net/Phytozome_info.php; 16 de janeiro de 2015).

Dentre os genomas completos disponíveis de leguminosas (Tabela 1) destacam-se o de *P. vulgaris* e *M. truncatula*, *Glycine max* (L.) Merr e *Lotus japonicus* (Regel) K. Larsen. Estes genomas apresentam-se como fontes de dados interessantes para análises comparativas entre as leguminosas de um modo geral, como pode ser visto por exemplo para *V. unguiculata* com *G. max*, *P. vulgaris* ou *M. truncatula* em <http://harvest.ucr.edu> (HarvEST: Cowpea).

De forma complementar, bancos de dados de acesso restrito agregam informações sobre espécies de interesse local, como o consórcio NordEST, que foi criado em 2005 objetivando análises genômicas, funcionais e estruturais do feijão-caupi (*V. unguiculata*). Com o intuito de identificar genes para melhoramento foram usadas ferramentas de genômica expressa como EST (Etiquetas de Sequência Expressa, *Expressed Sequence Target*), SuperSage e LongSage além de RNAseq. Os bancos de dados gerados totalizam 27.453 ESTs gerados pelo NordEST, 21 milhões de sequências de Super/LongSage e aproximadamente 500 milhões de sequências de RNAseq, todas associadas a estresses bióticos e abióticos (Benko-Iseppon, comunicação pessoal).

Tabela 1: Espécies da família Fabaceae com genoma completo sequenciado e nome dos sites onde estão disponíveis.

Nome científico	Nome vulgar	Iniciativa	Disponibilizado	ID (NCBI)	Referência
<i>Arachis duranensis</i> Krapov. & W.C.Gregory	Amendoim	<i>The International Peanut Genome Initiative The International Peanut Genome Initiative</i>	PeanutBase/ NCBI*	12052	-
<i>Arachis ipaensis</i> Krapov. & W.C.Gregory	Amendoim	<i>The International Peanut Genome Initiative</i>	PeanutBase/ NCBI*	35711	-
<i>Cajanus cajan</i> (L.) Millsp.	Ervilha d'angola	<i>Pigeonpea genomics initiative (PGI)</i>	PGI*/ NCBI*	2878	Singh <i>et al.</i> ,2012
<i>Cicer arietinum</i> L.	Grão-de-bico	<i>Beijing Genomics Institute</i>	NCBI*	2992	Varshney <i>et al.</i> ,2013
<i>Glycine max</i> (L.) Merr.	Soja	<i>Joint Genome Institute (JGI)</i>	NCBI*/ Phytozome	5	Schmutz <i>et al.</i> , 2010
<i>Glycine soja</i> Siebold & Zucc.	Soja selvagem	<i>The Chinese University of Hong Kong</i>	NCBI*	13239	Qi <i>et al.</i> , 2014
<i>Lotus japonicus</i> (Regel) K. Larsen	-	<i>Kazusa DNA Research Institute</i>	NCBI*	89	Sato <i>et al.</i> , 2008
<i>Lupinus angustifolius</i> L.	Tremoço-de-folhas-estreitas	<i>Department of Agriculture and Food, WA government, Australia</i>	NCBI*	11024	Yang <i>et al.</i> , 2013
<i>Medicago truncatula</i> Gaertn.	Luzerna-cortada	<i>Medicago truncatula Consortium</i>	NCBI*/Phytozome	6	Young <i>et al.</i> , 2011
<i>Phaseolus vulgaris</i> L.	Feijão comum	<i>Joint Genome Institute (JGI)</i>	NCBI*/Phytozome	380	Guo <i>et al.</i> , 2007
<i>Trifolium pratense</i> L.	Trevo vermelho	<i>Masaryk University</i>	NCBI*	11112	Istváněk <i>et al.</i> , 2014
<i>Vigna angularis</i> (Willd.) Ohwi & H. Ohashi	Feijão-azuqui	<i>Seoul National University</i>	NCBI*	11109	-
<i>Vigna radiata</i> (L.) R.Wilczek	Feijão-mungo	<i>Seoul National University</i>	NCBI*	664	-

*NCBI: National Center for Biotechnology Information; PGI: Pigeonpea Genomics Initiative)

3 OBJETIVOS

3.1 Geral

Identificar e caracterizar elementos de Classe II (*CACTA* e *Mutator*) em feijão-caupi e feijão comum, comparativamente a *Medicago truncatula* em âmbito genômico e estrutural.

3.2 Específicos

- Inferir sobre a abundância de superfamílias de TEs nos genomas de *Phaseolus vulgaris* e *M. truncatula*, completamente sequenciados, observando possíveis agrupamentos e regiões preferencias de inserção.
- Identificar e caracterizar sequências de *CACTA* e *Mutator*, reconhecendo seus domínios conservados, estruturas e genes associados.
- Realizar inferências comparativas entre a distribuição dos domínios de *CACTA* e *Mutator* e das principais sequências associadas, mediante ancoragem de sequências *Mutator* de *V. unguiculata* em cromossomos de *P. vulgaris*.
- Avaliar diferenças na constituição de sequências TEs expressas em cultivares de *V. unguiculata* contrastantes para seca em condição de estresse.

4 MATERIAL E MÉTODOS

4.1 Mineração de sequências candidatas

Foram usadas todas as sondas de *CACTA* e *Mutator* descritas por Du *et al.* (2010a) (para lista de referências das sondas ver Anexo 1) contra sequências expressas (dados de RNAseq) de ‘Pingo de Ouro’ e ‘Santo Inácio’ (cultivares contrastantes de *Vigna unguiculata*) disponíveis no banco de dados da rede NordEST (<http://bioinfo03.ibi.unicamp.br/vigna/>). O alinhamento (BLASTn) foi realizado com ponto de corte (*cut-off*) $e\text{-value} = e^{-4}$, contra ‘*Chromosomes and scaffold sequences*’ e resultado em ‘*Hit table*’.

Paralelamente, foi feita uma análise com o RepeatMasker (identificando e caracterizando todas as superfamílias de TEs) contra os genomas completos de *M. truncatula*, *P. vulgaris* e sequências disponíveis de *V. unguiculata* entre sequências de EST e sequências genômicas disponíveis no banco de dados CGKB (<http://cropgenebank.sgrp.cgiar.org/>) que montadas formaram um banco coeso que abrangeu 1/3 do genoma de *V. unguiculata*, compreendendo a porção rica em genes. O RepeatMasker foi usado em plataforma LINUX, restringindo-se apenas a genomas de plantas e excluindo-se microssatélites. (Figura 8)

4.2 Recuperação e tradução das sequências

As sequências candidatas de *CACTA* e *Mutator*, referentes às cultivares contrastantes de *V. unguiculata*, foram recuperadas e salvas em um banco de dados local com gotDNA+.java (aplicativo em Java), utilizando as coordenadas obtidas após alinhamento contra as sequências de RNAseq.

Para a tradução nos seis *frames* (quadros de leitura) foi utilizado o script “seqs_processor_and_translator_bin_V118_AGCT.py” em plataforma Linux, disponível em linguagem Python (The Compositae Genome Project, <http://cgpdb.ucdavis.edu>).

4.3 Caracterização estrutural

As sequências de *CACTA* e *Mutator*, provenientes de dados de RNAseq, foram caracterizadas mediante ferramenta CD-Search Batch disponível no NCBI, alterando o número de ‘*Maximum number of hits*’ para 100.000, com resultados curados manualmente, avaliando-se a presença de domínio das transposases e sua integridade. Adicionalmente foi observada a possível presença de domínios de outras proteínas (Figura 8).

4.4 Inferências, genoma sequenciado X genoma real

Foi calculado, para os dados provenientes do RepeatMasker, a quantidade de Mb de elementos transponíveis por Mb do genoma sequenciado a partir de uma simples divisão, como descrito abaixo, considerando ambos os valores em Mb:

$$\frac{\textit{Total TEs}}{\textit{Genoma seq.}} = \frac{\textit{x TEs}}{1}$$

Descobrimos a quantidade de TEs por Mb do genoma sequenciado este valor foi multiplicado pelo tamanho total em Mb do genoma real, como descrito abaixo:

$$\frac{\textit{x TEs}}{1} * \textit{Genoma Real}$$

Esses valores foram calculados para transposons e retrotransposons, permitindo estimar o número real de elementos no genoma de cada espécie analisada.

4.5 Ancoragem

Para as sequências de transposase de *Mutator* provenientes das cultivares contrastantes de *V. unguiculata* foi feito um alinhamento múltiplo com a ferramenta MUSCLE do MEGA 5.2 utilizando parâmetros padrões. A região mais conservada das sequências contrastantes foi alinhada usando BLASTn contra o genoma de *P. vulgaris*, disponível no Phytozome.net, sem mascarar as sequências repetitivas. Além disso, três sequências de quinases relacionadas a esta superfamília foram analisadas contra o Phytozome.net, devido à sua abundância e presença vizinhança da maioria dos clusters gênicos. Após anotação manual das posições das transposases e das quinases, uma imagem foi gerada usando as sequências anotadas do Phytozome.net com o programa CIRCOS (Krzywinski *et al.*, 2009) em plataforma LINUX. (Figura 8)

A Figura 8 apresenta um fluxograma com as etapas da metodologia realizada neste trabalho, especificadas A.

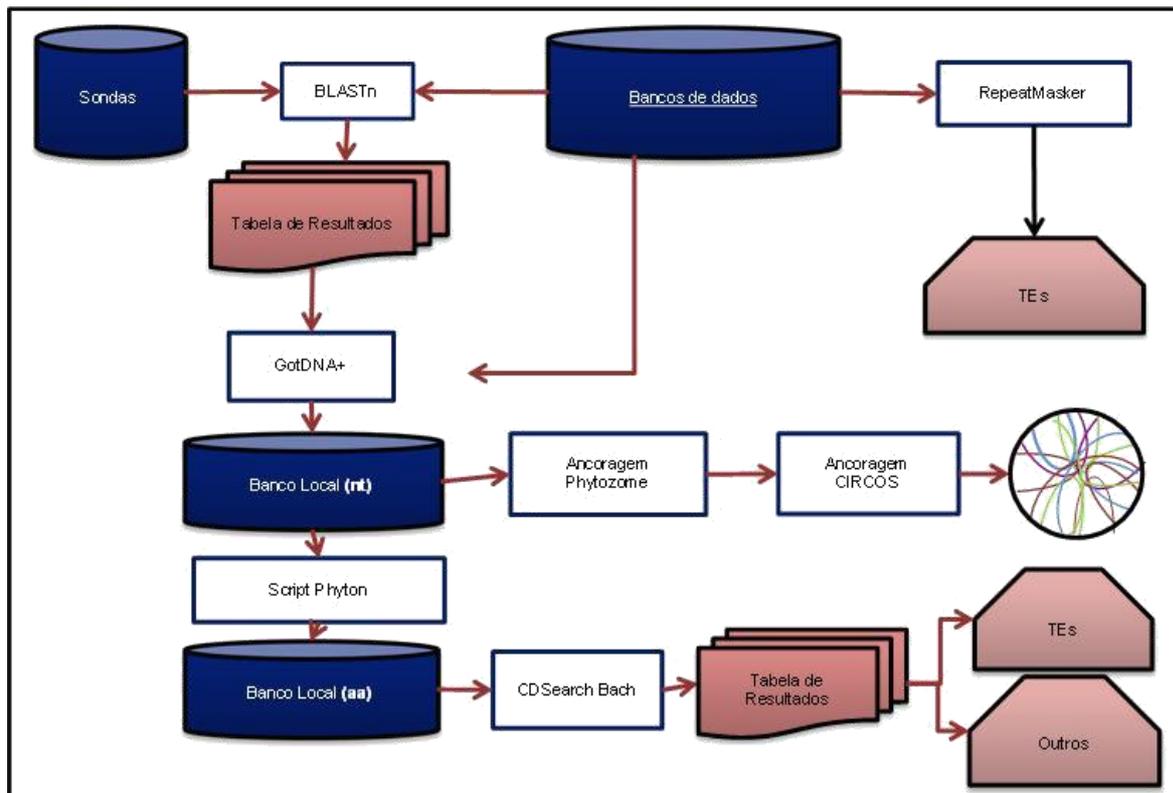


Figura 8: Fluxograma de etapas referentes à metodologia. Cilindros azuis representam os bancos de dados; quadros vazados, as ferramentas utilizadas em cada etapa; figuras avermelhadas, os resultados obtidos em cada etapa; o círculo, ilustra a imagem gerada com o CIRCOS.

5 RESULTADOS

Com o auxílio do RepeatMasker foram identificadas 54,07 Mb de sequências de TEs de *P. vulgaris* (*Pv*), sendo 5,64 Mb de DNA transposons e 48,42 Mb de retrotransposons e 65,24 Mb de TEs de *M. truncatula* (*Mt*), sendo 14,04 Mb referentes a DNA transposons e 51,20 relativos a retrotransposons. Já em 1/3 do genoma disponível de *V. unguiculata* (*Vu*) foram identificados ~11,73 Mb de sequências relativas a TEs, com DNA transposons e retrotransposons abrangendo respectivamente 2,38 Mb e 9,34 Mb. Todas as famílias de transposons já relatadas em plantas foram identificadas com o RepeatMasker, como descrito a seguir (Figura 9):

- *Phaseolus vulgaris*: Entre elementos de classe II foram observadas as seguintes famílias: *CACTA* (8333 sequências/3,01 Mb); *Mutator* (4555/1,14 Mb); *Helitron* (2185/0,63 Mb); *hAT* (1698/0,59 Mb); *PIF-Harbinger* (850/0,22 Mb) e *Tc-Mariners* (282/0,023 Mb). Dentre os elementos de classe I foram observadas: *LTRs* (83293/44,98 Mb); *LINEs* (7829/3,43 Mb) e *SINEs* (89/0,01 Mb).
- *Medicago truncatula*: Entre elementos de classe II foram identificadas: *Mutator* (32065/7,45 Mb); *PIF-Harbinger* (8550/1,49); *hATs* (6578/1,38 Mb); *Helitron* (5655/2,2 Mb); *TcMariner* (4103/0,79 Mb) e *CACTA* (4075/0,7 Mb). Quanto aos elementos de classe I, observaram-se *LTRs* (60594/39,85 Mb); *LINEs* (23983/10,95 Mb) e *SINEs* (3651/0,39 Mb).
- *Vigna unguiculata*: Entre os elementos de classe II foram identificadas: *CACTA* (6732/0,81 Mb); *Mutator* (5457/0,49 Mb); *Helitron* (3532/0,62Mb); *hATs* (1145/0,32 Mb); *TcMariner* (796/0,06 Mb) e *PIF-Harbinger* (700/0,06 Mb). Entre elementos de Classe I, observaram-se: *LTRs* (31729/9,20 Mb); *LINEs* (617/0,13

Mb); SINEs (71/0,008 Mb).

Contudo, as quantidades reais de cada classe e de cada superfamília devem ser maiores do que as observadas, levando em conta que os genomas sequenciados até o momento não representam os genomas completos, dada a dificuldade de sequenciar e montar as regiões repetitivas, nas quais também se incluem os TEs. Foi realizada uma estimativa considerando quantidades equivalentes por Mb do genoma para todas as famílias. Para *V. unguiculata* com sequenciamento de apenas 192,49 Mb disponível até esta análise e um genoma estimado de 613 Mb (Arumuganathan e Earle, 1991), estimou-se ~7,60 Mb de elementos de classe II e ~29,75 Mb de elementos de classe I. Para *P. vulgaris* com 472,5 Mb e genoma real estimado em 637 Mb (Arumuganathan e Earle, 1991) estima-se que para a classe II existam ~7,61 Mb e para classe I ~65,28 Mb. Já para *M. truncatula*, que apresenta genoma sequenciado composto por 300,6 Mb e genoma real estimado em 490 Mb (Arumuganathan e Earle, 1991) avalia-se que o número de TEs seja de aproximadamente 22,89 Mb para elementos de classe II e de 83,46 Mb para elementos de classe I. Em *M. truncatula* e *P. vulgaris* os TEs estão distribuídos em todos os cromossomos em quantidades equivalentes (Tabela 2), embora o cromossomo 6 (*Mv6* e *Pv6*) tenha apresentado uma menor quantidade em ambas as espécies.

Considerando a presença do domínio de *Mutator* para as leguminosas contrastantes de *V. unguiculata*, nenhuma sequência expressa foi encontrada com o domínio completo. Além disso, tanto para a cultivar tolerante (PO) quanto para a sensível (SI), quase todos os domínios de *Mutator* estavam incompletos na porção N-terminal, havendo apenas um fragmento em PO incompleto na porção C-terminal.

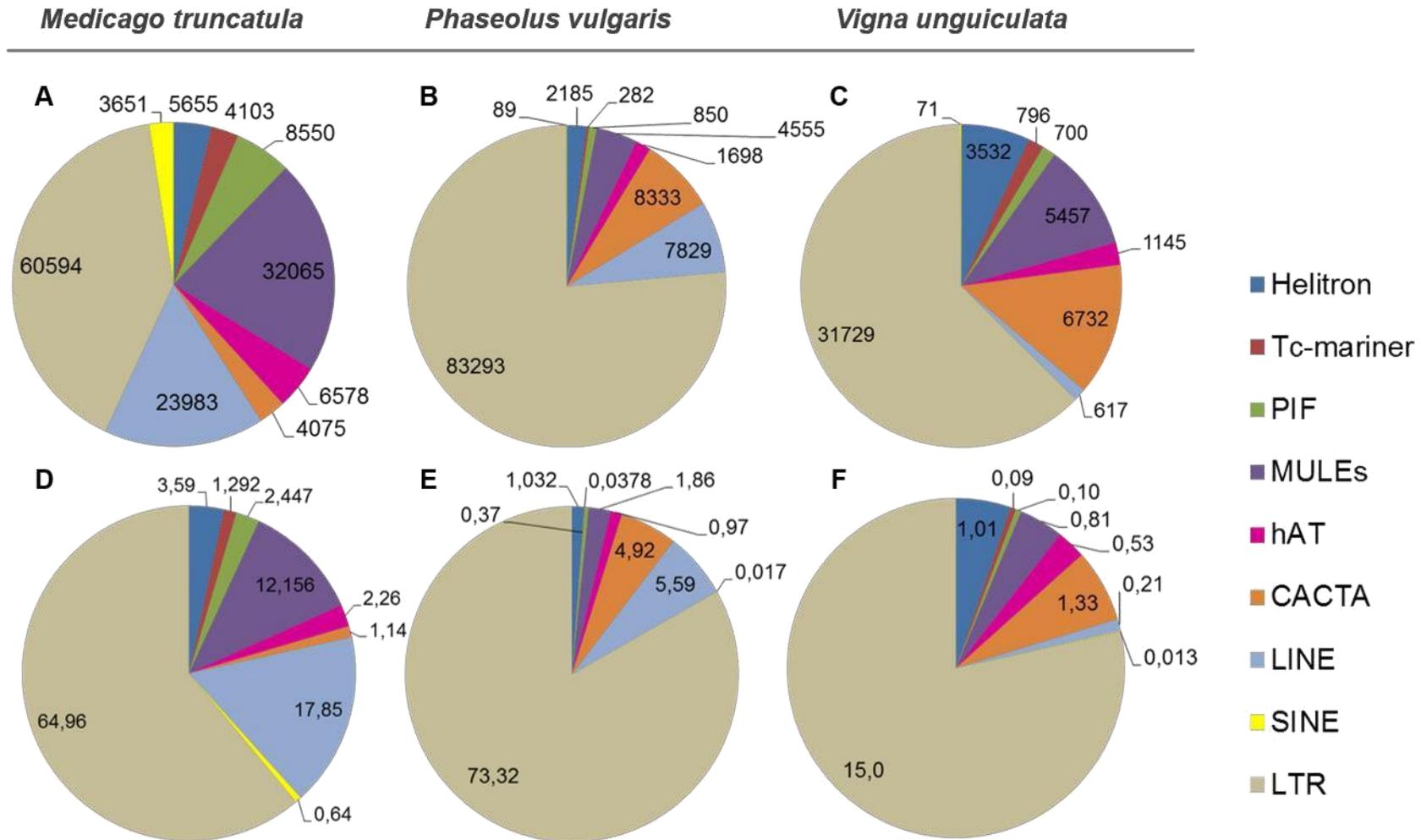


Figura 9: Representação gráfica da quantidade em número absoluto (A, B e C) e em Megabases* (Mb) (D, E e F) das superfamílias de elementos transponíveis por espécie [*Phaseolus vulgaris* (Pv), *Medicago truncatula* (Mt) e *Vigna unguiculata* (Vu)] identificados utilizando-se o RepeatMasker. Informações adicionais disponíveis no Anexo 2. *Valores em Mb estimados para o genoma total.

Tabela 2: Quantidades em números absolutos de elementos transponíveis por cromossomo identificados em *Phaseolus vulgaris* (Pv) e *Medicago truncatula* (Mt) com o programa RepeatMasker. "Scaffolds" representam sequências identificadas nos genomas sem associação a cromossomos.

CROMOSSOMO	NÚMERO DE TEs POR CROMOSSOMO	
	Pv	Mt
01	11085	10318
02	9419	12289
03	10592	14591
04	10114	15481
05	9077	17916
06	6079	7989
07	10844	14165
08	13058	10729
09	6745	0
10	9556	0
11	11053	0
Scaffolds	1492	45780

Para *CACTA* e seus três domínios (TNP2, En/Spm e TNP_assoc), as duas cultivares apresentaram-se diferentemente. Em PO, 74 domínios TNP2 apresentaram-se completos, 133 incompletos na porção C-terminal e oito incompletos em ambos os terminais. Já En/Spm não apresentou nenhuma sequência expressa com domínio completo, observando-se 44 fragmentos incompletos na porção N-terminal e 66 incompletos em ambos os terminais.

Por sua vez, TNP_assoc apresentou 41 fragmentos incompletos N-terminal e apenas dois incompletos na porção C-terminal, além de três incompletos em ambas as extremidades. Na cultivar sensível (SI), TNP2 não teve nenhum representante

completo; apresentou 33 incompletos nas duas terminações, 29 incompletos na porção C-terminal e três incompletos na terminação N (Quadro 1).

Superfamília	Domínio	Integridade	Cultivar de <i>Vigna unguiculata</i>	
			PO	SI
CACTA	TNP_Assoc	N	41	0
		C	2	2
		NC	3	0
		-	0	0
	TNP2	N	0	3
		C	133	29
		NC	8	33
		-	74	16
	Transposase_24	N	44	11
		C	0	0
		NC	66	44
		-	0	0
TOTAL			371	137
Mutator	MULE	N	1146	396
		C	1	0
		NC	0	0
		-	0	0
TOTAL			1146	396

Quadro 1: Quantidade de sequências de *Mutator* identificadas, provenientes do banco de RNAseq de *Vigna unguiculata* (banco NordEST) e integridade de seus domínios em relação às cultivares contrastantes para tolerância à seca ‘Pingo de Ouro’ (PO) e ‘Santo Inácio’ (SI).

É interessante observar a prevalência de *frames* positivos na caracterização dos domínios, mostrando que a maioria está sendo expressa na fita ‘sense’. Os domínios de transposase de *Mutator* estão todos em *frames* positivos, sendo encontrados apenas no *frame* +1 em SI e nos três *frames* em PO. CACTA, por sua vez, tem a maior parte de seus domínios em *frames* positivos, apesar de apresentar alguns domínios expressos em *frames* negativos indicando que são expressos na

fita 'anti-sense' ou que suas sequências estão invertidas, como pode ser observado no Quadro 2.

Superfamília	Domínio	Cultivar de <i>Vigna unguiculata</i>			
		PO		SI	
		Frame	Quantidade	Frame	Quantidade
Mutator	MULE	fr1	1098	fr1	396
		fr2	24	fr2	0
		fr3	25	fr3	0
		fr4	0	fr4	0
		fr5	0	fr5	0
		fr6	0	fr6	0
CACTA	Transpos_assoc	fr1	2	fr1	0
		fr2	39	fr2	0
		fr3	3	fr3	0
		fr4	0	fr4	0
		fr5	2	fr5	0
		fr6	0	fr6	2
	TNP2	fr1	81	fr1	21
		fr2	5	fr2	12
		fr3	120	fr3	4
		fr4	0	fr4	17
		fr5	3	fr5	18
		fr6	6	fr6	9
	Transposase_24	fr1	0	fr1	0
		fr2	110	fr2	55
		fr3	0	fr3	0
		fr4	0	fr4	0
		fr5	0	fr5	0
		fr6	0	fr6	0

Quadro 2: Distribuição dos domínios por *frame* de leitura. Distribuição dos domínios de tranposase de *Mutator* e *CACTA* por *frame* de leitura (*Open Reading Frame*).

Com essa análise também foi possível observar que as sequências gênicas relacionadas à superfamília *Mutator* nas duas cultivares também foram bem distintas. Em PO, os quatro grupos de sequências mais encontrados relacionando-

se com *Mutator* foram quinases, *Plant Mobile Domain* (PMD), cinesinas e transcriptases reversas (RT), nessa ordem (Figura 10). Já em SI, apesar das quinases e PMD figurarem também em primeiro e segundo lugar, em sequência estavam as proteínas chaperoninas e proteínas de ligação ao DNA (*DNA-binding*) (Figura 11). Para a superfamília *CACTA*, o domínio mais encontrado nas duas cultivares foi a família gênica ABC. Entretanto em PO, os domínios DUF4216 (proteína de função desconhecida geralmente encontrada próxima a TNP2), ATPase e transportador de cobalto se seguiram como mais abundantes (Figura 12), enquanto que em SI foram encontrados os domínios transportador de cobalto, transportador de fosfato e chaperoninas, em ordem de abundância (Figura 13).

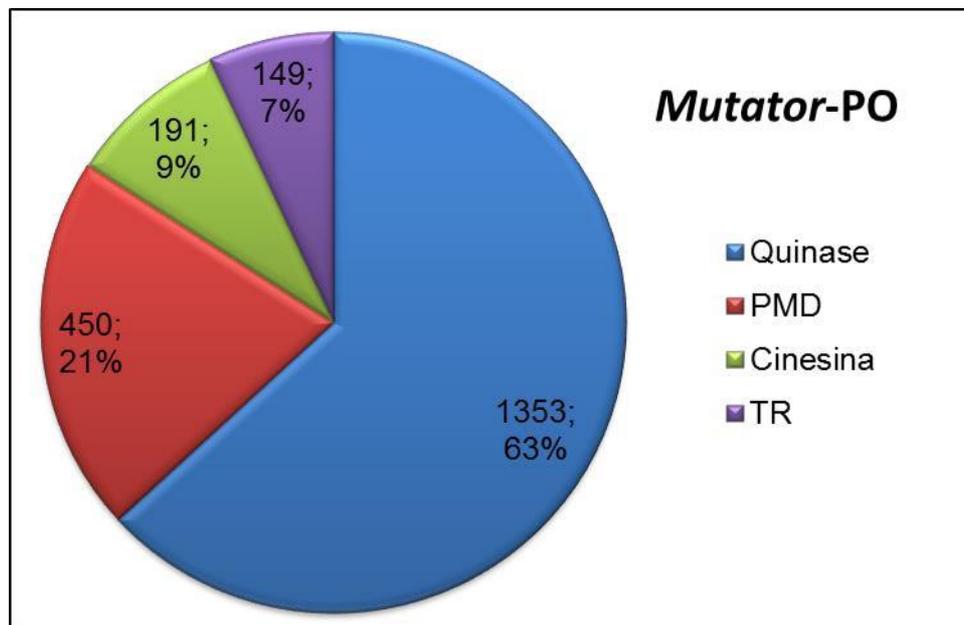


Figura 10: Domínios gênicos mais abundantes encontrados em associação com transposons *Mutator* na cultivar tolerante 'Pingo de Ouro' de *Vigna unguiculata* [1353 quinases; 450 PMD; 191 cinesinas; 149 transcriptases reversas (TR)].

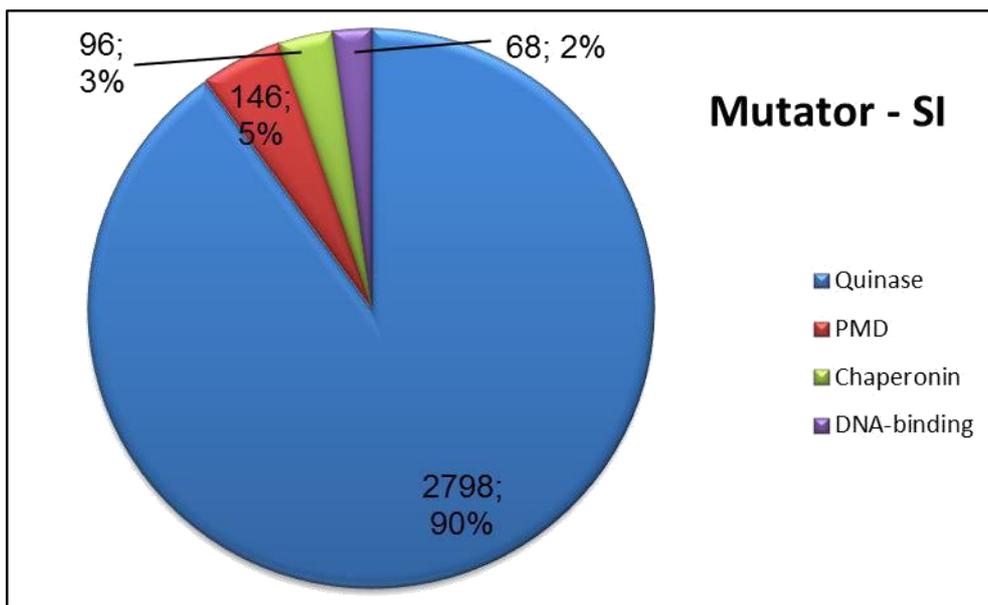


Figura 11: Domínios gênicos mais abundantes encontrados em associação com transposons *Mutator* na cultivar sensível 'Santo Inácio' de *Vigna unguiculata* [2798 quinases; 146 PMD; 96 chaperoninas; 68 domínios de ligação ao DNA (*DNA-binding*)].

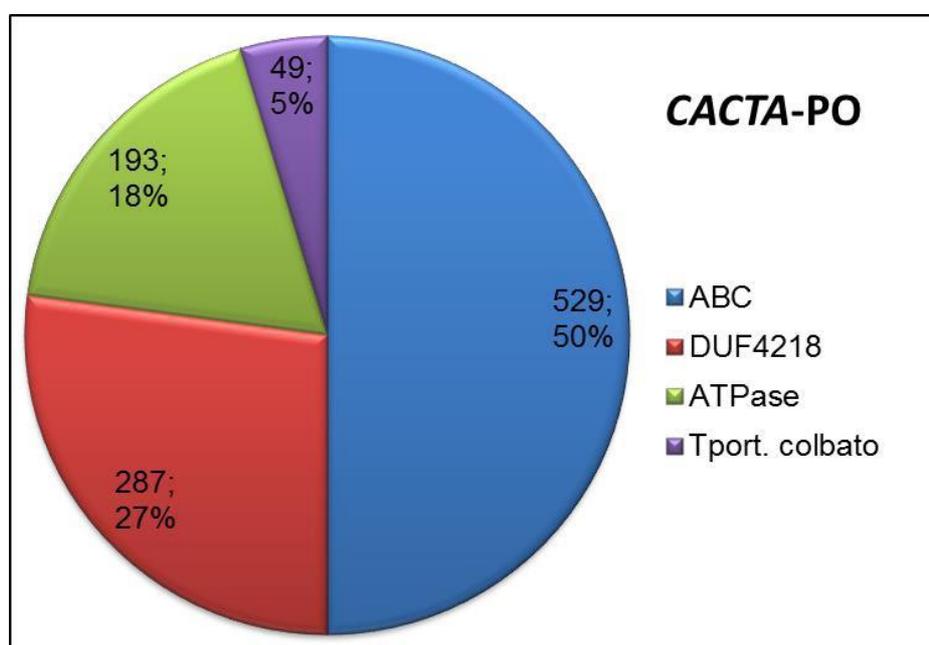


Figura 12: Domínios gênicos mais abundantes encontrados em associação com transposons *CACTA* na cultivar tolerante 'Pingo de Ouro' de *Vigna unguiculata* [529 ABCs; 287 DUF4218 (proteína de função desconhecida); 193 ATPases; 49 transportadores de cobalto (Tport → transportador)].

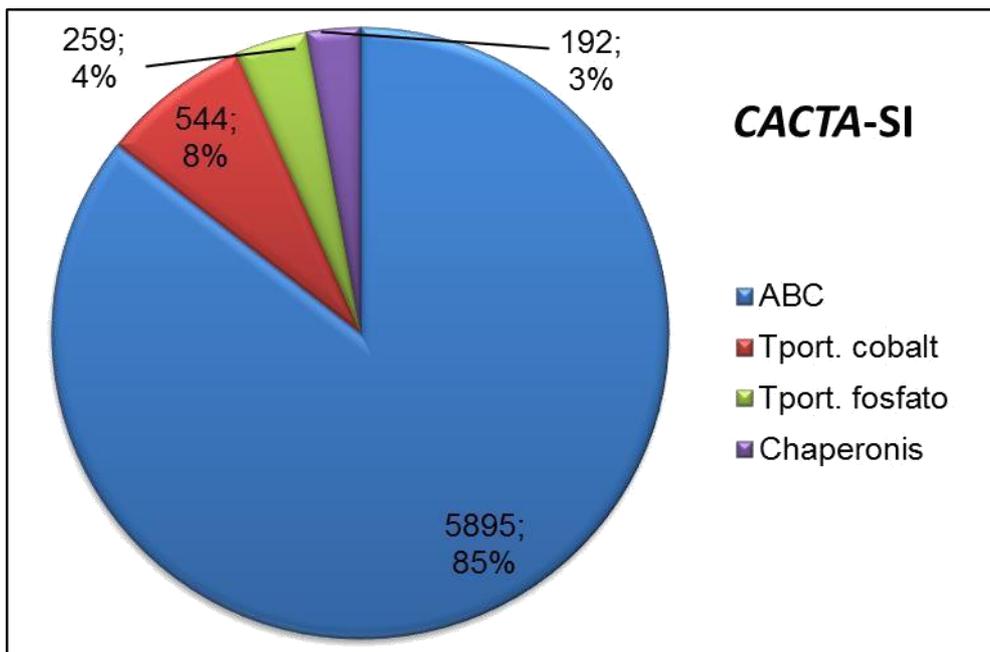


Figura 13: Domínios gênicos mais abundantes encontrados em associação com transposons *CACTA* na cultivar sensível ‘Santo Inácio’ de *Vigna unguiculata* [529 ABCs; 287 DUF4218 (proteína de função desconhecida); 193 ATPases; 49 transportadores de cobalto (Tport → transportador)].

De acordo com o alinhamento das sequências expressas referentes ao domínio da TNP de *Mutator* no MEGA 5.2, tanto PO como SI apresentaram uma região bastante conservada, de maior extensão em PO (Figura 14). Adicionalmente, destacou-se um polimorfismo de um único nucleotídeo na região conservada entre ambas as cultivares, observando-se uma timina (T) em PO e uma citosina (C) em SI.

Considerando os domínios referentes às transposases de *Mutator* e das três quinases, a ancoragem das sequências de *V. unguiculata* no genoma de *P. vulgaris* com o programa CIRCOS (Figura 15) permitiu observar que as sequências provenientes de PO e de SI apresentaram diferenças em sua localização. Contudo, ambas mostraram proximidade com sequências de quinase, domínio mais abundante dentre os relacionados a este elemento em *V. unguiculata*. Dessa forma,

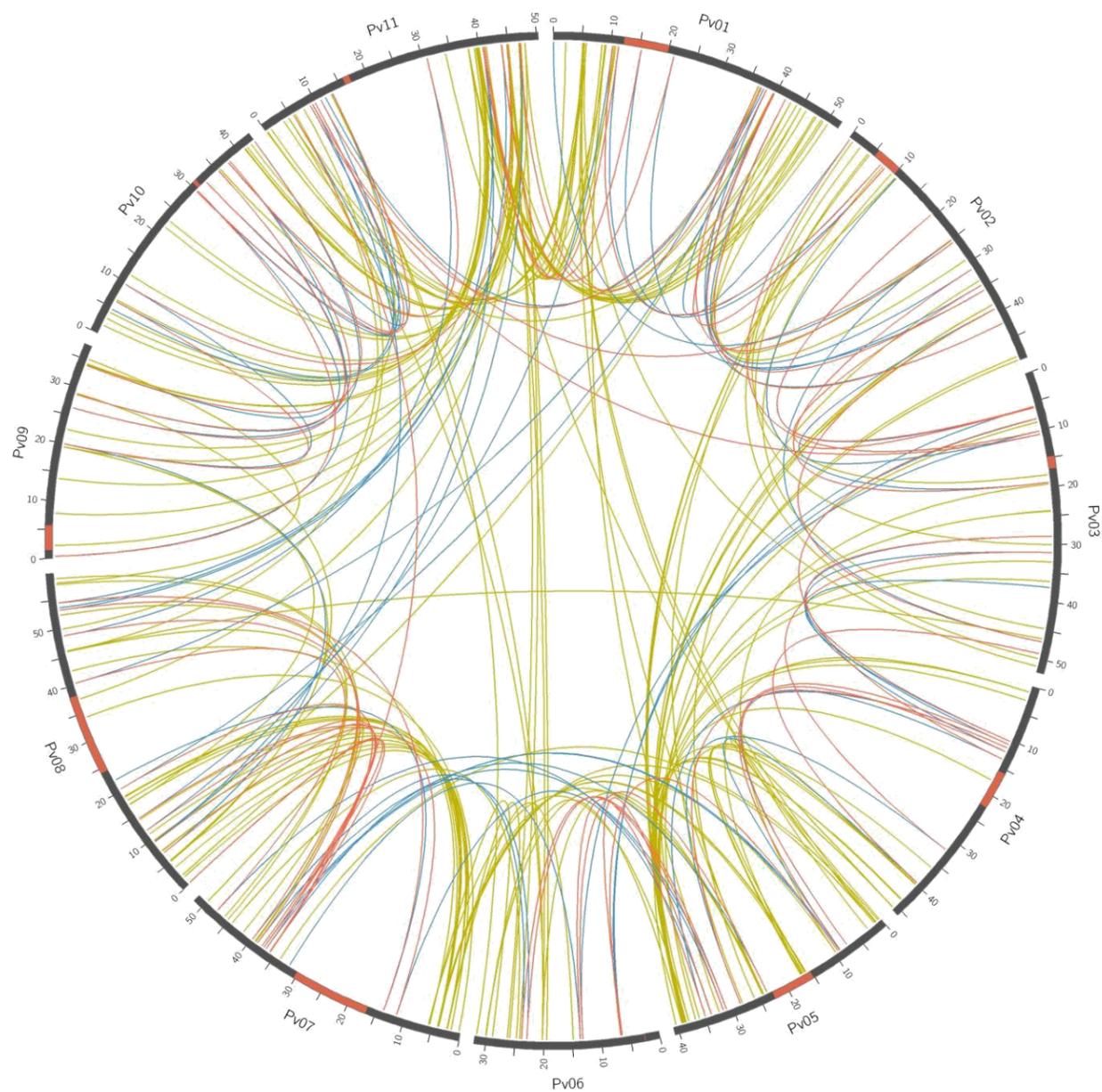


Figura 15: Disposição cromossômica das transposases (TNPs) e quinases de *Vigna unguiculata* nos cromossomos de *Phaseolus vulgaris* (em azul claro). Linhas ligantes da mesma cor representam a mesma sequência ligada a seus semelhantes; linhas amarelas representam quinases; linhas azuis TNPs de *Mutator* de 'Pingo de ouro' (cultivar de *V. unguiculata* tolerante à seca); linhas vermelhas TNPs de *Mutator* de 'Santo Inácio' (cultivar de *V. unguiculata* sensível à seca).

6 DISCUSSÃO

6.1 Abundância de TEs nos genomas

Conhecer o conteúdo de elementos transponíveis (TEs) e sua localização no genoma é etapa crucial para os que pretendem compreender a estrutura e a dinâmica genética de um organismo. Para as três leguminosas aqui estudadas (*M. truncatula*, *P. vulgaris* e *V. unguiculata*), assim como para outras dicotiledôneas e monocotiledôneas relatadas na revisão de Oliver *et al.* (2013), foi constatada a prevalência de retrotransposons em relação aos transposons, justificada por seu mecanismo de expansão, capaz de gerar diversas cópias no genoma (Wicker *et al.*, 2007).

Contudo, as quantidades de TEs encontrados foram diferentes daqueles encontrados por Young *et al.* (2011) para *M. truncatula* (*Mt*), apesar da utilização de metodologia semelhante (RM). As quantidades de TEs identificadas no presente trabalho para cada superfamília foram, em geral, cerca de duas vezes maiores que no trabalho anterior, com exceção de: *PIF*, que apresentou quantidade igual (em torno de 0,5%); *Tc-mariner*, que não foi previamente identificado (0,26% no presente trabalho) e os *LTRs*, que foram observados em 13,2% do genoma (64,96 Mb) neste trabalho, reportados anteriormente em 24,2% pelos referidos autores. A diferença entre as duas pesquisas deve-se provavelmente a diferenças nos bancos de dados utilizados e às atualizações sofridas no que tange aos elementos depositados nesses bancos. Contudo, no caso do *LTR*, esta diferença deve ter ocorrido principalmente devido ao uso da ferramenta LTR-STRUC pelos referidos autores.

Para *P. vulgaris*, Schumutz *et al.* (2014) apresentaram valor total de sequências maior para todas as superfamílias de classe II e de elementos não *LTRs*. Contudo, a metodologia adotada pelos autores baseada apenas em BLAST é considerada menos robusta que o RM por não incluir uma curagem incluindo a identificação de falsos positivos e por não promover a clusterização de sequências diminuindo a redundância (Buisine *et al.*, 2008; Li *et al.*, 2008), sendo provavelmente essa a razão dos diferentes resultados. Com relação aos elementos *LTR*, foram identificados cerca de 36,6% do genoma (Schumutz *et al.*, 2014), enquanto no presente trabalho sua abrangência foi de aproximadamente 15% (14,96 Mb). É provável que a diferença deva-se principalmente ao uso do LTR_Finder por Schumutz *et al.* (2014), que promove a busca *de novo* desta superfamília. Contudo, os dois trabalhos concordam que para esta leguminosa as duas superfamílias mais abundantes de transposons são *CACTA* e *Mutator*. Por outro lado, o trabalho de Schumutz *et al.* (2014) não reporta nenhum elemento *Tc-mariner* ao contrário do que foi identificado neste trabalho (0,0077% do genoma, equivalente 0,04 Mb).

Ao comparar *Mt* (25%) e *Pv* (21%) com outros dados de dicotiledôneas se observa que ambas possuem quantidade de TEs dentro da variação encontrada para as dicotiledôneas, entre 12% e 63%, embora com valor abaixo da média (37,98%) (Tabela 3). Ao contrário do esperado, *M. truncatula* apresentou maior cobertura de TEs que *P. vulgaris* e *V. unguiculata* a despeito do seu genoma menor. Em *V. unguiculata*, por sua vez, a quantidade de TEs estimada para seu genoma foi baixa, comparável à de *Arabidopsis thaliana* (The *Arabidopsis* Initiative, 2000), que possui um genoma aproximadamente quatro vezes menor que o seu [145 Mb e 613 Mb respectivamente - Arumuganathan e Earle (1991)]. Isto ocorreu provavelmente pelos dados de *V. unguiculata* estarem subestimados, devido à porção do genoma

analisado, que consistia em regiões hipometiladas, ou seja, ricas em genes, na qual se espera baixa quantidade de TEs.

Para todos os organismos comparados, os *LTRs* figuraram como os mais abrangentes dentre todos os TEs. Os genomas de *P. vulgaris* (64,96 Mb – 14,96% do genoma total) e *M. truncatula* (73,32 Mb – 13,25%), por exemplo, apresentaram quantidades similares de *LTRs* identificados, mostrando que a maior diferença entre estes dois genomas está na cobertura de elementos de classe II. Por outro lado, as quantidades de *LTRs* identificadas para as três leguminosas estudadas (*Vu*, *Pv* e *Mt*) são menores que as encontradas na maioria dos organismos como *Cajanus cajan* – 22% (Singh *et al.*, 2012), *Trifolium pratense* – 21% (Ištvánek *et al.*, 2014), *G. max* ~42% (Schumutz *et al.* 2010), *Cicer arietinum* ~45% (Varshney *et al.* 2013), *Triticum aestivum* ~61% de *LTRs* (Oliver *et al.*, 2013) e *Zea mays* ~69%, considerado um dos organismos eucariotos com maior quantidade de TEs (Oliver *et al.*, 2013).

Com relação aos transposons, embora o genoma de *M. truncatula* seja quase a metade dos genomas de *Solanum lycopersicum* e de *Solanum tuberosum*, a quantidade de transposons foi bem maior, com 22,89 Mb abrangendo 4,67%, enquanto em *S. lycopersicum* e em *S. tuberosum* foi, respectivamente, de 0,9% e 1,2%. Também, *A. thaliana*, espécie com genoma menor que as três espécies anteriores [145 Mb - Arumuganathan e Earle, (1991)], possui quantidade de transposons (aproximadamente 6 Mb) (The *Arabidopsis* Initiative, 2000) maior que a maioria das espécies apresentadas na Tabela 3.

Tabela 3: Comparação da composição de elementos transponíveis (TEs) entre angiospermas. Em vermelho estão destacados os dados conseguidos no presente trabalho.

Espécie	Genoma total (Mb)	Transp. (%)	Retrotra. (%)	Total TEs (%)	Referência
Dicotiledôneas					
<i>Arabidopsis</i>	145	6,00	6,00	12,00	The <i>Arabidopsis</i> Initiative, 2000
<i>Cajanus cajan</i>	511	2,99	23,59	26,58	Singh <i>et al.</i> , 2012
<i>Cicer arietinum</i>	738	9,32	45,64	54,96	Varshney <i>et al.</i> , 2013
<i>Fragaria vesca</i>	240	5,20	14,70	19,90	Oliver <i>et al.</i> , 2013
<i>Glycine max</i>	1115	17,00	42,00	59,00	Schumutz <i>et al.</i> , 2010
<i>Malus x domestica</i>	742	0,90	37,60	38,50	Oliver <i>et al.</i> , 2013
<i>Medicago truncatula</i>	490	2,48	22,40	24,88	Presente trabalho
<i>Phaseolus vulgaris</i>	637	1,87	19,19	21,06	Presente trabalho
<i>Solanum lycopersicum</i>	900	0,90	62,30	63,20	Oliver <i>et al.</i> , 2013
<i>Solanum tuberosum</i>	1597	1,20	53,20	54,40	Oliver <i>et al.</i> , 2013
<i>Trifolium pratense</i>	468	6,70	23,81	30,51	Ištvánek <i>et al.</i> , 2014
<i>Vigna unguiculata</i>	613	0,79	4,28	5,07	Presente trabalho
<i>Vitis vinifera</i>	483	1,40	19,40	20,80	Oliver <i>et al.</i> , 2013
Monocotiledôneas					
<i>Brachypodium distachyon</i>	272	4,80	23,30	28,10	Oliver <i>et al.</i> , 2013
<i>Hordeum vulgare</i>	5100	5,00	52,70	57,70	Oliver <i>et al.</i> , 2013
<i>Musa acuminata</i>	873	1,30	42,40	43,70	Oliver <i>et al.</i> , 2013
<i>Oryza sativa</i>	410	13,70	25,80	39,50	Oliver <i>et al.</i> , 2013
<i>Setaria itálica</i>	423	9,40	31,60	41,00	Oliver <i>et al.</i> , 2013
<i>Sorghum bicolor</i>	748	7,50	54,50	62,00	Oliver <i>et al.</i> , 2013
<i>Triticum aestivum</i>	2118	14,90	63,70	78,60	Oliver <i>et al.</i> , 2013
<i>Zea mays</i>	2300	8,60	75,60	84,20	Oliver <i>et al.</i> , 2013

Dentre os transposons encontrados no presente trabalho, *CACTA* foi a superfamília mais representativa em *Pv* e *Vu*, seguida respectivamente por *Mutator* e *Helitron*. Diferentemente, em *Mt*, *Mutator* figurou como elemento de maior cobertura seguido por *Helitron*. Complementarmente, de acordo com a revisão de Oliver *et al.* (2013), *CACTA* predomina também em *Fragaria vesca* e *Mutator* em *S. lycopersicum* e *S. tuberosum*. Em *A. thaliana* (The *Arabidopsis* Initiative, 2000) e *C.*

cajan (Singh *et al.*, 2012) *Mutator* e *hAT* figuram, respectivamente, como elementos mais abundantes seguidos, por *CACTA* e *Mutator*, respectivamente. Para *G. max* (17% de transposons, Schmutz *et al.* 2010), por sua vez, as superfamílias *Tc-mariner*, *hAT*, *Mutator*, *PIF*, *Pong*, *CACTA* e *Helitrons* foram identificadas. Dentre essas, a superfamília *Pong* foi a única não identificada em nenhuma das três leguminosas (*Mt*, *Pv* e *Vu*) no presente trabalho, contudo, como sugerido por Zhang *et al.* (2004), *Pong* está incluída na superfamília *PIF*.

Com exceção de *Musa acuminata*, em sete das oito monocotiledôneas (*T. aestivum*, *Hordeum vulgare*, *Z. mays*, *Sorghum bicolor*, *Oryza sativa*, *Brachypodium distachyon* e *Setaria italica*) apresentadas na revisão de Oliver *et al.* (2013), a prevalência foi dos elementos *CACTA*, com *Mutator* figurando em segundo lugar para quatro delas (*T. aestivum*, *H. vulgare*, *O. sativa*, *B. distachyon*), e em terceiro em três delas (*Z. mays*, *S. bicolor* e *S. itálica*), sugerindo que *CACTA* e *Mutator* apesar de não constituírem uma regra quanto a sua predominância apresentam uma tendência natural a se manter entre as superfamílias mais presentes entre as angiospermas.

6.2 Dispersão no genoma

Tanto os elementos de classe I quanto os de classe II foram identificados por todos os cromossomos tanto para *M. truncatula* quanto para *P. vulgaris*. Contudo, enquanto em *Pv* as sequências não ancoradas no genoma (*scaffolds*) apresentaram o número mais baixo de TEs, para *Mt* esses *scaffolds* apresentaram a maior quantidade de TEs identificada. Entre os pseudocromossomos analisados em *Pv* o cromossomo *Pv8* apresentou a maior quantidade de sequências e *Pv6* a menor. De

modo semelhante, em *Mt* o cromossomo *Mt6* figurou com a menor quantidade de sequências, embora *Mt5* tenha apresentado a maior quantidade. Em *C. cajan*, de acordo com Singh *et al.* (2012), o cromossomo *Cc10* mostrou o maior número de elementos e *Cc11* o menor e, de modo semelhante a *Pv*, as sequências *scaffold* agruparam ínfima parte dos elementos identificados.

Os elementos *Mutator* de *V. unguiculata* observados em cromossomos de *P. vulgaris* localizaram-se preferencialmente na região distal dos cromossomos, região de eucromatina, assim como observado para *Mutator* em milho (Liu *et al.* 2009) e *Pack-MULEs* em arroz (Jiang *et al.*, 2004). Este fato está relacionado a uma inserção preferencial destes elementos na vizinhança de genes (regiões de eucromatina), tanto para a superfamília *Mutator* quanto para *CACTA* e *Helitrons* (Wicker *et al.*, 2007)

Por outro lado, outras superfamílias, principalmente as de classe I, parecem apresentar distribuição preferencialmente pericentromérica, associadas aparentemente a regiões de heterocromatina. Em *Jatropha curcas*, por exemplo, a distribuição de elementos *Gypsy-like* foi preferencialmente em regiões pericentroméricas com retrotransposons também acumulados em regiões centroméricas e teloméricas (Alipour *et al.*, 2014). De modo semelhante, as regiões heterocromáticas pericentroméricas e centroméricas de *O. sativa*, *P. vulgaris*, *G. max* e *V. unguiculata* mostraram-se ricas em TEs (Bortoleti, 2010; Du *et al.*, 2010b; Bortoleti *et al.*, 2012).

6.3 Expressão

Ao analisar as sequências expressas das cultivares contrastantes de *V. unguiculata* que foram alinhadas contra o genoma de *P. vulgaris*, permitindo visualizar a posição de sequências similares, pôde-se confirmar a sua proximidade com quinases, formando pequenos *clusters*. Estes possivelmente interferem na expressão das quinases, aumentando ou reduzindo sua expressão. Em arroz, por exemplo, sequências mais próximas de TEs tendem a ser superexpressas em relação às mais distantes (Hanada *et al.*, 2009), enquanto, em três cultivares de *A. thaliana*, TEs reduziram a expressão de genes próximos. Adicionalmente, em uma avaliação mais detalhada em *A. thaliana*, foi observado que genes próximos aos elementos *CACTA* tiveram maior expressão que aqueles próximos a retrotransposons (Wang *et al.*, 2013).

É interessante observar que *V. unguiculata* cv. PO (cultivar tolerante) apresentou uma lista muito maior de genes relacionados a *Mutator* do que cv. SI (cultivar sensível), com 853 domínios e 166 domínios, respectivamente além de maior número de sequências expressas de *Mutator* e *CACTA*. Entre os domínios relacionados, as quinases se destacaram para elementos *Mutator* e domínios da família de proteínas ABC para elementos *CACTA*. Foi possível também observar associações com fatores de transcrição, sequências ribossomais, proteínas características de bactérias, entre outras. Considerando as respostas mais eficientes aos estresses na cultivar tolerante e o aumento da expressão desses elementos, é possível supor que estes atuem como ativadores ou promotores de genes, melhorando a expressão de genes interessantes para o momento da planta e

incrementando seus mecanismos de resposta à seca, ao contrário do observado para a cultivar sensível.

Em arroz, os elementos *MULEs* autônomos também tenderam a ser mais expressos em raízes e folhas jovens, durante estresse de seca e salinidade, bem como em raízes jovens e em folhas no caso de estresse de frio (4°C). Por outro lado, *Pack-MULEs* foram mais expressos em raízes e folhas não estressadas e em tecido meristemático (Hanada *et al.*, 2009). Em milho, 20 famílias de TEs superexpressas estão associadas a genes responsivos a estresse abiótico tais como salinidade, frio, raios UV e alguns respondendo a mais de um estresse, enquanto apenas três famílias de TEs estavam associadas a genes inibidos (Makarevitch *et al.*, 2014). Adicionalmente, vale destacar que no presente estudo *Mutator* e os três domínios *CACTA* apresentaram consideravelmente mais transcritos em *frames* positivos, indicando sua transcrição na fita 'sense'. Em arroz, os *Pack-MULEs* expressos em material não estressado, os éxons também eram em sua maioria da fita 'sense', apesar da possibilidade de transcrição nas duas direções

Como dito anteriormente, além da abundância de quinases (associadas a *Mutator*) e ABCs (associadas a *CACTA*), uma grande diversidade de genes foi observada relacionada a essas duas superfamílias estudadas em *V. unguiculata* (como pode ser visto para *Mutator* no Anexo 4), refletindo na influência e na dinâmica genômica das espécies. Este fato reforça a ideia de que todas as famílias gênicas relacionadas a TEs, dependendo de sua proximidade, podem sofrer influências regulatórias, especialmente em situações de estresse. Uma das formas de interferência nesta dinâmica ocorre quando o TE funciona como promotor para o gene (Cui e Cao, 2014; Makarevitch *et al.*, 2014), podendo ainda gerar pequenos

RNAs (sRNAs) que podem interferir no controle pré e pós-transcricional de outros genes (Cantu *et al.*, 2010).

Com base nos presentes resultados, outras perguntas podem ser feitas. Essas sequências podem influenciar o sistema de defesa contra patógenos ou a resposta a outros tipos de estresse? Elas podem promover a manutenção de quantidades mínimas de siRNAs capazes de silenciar determinados genes como resposta a situações de estresse? Estudos funcionais adicionais incluindo transformação genética, análises integradas das respostas em nível transcricional, pós-transcricional e de proteômica podem lançar luz a estas questões, possivelmente confirmando o papel funcional relevante desses elementos na evolução e no funcionamento dos seres vivos. A julgar pelo observado em feijão-caupi, a prevalência de transcritos na cultivar tolerante pode ser um indicativo da relevância desses elementos como fatores promotores de uma maior adaptabilidade ao estresse hídrico.

7 CONCLUSÕES

- 7.1** Os elementos *CACTA* e *Mutator* apresentam-se entre os mais abundantes elementos de classe II nos genomas das espécies *M. truncatula* e *P. vulgaris*.
- 7.2** Os *LTRs* são as sequências mais abundantes dentre todos os TEs nos genomas das três espécies estudadas (*V. unguiculata*, *M. truncatula* e *P. vulgaris*).
- 7.3** Ao contrário do esperado, *M. truncatula*, a despeito de seu genoma menor, apresenta maior número de elementos identificados, indicando maior tolerância às inserções ou menor controle sob atividade de TEs.
- 7.4** As sequências de *Mutator* ancoradas em *P. vulgaris* apresentam aglomerados nas regiões eucromáticas, ricas em genes, sendo que esta superfamília figura como uma das mais abundantes em *V. unguiculata*, corroborando com a sua preferência por regiões gênicas.
- 7.5** As sequências de *Mutator* ancoradas em *P. vulgaris* se mantêm preferencialmente próximas a quinases. Contudo, não há predileção por nenhum cromossomo quanto à sua inserção de forma geral.
- 7.6** As diferenças entre as cultivares contrastantes de *V. unguiculata* sugerem uma influência positiva de *Mutator* e *CACTA* dada a diferença em sua expressão entre PO e SI. A resposta de PO parece manter em equilíbrio a atividade destes elementos de forma a beneficiá-la ou não interferir em seus mecanismos de defesa e nas funções vitais em condições de estresse.

7.7 A grande diversidade de domínios relacionados a *CACTA* e *Mutator* nessas cultivares são indícios do alcance e da influência desses elementos sob outras sequências.

REFERÊNCIAS

- Alipour A, Tsuchimoto S, Sakai H, Ohmido N e Fukui K (2013) Structural characterization of *Copia-type* retrotransposon leads to insights into the marker development in a biofuel crop, *Jatropha curcas* L. *Biotechnol Biofuels*. 6:129-141
- Alipour A, Cartagena JA, Tsuchimoto S, Sakai H, Ohmido N e Fukui K (2014) Identification and characterization of novel *Gypsy-Type* retrotransposons in a biodiesel crop, *Jatropha curcas* L. *Plant Mol Biol*. 32:923-930
- Arensburger P *et al.*, (2011). Phylogenetic and functional characterization of the hAT transposon superfamily. *Genetics* 188(1):45–57
- Arumuganathan K e Earle ED (1991) Nuclear DNA content of some important plant species. *Plant Mol Biol Rep*. 9(4):208-218
- Bennetzen JL (2000). Transposable element contributions to plant gene and genome evolution. *Plant Mol Biol*. 42(1):251–269
- Bennetzen JL e Wang H (2014). The Contributions of Transposable Elements to the Structure, Function, and Evolution of Plant Genomes. *Anu Rev Plant Biol*, 65, 505–530
- Biémont C e Vieira C (2006). Genetics: junk DNA as an evolutionary force. *Nature*. 443(7111): 521–524
- Bortoleti KCA (2010) Mapeamento cromossômico comparativo em espécies de *Glycine* Willd., *Phaseolus* L. e *Vigna* Savi. Tese (Doutorado em Genética) – Universidade Federal de Pernambuco, Centro de Ciências Biológicas (CCB).
- Bortoleti KCA, Benko-Iseppon AM, Melo NF e Brasileiro-Vidal AC (2012) Chromatin differentiation between *Vigna radiata* (L.) R. Wilczek and *V. unguiculata* (L.) Walp. (Fabaceae). *Plant Syst Evol*. 298(3):689-693
- Brito R, Lopes HM, Fernandes MCA, Aguiar LA E Ceará PS (2013). Avaliação da qualidade fisiológica e sanitária de sementes de feijão-vagem (*Phaseolus vulgaris* L). produzidas sob manejo orgânico e submetidas ao congelamento. *ABA-Agroecologia*. 8(3):131–140
- Buisine N, Quesneville H, e Colot V (2008). Improved detection and annotation of transposable elements in sequenced genomes using multiple reference sequence sets. *Genomics*. 91(5):467–75
- Cantu D *et al.*, (2010). Small RNAs, DNA methylation and transposable elements in wheat. *BMC Genomics*. 11:408–423
- Chen S e Li X (2008). Molecular characterization of the first intact Transib transposon from *Helicoverpa zea*. *Gene*. 408(1-2): 51–63

- Chénais B, Caruso A, Hiard S e Casse N (2012). The impact of transposable elements on eukaryotic genomes: from genome size increase to genetic adaptation to stressful environments. *Gene*. 509(1):7–15
- Conab: Companhia Nacional de Abastecimento. (2015). Acompanhamento da Safra Brasileira Grãos. V2(N4):1-95
- Cook DR (1999). *Medicago truncatula* - a model in the making! *Curr Opin Plant Biol*. 2: 301–304
- Cui X e Cao X (2014) Epigenetic regulation and functional exaptation of transposable elements in higher plants. *Plant Biol*. 21:83-88
- Diao, X-M e Lisch D (2006). *Mutator* transposon in maize and MULEs in the plant genome. *Acta Genetica Sin*. 33(6):477–87
- Dong Y *et al.*, (2011). Structural characterization of helitrons and their stepwise capturing of gene fragments in the maize genome. *BMC Genomics*, 12(1):609
- Dos-Santos A, Martinho-Correa A, de Melo CLP, Durante-Yock LG e Carneiro T (2011). Desempenho agrônomo de genótipos de feijão comum cultivados no período “da seca” em Aquidauana-MS. *Revista Agrarian*, 4 n11, 33–42
- Du J, Grant D, Tian Z, Nelson RT, Zhu L, Shoemaker RC e Ma J (2010a). SoyTEdb: a comprehensive database of transposable elements in the soybean genome. *BMC Genomics*, 11:113-120
- Du J, *et al.*, (2010b). Evolutionary conservation, diversity and specificity of LTR-retrotransposons in flowering plants: insights from genome-wide analysis and multi-specific comparison. *Plant J*. 63(4):584–598
- Edlefsen PT e Liu JS (2010). Transposon identification using profile HMMs. *BMC Genomics*. 11(Suppl 1):S10
- Evgen'ev MB e Arkhipova IR (2005). Penelope-like elements--a new class of retroelements: distribution, function and possible evolutionary significance. *Cytogenet Genome Res*. 110(1-4):510–521
- Fernández-Suárez XM e Galperin MY (2013). The 2013 Nucleic Acids Research Database Issue and the online molecular biology database collection. *Nucl Acids Res*. 41(Database issue):1–7
- Feschotte C (2004). Merlin, a new superfamily of DNA transposons identified in diverse animal genomes and related to bacterial IS1016 insertion sequences. *Mol Biol Evol*. 21(9):1769–1780
- Freire SM, Wetzel MMVS, Falad MGR e Freire AB (1999) Germoplasma de caupi: coleção ativa e de base. 1ª edição. Recursos genéticos e melhoramento de plantas para o

- nordeste brasileiro. Petrolina, PE: Embrapa Semi-Árido; Brasília, DF: Embrapa Recursos Genéticos e Biotecnologia
- Freire-Filho FR, Ribeiro VQ, Rocha MM, SilvaK JD, Nogueira MSR, Rodrigues EV (2011). Feijão-caupi no Brasil: Produção, melhoramento genético, avanços e desafios. 1ª edição. Dados Internacionais de Catalogação na Publicação (CIP). Teresina, PI: Embrapa Meio-Norte
- Gemayel R, Cho J, Boeynaems S e Verstrepen KJ (2012). Beyond Junk-Variable Tandem Repeats as Facilitators of Rapid Evolution of Regulatory and Coding Sequences. *Genes*. 3(4): 461–480
- Gifford WD, Pfaff SL e Macfarlan TS (2013). Transposable elements as genetic regulatory substrates in early development. *Trends Cell Biol*. 23(5):218–26
- Gil E, Bosch A, Lampe D, Lizcano JM, Perales JC, Danos O e Chillon M (2013). Functional characterization of the human mariner transposon Hsmar2. *PLoS One*. 8(9):e73227-7339
- Gomes EP, Biscaro GA, Avila MR, Loosli FS, Vieira CV e Barbosa AP (2012). Desempenho agrônomo do feijoeiro comum de terceira safra sob irrigação na região Noroeste do Paraná. *Semina Ci Agr*. 33(3):899–910
- Goodwin TJD e Poulter RTM (2004). A new group of tyrosine recombinase-encoding retrotransposons. *Mol Biol Evol*. 21(4):746–59
- Grzebelus D, Lasota S, Gambin T, Kucherov G e Gambin A (2007). Diversity and structure of PIF/Harbinger-like elements in the genome of *Medicago truncatula*. *BMC Genomics*. 8:409-423
- Guermonprez H, Hénaff E, Cifuentes M e Casacuberta JM (2012) MITEs, Miniature elements with a major role in plant genome evolution. *Plant Transposable Elements*. 24:113-124
- Hanada K, *et al.*, (2009). The functional role of *pack-MULEs* in rice inferred from purifying selection and expression profile. *The Plant Cell*. 21(1):25–38
- Ištvánek J, Jaroš M, Křenek A e Řepková J (2014). Genome assembly and annotation for red clover (*Trifolium pratense*; Fabaceae). *Am J Bot*. 101(2)-327-337
- Jiang N, Bao Z, Zhang X, Eddy SR e Wessler SR (2004). Pack-MULE transposable elements mediate gene evolution in plants. *Nature*. 431(September):569–573
- Joly-Lopez Z e Bureau TE (2014). Diversity and evolution of transposable elements in *Arabidopsis*. *Chromosome Res*. 22(2):203–16
- Kaufman PD e Riot DC (1992). P Element Transposition In Vitro Proceeds by a Cut-and-Paste Mechanism and Uses GTP as a Cofactor. *Cell Press*. 69(1):27–39
- Kojima KK e Jurka J (2011). Crypton transposons: identification of new diverse families and ancient domestication events. *Mob DNA*. 2(1):12-29

- Krupovic M, Bamford DH e Koonin EV (2014). Conservation of major and minor jelly-roll capsid proteins in Polinton (Maverick) transposons suggests that they are bona fide viruses. *Biol Direct.* 9(1):6- 13
- Kurdyukov S, Song Y, Sheahan MB e Rose RJ (2014). Transcriptional regulation of early embryo development in the model legume *Medicago truncatula*. *Plant Cell Rep.* 33(2):349–362
- Krzywinski ML *et al.*, (2009) Circos: An information aesthetic for comparative genomics. *Genome Res.* 19(9):1639-1645
- Lee S-I, Park K-C, Ha M-W, Kim K-S, Jang Y-S, e Kim N-S (2012). CACTA transposon-derived Ti-SCARs for cultivar fingerprinting in rapeseed. *Genes Genom.* 34(5):575–579
- Lerat E (2010). Identifying repeats and transposable elements in sequenced genomes : how to find your way through the dense forest of programs. *Heredity.* 104(6):520–533
- Li Q, Li L, Dai J, Li J e Yan J (2009). Identification and characterization of CACTA transposable elements capturing gene fragments in maize. *Chinese Sci Bull.* 54(4):642–651
- Li X, Kahveci T e Settles AM (2008). Sequence analysis A novel genome-scale repeat finder geared towards transposons. *Bioinformatics.* 24(4):468–476
- Liu S, *et al.*, (2009). Mu transposon insertion sites and meiotic recombination events co-localize with epigenetic marks for open chromatin across the maize genome. *PLoS Gen.* 5(11): e1000733-1000746
- Mao *et al.*, (2015) Rice Transposable element: A survey of 73,000 Sequence-Tagged-Connectors. CSHL Press. 10:982-990
- Masson P, Strem M e Fedoroff N (1991). The tnpA and tnpD gene products of the Spm element are required for transposition in tobacco. *The Plant Cell.* 3(1):73–85
- Makarevitch I, Waters AJ e West PT (2014) Transposable element contribute to activation of maize genes in response to abiotic stress. *Plos Genet.* 11(1):e1004915
- Mayer K *et al.*, (1999) Sequence and analysis of chromosome 4 of the plant *Arabidopsis thaliana*. *Nature.* 402:769-777
- McClintock B (1951). Chromosome organization and genic expression. *Cold Spring Harb Sym.* 16:13–47
- McClintock B (1950). The origin and behaviour of mutable loci in maize. *Proc N A S.* 36(6):344–355
- Mingotte FLC, Guarnier CC, Farinelli RDO e Lemos LB (2013). Desempenho produtivo e qualidade pós-colheita de genótipos de feijão do grupo comercial carioca cultivados na época de inverno-primavera. *Biosci J.* 29(2010):1101–1110

- Nakamura Y, Cochrane G e Karsch-Mizrachi I (2013). The International Nucleotide Sequence Database Collaboration. *Nucleic Acids Res.* 41(Database issue):21–24
- Nobelprizeorg NMA - 2014 W 18 nov (2014). “Barbara McClintock - Nobel Lecture: The Significance of Responses of the Genome to Challenge” Web 18 Nov 2014 Retrieved from http://www.nobelprize.org/nobel_prizes/medicine/laureates/1983/mcclintock-lecture.html
- Oliver KR, McComb JA e Greene WK (2013). Transposable Elements: Powerful contributors to angiosperm evolution and diversity. *GBE.* 5(10):1886-1901
- Penton EH, Sullender BW e Crease TJ (2002). Pokey, a new DNA transposon in *Daphnia* (cladocera: crustacea). *J Mol Evol.* 55(6):664–73
- Pimentel C, *et al.*, (2002). Tolerância protoplasmática foliar à seca, em dois genótipos de Caupi cultivados em campo. *RCV.* 22(1):7–14
- Ping-fang T (2006). Progress in Plant CACTA Elements. *Acta Genetica Sin.* 33(30125030):765–774
- Quesneville H, Bergman CM, Andrieu O, Autard D, Nouaud D, Ashburner M e Anxolabehere D (2005). Combined evidence annotation of transposable elements in genome sequences. *PLoS Comp Biol.* 1(2):166–175
- Rodrigues AC, Emanuel J e Antunes L (2012). Resposta da co-inoculação de bactérias promotoras de crescimento em plantas e *Bradyrhizobium Sp.* IN COWPEA. 28(1):196–202
- Saze H, Tsugane K, Kanno T e Nishimura T (2012). DNA methylation in plants: relationship to small RNAs and histone modifications, and functions in transposon inactivation. *Plant Cell Physiol.* 53(5). 766–84
- Schmidt T (1999). LINEs , SINEs and repetitive DNA : non-LTR retrotransposons in plant genomes. *Plant Mol Biol,* 40(6).903–910
- Schumutz *et al.*, (2010). Genome sequence of the palaeopolyploid soybean. *Nature.* 463(7278):178-183
- Schumutz *et al.*, (2014). A reference genome for common bean and genome-wide analysis of dual domestications. *Nat Genet.* 46(7):707-713
- Singh NK *et al.*, (2012). The first draft of the pigeonpea genome sequence. *J Plant Biochem Biotechnol.* 21(1):98-112
- Sinzelle L, Kapitonov VV, Grzela DP, Jursch T, Jurka J e Izsva Z (2008). Transposition of a reconstructed Harbinger element in human cells and functional homology with two transposon-derived cellular genes. *PNAS,* 105(12):4715–4720

- Skibbe DS, Fernandes JF e Walbot, V (2012). Mu killer-Mediated and Spontaneous Silencing of *Zea mays Mutator* Family Transposable Elements Define Distinctive Paths of Epigenetic Inactivation. *Front Plant Sci.* 3(September):212
- Takahashi R, Morit, Y, Nakayama M, Kanazawa A e Abe J (2012). An Active CACTA-Family Transposable Element is Responsible for Flower Variegation in Wild Soybean. *TPG*, 5(2):62-70
- Tempel S (2012). Using and Understanding RepeatMasker. *MGE.* 859:29-51.
- The Arabidopsis Initiative G (2000). Analysis of the genome sequence of the flowering plant *Arabidopsis thaliana*. *Nature.* 408(6814):796–815
- Tóth G, Deák G, Barta E e Kiss GB (2006). PLOTREP: a web tool for defragmentation and visual analysis of dispersed genomic repeats. *Nucl Acids Res.* 34(Web Server issue):W708–13
- Varshney *et al.*, (2013) Draft genome sequence of chickpea (*Cicer arietinum*) provides a resource for trait improvement. *Nat Biot.* 31(3):240-246
- Wang W, Basia V, Arie A (2003). Plant responses to drought, salinity and extreme temperatures: towards genetic engineering for stress tolerance. *Planta.* 218(1):1–14
- Wang X, Weigel D e Smith Lm (2013). Transposon variants and their effects on gene expression in *Arabidopsis*. *Plos Genet.* 9(2):e1003255
- Wicker T, Guyot R, Yahiaoui N e Keller B (2003). CACTA Transposons in Triticeae A Diverse Family of High-Copy Repetitive Elements. *Plant Physiol.* 132(May):52–63
- Wicker T, et al., (2007). A unified classification system for eukaryotic transposable elements. *Nat Rev Genet.* 8(12):973–82
- Young ND *et al.*, (2011). The *Medicago* genome provides insight into the evolution of rhizobial symbioses. *Nature.* 498(7378):520-524
- Zeh DW, Zeh JA e Ishida Y (2009). Transposable elements and an epigenetic basis for punctuated equilibria. *Bioessays.* 31(7):715–26
- Zhang X, Jiang N, Feschotte C e Wessler SR (2004). *PIF*- and *Pong*-like transposable element: distribution, evolution and relationship with miniature-inverted-repeat transposable element. *GSA.* 166:971-986

Anexos

Anexo 1: Lista das sondas utilizadas neste trabalho.

Sondas de *Mutator* (Du et al., 2010a)

DTM_uuu_Gm10-1	DTM_uuu_Gm10-34	DTM_uuu_Gm10-9	DTM_uuu_Gm11-28	DTM_uuu_Gm11-79
DTM_uuu_Gm10-10	DTM_uuu_Gm10-35	DTM_uuu_Gm10-90	DTM_uuu_Gm11-29	DTM_uuu_Gm1-18
DTM_uuu_Gm10-100	DTM_uuu_Gm10-36	DTM_uuu_Gm10-91	DTM_uuu_Gm11-13	DTM_uuu_Gm11-8
DTM_uuu_Gm10-101	DTM_uuu_Gm10-37	DTM_uuu_Gm10-92	DTM_uuu_Gm11-3	DTM_uuu_Gm11-80
DTM_uuu_Gm10-102	DTM_uuu_Gm10-38	DTM_uuu_Gm10-93	DTM_uuu_Gm11-30	DTM_uuu_Gm11-81
DTM_uuu_Gm10-103	DTM_uuu_Gm10-39	DTM_uuu_Gm10-94	DTM_uuu_Gm11-31	DTM_uuu_Gm11-82
DTM_uuu_Gm10-104	DTM_uuu_Gm10-4	DTM_uuu_Gm10-95	DTM_uuu_Gm11-32	DTM_uuu_Gm1-19
DTM_uuu_Gm10-105	DTM_uuu_Gm10-40	DTM_uuu_Gm10-96	DTM_uuu_Gm11-33	DTM_uuu_Gm11-9
DTM_uuu_Gm10-106	DTM_uuu_Gm10-41	DTM_uuu_Gm10-97	DTM_uuu_Gm11-34	DTM_uuu_Gm1-2
DTM_uuu_Gm10-107	DTM_uuu_Gm10-42	DTM_uuu_Gm10-98	DTM_uuu_Gm11-35	DTM_uuu_Gm1-20
DTM_uuu_Gm10-108	DTM_uuu_Gm10-43	DTM_uuu_Gm10-99	DTM_uuu_Gm11-36	DTM_uuu_Gm1-21
DTM_uuu_Gm10-109	DTM_uuu_Gm10-44	DTM_uuu_Gm1-1	DTM_uuu_Gm11-37	DTM_uuu_Gm12-1
DTM_uuu_Gm10-11	DTM_uuu_Gm10-45	DTM_uuu_Gm1-10	DTM_uuu_Gm11-38	DTM_uuu_Gm12-10
DTM_uuu_Gm10-110	DTM_uuu_Gm10-46	DTM_uuu_Gm1-100	DTM_uuu_Gm11-39	DTM_uuu_Gm12-100
DTM_uuu_Gm10-111	DTM_uuu_Gm10-47	DTM_uuu_Gm1-101	DTM_uuu_Gm1-14	DTM_uuu_Gm12-101
DTM_uuu_Gm10-112	DTM_uuu_Gm10-48	DTM_uuu_Gm1-102	DTM_uuu_Gm11-4	DTM_uuu_Gm12-102
DTM_uuu_Gm10-113	DTM_uuu_Gm10-49	DTM_uuu_Gm1-103	DTM_uuu_Gm11-40	DTM_uuu_Gm12-103
DTM_uuu_Gm10-114	DTM_uuu_Gm10-5	DTM_uuu_Gm1-104	DTM_uuu_Gm11-41	DTM_uuu_Gm12-104
DTM_uuu_Gm10-115	DTM_uuu_Gm10-50	DTM_uuu_Gm1-105	DTM_uuu_Gm11-42	DTM_uuu_Gm12-105
DTM_uuu_Gm10-116	DTM_uuu_Gm10-51	DTM_uuu_Gm1-106	DTM_uuu_Gm11-43	DTM_uuu_Gm12-11
DTM_uuu_Gm10-117	DTM_uuu_Gm10-52	DTM_uuu_Gm1-107	DTM_uuu_Gm11-44	DTM_uuu_Gm12-12
DTM_uuu_Gm10-118	DTM_uuu_Gm10-53	DTM_uuu_Gm1-108	DTM_uuu_Gm11-45	DTM_uuu_Gm12-13
DTM_uuu_Gm10-119	DTM_uuu_Gm10-54	DTM_uuu_Gm1-109	DTM_uuu_Gm11-46	DTM_uuu_Gm12-14
DTM_uuu_Gm10-12	DTM_uuu_Gm10-55	DTM_uuu_Gm1-11	DTM_uuu_Gm11-47	DTM_uuu_Gm12-15
DTM_uuu_Gm10-120	DTM_uuu_Gm10-56	DTM_uuu_Gm11-1	DTM_uuu_Gm11-48	DTM_uuu_Gm12-16
DTM_uuu_Gm10-121	DTM_uuu_Gm10-57	DTM_uuu_Gm1-110	DTM_uuu_Gm11-49	DTM_uuu_Gm12-17
DTM_uuu_Gm10-122	DTM_uuu_Gm10-58	DTM_uuu_Gm11-10	DTM_uuu_Gm1-15	DTM_uuu_Gm12-18
DTM_uuu_Gm10-123	DTM_uuu_Gm10-59	DTM_uuu_Gm1-111	DTM_uuu_Gm11-5	DTM_uuu_Gm12-19
DTM_uuu_Gm10-124	DTM_uuu_Gm10-6	DTM_uuu_Gm11-11	DTM_uuu_Gm11-50	DTM_uuu_Gm1-22
DTM_uuu_Gm10-125	DTM_uuu_Gm10-60	DTM_uuu_Gm1-112	DTM_uuu_Gm11-51	DTM_uuu_Gm12-2
DTM_uuu_Gm10-126	DTM_uuu_Gm10-61	DTM_uuu_Gm11-12	DTM_uuu_Gm11-52	DTM_uuu_Gm12-20
DTM_uuu_Gm10-127	DTM_uuu_Gm10-62	DTM_uuu_Gm1-113	DTM_uuu_Gm11-53	DTM_uuu_Gm12-21
DTM_uuu_Gm10-128	DTM_uuu_Gm10-63	DTM_uuu_Gm11-13	DTM_uuu_Gm11-54	DTM_uuu_Gm12-22
DTM_uuu_Gm10-129	DTM_uuu_Gm10-64	DTM_uuu_Gm1-114	DTM_uuu_Gm11-55	DTM_uuu_Gm12-23
DTM_uuu_Gm10-13	DTM_uuu_Gm10-65	DTM_uuu_Gm11-14	DTM_uuu_Gm11-56	DTM_uuu_Gm12-24
DTM_uuu_Gm10-130	DTM_uuu_Gm10-66	DTM_uuu_Gm1-115	DTM_uuu_Gm11-57	DTM_uuu_Gm12-25
DTM_uuu_Gm10-131	DTM_uuu_Gm10-67	DTM_uuu_Gm11-15	DTM_uuu_Gm11-58	DTM_uuu_Gm12-26
DTM_uuu_Gm10-132	DTM_uuu_Gm10-68	DTM_uuu_Gm1-116	DTM_uuu_Gm11-59	DTM_uuu_Gm12-27
DTM_uuu_Gm10-133	DTM_uuu_Gm10-69	DTM_uuu_Gm11-16	DTM_uuu_Gm1-16	DTM_uuu_Gm12-28
DTM_uuu_Gm10-14	DTM_uuu_Gm10-7	DTM_uuu_Gm1-117	DTM_uuu_Gm11-6	DTM_uuu_Gm12-29
DTM_uuu_Gm10-15	DTM_uuu_Gm10-70	DTM_uuu_Gm11-17	DTM_uuu_Gm11-60	DTM_uuu_Gm1-23
DTM_uuu_Gm10-16	DTM_uuu_Gm10-71	DTM_uuu_Gm1-118	DTM_uuu_Gm11-61	DTM_uuu_Gm12-3
DTM_uuu_Gm10-17	DTM_uuu_Gm10-72	DTM_uuu_Gm11-18	DTM_uuu_Gm11-62	DTM_uuu_Gm12-30
DTM_uuu_Gm10-18	DTM_uuu_Gm10-73	DTM_uuu_Gm1-119	DTM_uuu_Gm11-63	DTM_uuu_Gm12-31
DTM_uuu_Gm10-19	DTM_uuu_Gm10-74	DTM_uuu_Gm11-19	DTM_uuu_Gm11-64	DTM_uuu_Gm12-32
DTM_uuu_Gm10-2	DTM_uuu_Gm10-75	DTM_uuu_Gm1-12	DTM_uuu_Gm11-65	DTM_uuu_Gm12-33
DTM_uuu_Gm10-20	DTM_uuu_Gm10-76	DTM_uuu_Gm11-2	DTM_uuu_Gm11-66	DTM_uuu_Gm12-34
DTM_uuu_Gm10-21	DTM_uuu_Gm10-77	DTM_uuu_Gm1-120	DTM_uuu_Gm11-67	DTM_uuu_Gm12-35
DTM_uuu_Gm10-22	DTM_uuu_Gm10-78	DTM_uuu_Gm11-20	DTM_uuu_Gm11-68	DTM_uuu_Gm12-36
DTM_uuu_Gm10-23	DTM_uuu_Gm10-79	DTM_uuu_Gm1-121	DTM_uuu_Gm11-69	DTM_uuu_Gm12-37
DTM_uuu_Gm10-24	DTM_uuu_Gm10-8	DTM_uuu_Gm11-21	DTM_uuu_Gm1-17	DTM_uuu_Gm12-38
DTM_uuu_Gm10-25	DTM_uuu_Gm10-80	DTM_uuu_Gm1-122	DTM_uuu_Gm11-7	DTM_uuu_Gm12-39
DTM_uuu_Gm10-26	DTM_uuu_Gm10-81	DTM_uuu_Gm11-22	DTM_uuu_Gm11-70	DTM_uuu_Gm1-24
DTM_uuu_Gm10-27	DTM_uuu_Gm10-82	DTM_uuu_Gm1-123	DTM_uuu_Gm11-71	DTM_uuu_Gm12-4
DTM_uuu_Gm10-28	DTM_uuu_Gm10-83	DTM_uuu_Gm11-23	DTM_uuu_Gm11-72	DTM_uuu_Gm12-40
DTM_uuu_Gm10-29	DTM_uuu_Gm10-84	DTM_uuu_Gm1-124	DTM_uuu_Gm11-73	DTM_uuu_Gm12-41
DTM_uuu_Gm10-3	DTM_uuu_Gm10-85	DTM_uuu_Gm11-24	DTM_uuu_Gm11-74	DTM_uuu_Gm12-42
DTM_uuu_Gm10-30	DTM_uuu_Gm10-86	DTM_uuu_Gm1-125	DTM_uuu_Gm11-75	DTM_uuu_Gm12-43
DTM_uuu_Gm10-31	DTM_uuu_Gm10-87	DTM_uuu_Gm11-25	DTM_uuu_Gm11-76	DTM_uuu_Gm12-44
DTM_uuu_Gm10-32	DTM_uuu_Gm10-88	DTM_uuu_Gm11-26	DTM_uuu_Gm11-77	DTM_uuu_Gm12-45
DTM_uuu_Gm10-33	DTM_uuu_Gm10-89	DTM_uuu_Gm11-27	DTM_uuu_Gm11-78	DTM_uuu_Gm12-46

DTM_uuu_Gm8-30	DTM_uuu_Gm8-71	DTM_uuu_Gm9-111	DTM_uuu_Gm9-35	DTM_uuu_Gm9-76
DTM_uuu_Gm8-31	DTM_uuu_Gm8-72	DTM_uuu_Gm9-112	DTM_uuu_Gm9-36	DTM_uuu_Gm9-77
DTM_uuu_Gm8-32	DTM_uuu_Gm8-73	DTM_uuu_Gm9-113	DTM_uuu_Gm9-37	DTM_uuu_Gm9-78
DTM_uuu_Gm8-33	DTM_uuu_Gm8-74	DTM_uuu_Gm9-114	DTM_uuu_Gm9-38	DTM_uuu_Gm9-79
DTM_uuu_Gm8-34	DTM_uuu_Gm8-75	DTM_uuu_Gm9-115	DTM_uuu_Gm9-39	DTM_uuu_Gm9-8
DTM_uuu_Gm8-35	DTM_uuu_Gm8-76	DTM_uuu_Gm9-116	DTM_uuu_Gm9-4	DTM_uuu_Gm9-80
DTM_uuu_Gm8-36	DTM_uuu_Gm8-77	DTM_uuu_Gm9-117	DTM_uuu_Gm9-40	DTM_uuu_Gm9-81
DTM_uuu_Gm8-37	DTM_uuu_Gm8-78	DTM_uuu_Gm9-118	DTM_uuu_Gm9-41	DTM_uuu_Gm9-82
DTM_uuu_Gm8-38	DTM_uuu_Gm8-79	DTM_uuu_Gm9-119	DTM_uuu_Gm9-42	DTM_uuu_Gm9-83
DTM_uuu_Gm8-39	DTM_uuu_Gm8-8	DTM_uuu_Gm9-12	DTM_uuu_Gm9-43	DTM_uuu_Gm9-84
DTM_uuu_Gm8-4	DTM_uuu_Gm8-80	DTM_uuu_Gm9-120	DTM_uuu_Gm9-44	DTM_uuu_Gm9-85
DTM_uuu_Gm8-40	DTM_uuu_Gm8-81	DTM_uuu_Gm9-121	DTM_uuu_Gm9-45	DTM_uuu_Gm9-86
DTM_uuu_Gm8-41	DTM_uuu_Gm8-82	DTM_uuu_Gm9-122	DTM_uuu_Gm9-46	DTM_uuu_Gm9-87
DTM_uuu_Gm8-42	DTM_uuu_Gm8-83	DTM_uuu_Gm9-123	DTM_uuu_Gm9-47	DTM_uuu_Gm9-88
DTM_uuu_Gm8-43	DTM_uuu_Gm8-84	DTM_uuu_Gm9-124	DTM_uuu_Gm9-48	DTM_uuu_Gm9-89
DTM_uuu_Gm8-44	DTM_uuu_Gm8-85	DTM_uuu_Gm9-125	DTM_uuu_Gm9-49	DTM_uuu_Gm9-9
DTM_uuu_Gm8-45	DTM_uuu_Gm8-86	DTM_uuu_Gm9-126	DTM_uuu_Gm9-5	DTM_uuu_Gm9-90
DTM_uuu_Gm8-46	DTM_uuu_Gm8-87	DTM_uuu_Gm9-127	DTM_uuu_Gm9-50	DTM_uuu_Gm9-91
DTM_uuu_Gm8-47	DTM_uuu_Gm8-88	DTM_uuu_Gm9-128	DTM_uuu_Gm9-51	DTM_uuu_Gm9-92
DTM_uuu_Gm8-48	DTM_uuu_Gm8-89	DTM_uuu_Gm9-129	DTM_uuu_Gm9-52	DTM_uuu_Gm9-93
DTM_uuu_Gm8-49	DTM_uuu_Gm8-9	DTM_uuu_Gm9-13	DTM_uuu_Gm9-53	DTM_uuu_Gm9-94
DTM_uuu_Gm8-5	DTM_uuu_Gm8-90	DTM_uuu_Gm9-130	DTM_uuu_Gm9-54	DTM_uuu_Gm9-95
DTM_uuu_Gm8-50	DTM_uuu_Gm8-91	DTM_uuu_Gm9-14	DTM_uuu_Gm9-55	DTM_uuu_Gm9-96
DTM_uuu_Gm8-51	DTM_uuu_Gm8-92	DTM_uuu_Gm9-15	DTM_uuu_Gm9-56	DTM_uuu_Gm9-97
DTM_uuu_Gm8-52	DTM_uuu_Gm8-93	DTM_uuu_Gm9-16	DTM_uuu_Gm9-57	DTM_uuu_Gm9-98
DTM_uuu_Gm8-53	DTM_uuu_Gm8-94	DTM_uuu_Gm9-17	DTM_uuu_Gm9-58	DTM_uuu_Gm9-99
DTM_uuu_Gm8-54	DTM_uuu_Gm8-95	DTM_uuu_Gm9-18	DTM_uuu_Gm9-59	DTM_uuu_Scaffold22-1
DTM_uuu_Gm8-55	DTM_uuu_Gm8-96	DTM_uuu_Gm9-19	DTM_uuu_Gm9-6	DTM_uuu_Scaffold22-2
DTM_uuu_Gm8-56	DTM_uuu_Gm8-97	DTM_uuu_Gm9-2	DTM_uuu_Gm9-60	DTM_uuu_Scaffold23-1
DTM_uuu_Gm8-57	DTM_uuu_Gm8-98	DTM_uuu_Gm9-20	DTM_uuu_Gm9-61	DTM_uuu_Scaffold420-1
DTM_uuu_Gm8-58	DTM_uuu_Gm8-99	DTM_uuu_Gm9-21	DTM_uuu_Gm9-62	DTM_uuu_Scaffold80-1
DTM_uuu_Gm8-59	DTM_uuu_Gm9-1	DTM_uuu_Gm9-22	DTM_uuu_Gm9-63	DTM_uuu_Scaffold84-1
DTM_uuu_Gm8-6	DTM_uuu_Gm9-10	DTM_uuu_Gm9-23	DTM_uuu_Gm9-64	DTM_uuu_Scaffold84-2
DTM_uuu_Gm8-60	DTM_uuu_Gm9-100	DTM_uuu_Gm9-24	DTM_uuu_Gm9-65	DTM_uuu_Scaffold84-3
DTM_uuu_Gm8-61	DTM_uuu_Gm9-101	DTM_uuu_Gm9-25	DTM_uuu_Gm9-66	DTM_uuu_Scaffold90-1
DTM_uuu_Gm8-62	DTM_uuu_Gm9-102	DTM_uuu_Gm9-26	DTM_uuu_Gm9-67	DTM_uuu_Scaffold90-2
DTM_uuu_Gm8-63	DTM_uuu_Gm9-103	DTM_uuu_Gm9-27	DTM_uuu_Gm9-68	DTM_uuu_Scaffold92-1
DTM_uuu_Gm8-64	DTM_uuu_Gm9-104	DTM_uuu_Gm9-28	DTM_uuu_Gm9-69	DTM_uuu_Scaffold92-2
DTM_uuu_Gm8-65	DTM_uuu_Gm9-105	DTM_uuu_Gm9-29	DTM_uuu_Gm9-7	
DTM_uuu_Gm8-66	DTM_uuu_Gm9-106	DTM_uuu_Gm9-3	DTM_uuu_Gm9-70	
DTM_uuu_Gm8-67	DTM_uuu_Gm9-107	DTM_uuu_Gm9-30	DTM_uuu_Gm9-71	
DTM_uuu_Gm8-68	DTM_uuu_Gm9-108	DTM_uuu_Gm9-31	DTM_uuu_Gm9-72	
DTM_uuu_Gm8-69	DTM_uuu_Gm9-109	DTM_uuu_Gm9-32	DTM_uuu_Gm9-73	
DTM_uuu_Gm8-7	DTM_uuu_Gm9-11	DTM_uuu_Gm9-33	DTM_uuu_Gm9-74	
DTM_uuu_Gm8-70	DTM_uuu_Gm9-110	DTM_uuu_Gm9-34	DTM_uuu_Gm9-75	

Sondas de CACTA (Du *et al.*, 2010a)

DTC_uuu_Gm10-1	DTC_uuu_Gm12-1	DTC_uuu_Gm1-3	DTC_uuu_Gm15-2
DTC_uuu_Gm10-2	DTC_uuu_Gm12-2	DTC_uuu_Gm1-4	DTC_uuu_Gm16-1
DTC_uuu_Gm10-3	DTC_uuu_Gm12-3	DTC_uuu_Gm14-1	DTC_uuu_Gm16-2
DTC_uuu_Gm10-4	DTC_uuu_Gm12-4	DTC_uuu_Gm14-2	DTC_uuu_Gm17-1
DTC_uuu_Gm1-1	DTC_uuu_Gm12-5	DTC_uuu_Gm14-3	DTC_uuu_Gm17-2
DTC_uuu_Gm11-1	DTC_uuu_Gm12-6	DTC_uuu_Gm1-5	DTC_uuu_Gm17-3
DTC_uuu_Gm1-2	DTC_uuu_Gm12-7	DTC_uuu_Gm15-1	DTC_uuu_Gm17-4
DTC_uuu_Gm17-5	DTC_uuu_Gm19-4	DTC_uuu_Gm3-4	DTC_uuu_Gm7-1
DTC_uuu_Gm18-1	DTC_uuu_Gm20-1	DTC_uuu_Gm4-1	DTC_uuu_Gm8-1
DTC_uuu_Gm18-2	DTC_uuu_Gm20-2	DTC_uuu_Gm4-2	DTC_uuu_Gm9-1
DTC_uuu_Gm18-3	DTC_uuu_Gm20-3	DTC_uuu_Gm4-3	DTC_uuu_Gm9-2
DTC_uuu_Gm18-4	DTC_uuu_Gm20-4	DTC_uuu_Gm5-1	DTC_uuu_Gm9-3
DTC_uuu_Gm18-5	DTC_uuu_Gm2-1	DTC_uuu_Gm5-2	DTC_uuu_Gm9-4
DTC_uuu_Gm18-6	DTC_uuu_Gm2-2	DTC_uuu_Gm5-3	DTC_uuu_Gm9
DTC_uuu_Gm19-1	DTC_uuu_Gm3-1	DTC_uuu_Gm6-1	
DTC_uuu_Gm19-2	DTC_uuu_Gm3-2	DTC_uuu_Gm6-2	
DTC_uuu_Gm19-3	DTC_uuu_Gm3-3	DTC_uuu_Gm6-3	

Anexo 2: Porcentagem (%) e quantidades em Megabases de sequências para cada superfamília calculada para o tamanho real dos genomas de *Medicago truncatula*, *Phaseolus vulgaris* e *Vigna unguiculata* (ver gráfico complementar na Figura 9 no corpo do texto). Valores ao lado das espécies representam o valor do genoma real de acordo com (Arumuganathan e Earle, 1991). Valores em negrito destacam as superfamílias com maior porcentagem nos genomas.

Classe	Superfamília	<i>M. truncatula</i> (490 Mb)		<i>P. vulgaris</i> (637 Mb)		<i>V. unguiculata</i> (613 Mb)	
		Mb	%	Mb	%*	Mb	%*
II	<i>Helitron</i>	3,596710845	0,734022621	1,032771989	0,210769794	1,011282735	0,206384232
	<i>Tc-mariner</i>	1,291409182	0,263552894	0,037804657	0,007715236	0,098911211	0,020185961
	<i>PIF</i>	2,443688357	0,49871191	0,373446507	0,076213573	0,103051597	0,021030938
	<i>MULEs</i>	12,15318506	2,48024185	1,865058017	0,380624085	0,813833599	0,16608849
	<i>hAT</i>	2,260567565	0,461340319	0,972092515	0,198386228	0,535078044	0,109199601
	<i>CACTA</i>	1,148880472	0,234465403	4,920393846	1,004162009	1,332289787	0,271895875
	Total	22,89444148	4,672334997	9,201567532	1,877870925	3,894446973	0,794785096
I	<i>LINE</i>	17,85101949	3,643065203	5,597242615	1,142294411	0,216377545	0,044158683
	<i>SINE</i>	0,645614937	0,13175815	0,017236394	0,003517631	0,013382901	0,002731204
	<i>LTR</i>	64,9676851	13,25871124	73,32525126	14,96433699	15,00111038	3,061451098
	Total	83,46431953	22,40167532	78,93973027	19,1905682	15,23087083	4,284221557

Anexo 4: Lista com todos os domínios de genes relacionados a sequências de *Mutator* identificados em *Vigna unguiculata*. Nomenclatura derivada do CDSearch Bach (NCBI).

Pingo de Ouro

PTKc_Tyro3	KISc_KLP2_like	PTZ00266	FH2 superfamily	STKc_LATS1	pknD
PTKc_Yes	Kri1	STKc_aPKC_iota	HEAT_EZ	STKc_LATS2	AMP-binding
PTKc_Zap-70	PHA02929	STKc_aPKC_zeta	superfamily	STKc_MRCK_alpha	CaiC
STKc_cGK_PKG	PRK01622	STKc_beta_ARK	MRP_assoc_pro	STKc_MRCK_beta	DUF4106
STKc_MAST	PTKc_EGFR	STKc_CDC2L6	ndhl	STKc_MSK_N	GAT_Gln-NAD-
STKc_MEKK3	PTKc_VEGFR	STKc_CDK6	Peptidase_C19D	STKc_MSK2_N	synth
STKc_MSK1_N	STKc_aPKC	STKc_CDK8	Peptidase_C19E	STKc_nPKC_delta	nitrilase
STKc_myosinIIIA	STKc_Nek6	STKc_CR1K	Peptidase_C19F	STKc_nPKC_epsilon	superfamily
STKc_myosinIIIB	STKc_nPKC_eta	STKc_DMPK_like	Peptidase_C19K	n	OGG_N
STKc_NDR1	STKc_nPKC_theta	STKc_ERK5	Peptidase_C19O	STKc_p70S6K	PLN02339
STKc_Nek11	delta	STKc_MAP4K4_6	Peptidase_C19R	STKc_PKB_alpha	PLN02479
STKc_Nek4	STKc_PKC	STKc_MASTL	PLN02647	STKc_PKB_beta	PLN03102
STKc_PFTAIRE2	zf-Di19	STKc_MEKK2	PLN03023	STKc_PKB_gamma	PRK08162
STKc_phototropin_	ATP-synt_C	STKc_MPK1	PLN03130	STKc_TDY_MAPK	ttLC_FACS_AEE21
like	DnaJ	STKc_Nek7	PLN03232	plant	_like
STKc_PKN	Peptidases_S8_S53	STKc_Nek9	PLN03232	TOMM_kin_cyc	UBN2
STKc_ROCK1	superfamily	STKc_NLK	PRK05192	xanthine_xdhB	eIF-4B
STKc_ROCK2	Peptidases_S8_Tri	STKc_nPKC_theta	PRK08310	XdhB	Fasciclin
STKc_Sck1_like	peptidyl_Amino pep	STKc_p38alpha_M	PRK09603	Alg14	GAT_1 superfamily
STKc_SGK	tidase_II	APK14	rpoC1	chap_CCT_alpha	PLN00140
STKc_SGK1	PRK08376	STKc_p38beta_MA	SSM4	chap_CCT_delta	PLN02481
STKc_SGK2	PTKc_Aatyk2	PK11	TAXI_C	chap_CCT_epsilon	PLN02832
STKc_SGK3	PTKc_PDGFR_beta	STKc_p38delta_MA	UBP12	chap_CCT_eta	PP2
STKc_SLK	GMC_oxred_N	PK13	zf-C3HC4	chaperonin_type_I	RAB
ATP-synt_C	PHA02988	STKc_p38gamma	Ald_Xan_dh_C2	II	Ras_like GTPase
superfamily	Rv0697	MAPK12	AspS	cpn60	superfamily
DUF605	AFD_class_I	STKc_PKB	COG2 superfamily	Cpn60_TCP1	Rho
V_ATP_synt_C	superfamily	STKc_Rim15_like	Dor1	DamX	SixA
Fer4_NifH	asnC	STKc_Sid2p_Dbf2p	Ferritin_2	FBOX	23DHB-AMP_Ig
ParA	Fmp27	STKc_Sty1_Hog1	PHA03212	F-box	2A0305
PMEI	Fmp27 superfamily	STKc_Tey_MAPK_	PLN00192	F-box-like	2A0306
PRK00771	Glyco_transf_8	plant	PLN02659	Glycosyltransferas	2A0308
PRK10416	Isy1	STKc_TNIK	PLN02906	e_GTB_type	AA_permease_2
PRK10867	nadB	UBCC superfamily	PLN03185	superfamily	AAA_assoc
PRK12724	Pat17_isozyyme_like	UQ_con	PLN03224	GroL	ANTH_AP180_CAL
PRK12726	Pat17_PNPLA8_PN	zf-rbx1	PTKc_HER2	GT1_ALG1_like	M
PRK14974	PLA9_like	chap_CCT_gamma	PTKc_Syk	HELICc	APP
SRP	Pat17_PNPLA8_PN	chaperonin_like	PTKc_Tie1	KISc_KID_like	APP_MetAP
SRP54_euk	PLA9_like1	superfamily	PTKc_Tie2	KISc_KIF23_like	superfamily
DUF260	Patatin_and_cPLA2	Condensation	PTKc_VEGFR2	Motor_domain	BFIT_BACH
PRK02106	superfamily	superfamily	PTKc_VEGFR3	PBP1	COG5127
TFIIA_alpha_beta_li	Peptidase_C19	PLN02663	PTZ00031	PDZ	DUF679
ke	PLN02603	PLN03157	PTZ00361	PHA01929	Dynein_light
Di19_C	PRK04176	PRK06558	PTZ00426	PHA03307	Dynein_light
FH2	PTZ00425	PRK13042	PTZ00454	PLN02275	superfamily
GGCT_like	TIGR00292	PRK14951	Ribosomal_L2	PLN03188	ENTH
superfamily	VHS	TCP1_gamma	rpl2	Prc	Histone
Peptidase_C19	VHS_ENTH_ANTH	Transferase	rplB_bact	PRK04654	hot_dog
superfamily	superfamily	PLN02708	RPT1	PRK10263	superfamily
PLN02221	Ycf1 superfamily	PMEI superfamily	SpoVK	PRK10856	IlvD
PLN02532	APC11	rplW	STKc_cPKC	PRK11192	ILVD_EDD
PLN02661	COG5078	DUF566	STKc_cPKC_alpha	PRK11634	superfamily
PLN02722	Exo70	eutB	STKc_cPKC_beta	PRK11776	lytB_ispH
rpoB	GT8_like_1	gpmA	STKc_GRK	PRK14954	superfamily
sdhA	PHA03207	PLN02217	STKc_GRK1	PTZ00110	Maf_Ham1
UCH	PHA03209	Ras	STKc_GRK3	PTZ00212	superfamily
60KD_IMP	PHA03211	Trp-synth-beta_II	STKc_GRK4	PTZ00424	PepP
superfamily	PLN02523	superfamily	STKc_GRK4_like	SrmB	PLN02734
KIP1	PLN02718	WD40	STKc_GRK5	TCP1_alpha	PLN03058
KISc_BimC_Eg5	PLN02742	Amidase	STKc_GRK6	TCP1_beta	PotE
KISc_C_terminal	PLN02769	superfamily	STKc_GRK7	TCP1_delta	PRK00566
KISc_KIF1A_KIF1B	PLN02829	ANK superfamily	STKc_JNK	TCP1_epsilon	PRK00911
KISc_KIF2_like	PLN02867	Ank_2	STKc_JNK1	TCP1_eta	PRK07656
KISc_KIF3	PLN02870	ANTH	STKc_JNK2	thermosome_arch	PRK13371
KISc_KIF4	PLN02910	ARM superfamily	STKc_JNK3	MRS6	PRK14367
KISc_KIF9_like	PTKc_Tie	Chloroa_b-bind	STKc_LATS	NPL4	PRK14906

PTZ00059 RNA_pol_Rpb2_1 RNA_pol_Rpb2_2 RNAP_beta'_N RNAP_largest_subunit_N superfamily rpoC_TIGR rpoC1_cyan RSN1_TM RSN1_TM superfamily Telomere_reg-2 V-ATPase_C superfamily VHS_ENTH_ANTH	EcAsnRS_like_N gag_pre-integrals kgd Mis12 PHA02882 PHA03210 PLN00162 PLN02560 PRK03992 PTZ00036 pup_AAA RNase_HI_RT_Ty1 rve_3 THUMP superfamily THUMP_THUMPD1	PRK10767 PRK14276 PRK14277 PRK14278 PRK14280 PRK14281 PRK14282 PRK14283 PRK14284 PRK14285 PRK14286 PRK14287 PRK14288 PRK14289 PRK14290 PRK14291	DUF3134 Glyco_transf_90 MPN superfamily MPN_NPL4 PLN02201 PLN02423 PLN02468 PLN02745 PLN03225 PMM superfamily PRK09722 PTZ00174 Retrotrans_gag Ribosomal_L23 rpl23 thiQ TIM_phosphate_biding superfamily ubiquitin UBQ 4CL DUF788 DUF974 EntE PH-like superfamily PLN02330 PLN02405 Prefoldin superfamily Prefoldin_2 PRK06187 PRK07529 PRK08275 PRK08276 PRK08316 RanBD_NUP50 Tmemb_161AB	ttLC_FACS_AiKk_ike ke ttLC_FACS_like UBN2_2 V-SNARE_C zf-Apc11 Cas8c_I-C Cdc42 COG1100 DUF4413 EF2 FusA gag-asp_proteas GPH_sucrose His_Phos_1 HP HP_HAP_like HP_PGM_like Miro1 PGAM pgm_1 phoE PK_STRAD_beta pk1 PLN00023 PLN03108 PLN03110 PLN03118 PRK01295 PRK07238 PRK07560 PRK10848 PRK15004 PTZ00122 PTZ00416 Rab1_Ypt1 Rab11_like Rab14	Rab15 Rab18 Rab19 Rab2 Rab21 Rab26 Rab27A Rab3 Rab30 Rab32_Rab38 Rab33B_Rab33A Rab35 Rab39 Rab4 Rab40 Rab5_related Rab6 Rab7 Rab8_Rab10_Rab13_like Rab9 RabL2 RabL3 Rac1_like retropepsin_like RGK Rho4_like RhoG ribazole_cobC Rit_Rin_Ric RPN6 RT_DIRS1 small_GTP small_GTPase Spg1 TAF6 WD40 superfamily Wrch_1
VHS_GGA VHS_Hrs_Vps27p VHS_STAM VHS_Tom1 26Sp45 Asp_Lys_Asn_RS_core Asp_Lys_Asn_RS_	_like Tra5 tRNA-synt_2 AMN1 Anoctamin CbpA chap_CCT_beta DnaJ superfamily	PRK14292 PRK14293 PRK14294 PRK14295 PRK14297 PRK14298 PRK14299 PRK14300	PRK14301 PRK14709 PT_UbiA superfamily PTZ00037 termin_org_DnaJ UBP5 zf-A20 zf-AN1 ZnF_A20 ZnF_AN1 ZUO1 3a0901s04IAP86 Dimer_Tnp_hAT		
N superfamily aspC AspRS_core aspS_bact aspS_nondisc AsxRS_core C2 superfamily C2B_MCTP_PRT_plant CDC48 CFTR_protein class_II_aaRS-like_core superfamily	DnaJ_bact DUF705 Kelch_3 superfamily Kelch_6 PDZ superfamily PDZ_CTP_protease PDZ_signaling PHA03398 PLN00012 PLN00049 PRK04537 PRK10266 PRK10590				

Santo Inácio

PLN00113 PMD PKc_like superfamily Pkinase_Tyr PTKc STYKc TyrKc PKc PTKc_Trk PTKc_Csk_like PTKc_TrkB PTKc_Ror S_TKc STKc_MAPKKK Pkinase PTKc_Csk PTKc_DDR_like PTKc_EphR PTKc_Itk PTKc_Jak_rpt2 PTKc_Tec_like STKc_Cdc7_like PLN02785 PTKc_Musk PTKc_TrkC PTK_Jak2_Jak3_rpt1 PTKc_c-ros TCP1_alpha TCP1_delta TCP1_eta thermosome_arch	PTKc_DDR1 PTKc_TrkA PKc_STE PTKc_Abl PTKc_Ack_like PTKc_Btk_Bmx PTKc_DDR PTKc_DDR2 PTKc_EGFR_like e PTKc_EphR_A PTKc_EphR_A2 PTKc_Lyn PTKc_Ror1 PTKc_Src_like PTKc_Tec_Rik PTKc_Tyk2_rpt2 STKc_ASK PRK02106 STKc_MEKK1 STKc_MEKK1_p lant STKc_MEKK3_ike ke STKc_PAK STKc_YSK4 Dimer_Tnp_hAT DUF3134 PTKc_EphR_B TCP 23DHB-AMP_Ig Peptidase_C19 superfamily	STKc_MAPKKK _Byr2_like BetA PTKc_CCK4 PTKc_Aatyk PTKc_ALK_LTK PTKc_Axl_like PTKc_EphR_A1 0 PTKc_Fer PTKc_Fes_like PTKc_Frk_like PTKc_HER4 PTKc_Jak2_Jak3_rpt2 PTKc_Lck_Blk PTKc_RET PTKc_Srm_Brk PTKc_Syk_like PTKc_Jak1_rpt2 PTKc_Ror2 STKc_PAK_II PK_STRAD PTKc_Axl STKc_MEKK3_ike_1 STKc_PAK6 PTK_Jak_rpt1 STKc_Nek PLN03202 PRK05192 RNase_H_like superfamily	PLN03185 STKc_PAK4 DUF4413 PTKc_HER3 Methyltransf_16 Auxin_resp NAD_binding_8 superfamily AP2 RVT_2 Smc GMC_oxred_N pepsin_retropepsin_like superfamily Rv0697 xylanase_inhibit or_I_like TAXI_N zf-BED ZnF_BED SMC_prok_B cnd41_like WEMBL PLN02661 PRK07656 flhF PLN00020 RNase_HI_RT_Ty1 UCH chap_CCT_eta Ffh	pepsin_A_like_p lant HLH RRM_SF superfamily RRM3_RBM19 RRM2_MRD1 PTKc_CSF-1R AFD_class_I superfamily PRK04176 STKc_CMGC Thi4 TIGR00292 Ycf1 AAA superfamily chap_CCT_gam ma chaperonin_like superfamily PTKc_FGFR PTKc_FGFR1 PTKc_InsR_like STKc_CDKL2_3 STKc_MAPK4_6 STKc_TAO2 STKc_TAO3 FtsY PRK00771 PRK10416 PRK10867	TCP1_gamma PTZ00212 MRS6 AP2 superfamily rve STKc_CCRK STKc_MAK_like chap_CCT_alpha a chap_CCT_delta chap_CCT_epsilon chaperonin_type_I_II cpn60 Cpn60_TCP1 GroL PTKc_Aatyk1_A atyk3 PTKc_FGFR3 PTKc_InsR PTKc_Kit PTKc_PDGFR SPS1 STKc_CDK9_lik e STKc_MEKK3 STKc_PAK2 STKc_TAO STKc_TAO1 PRK12724 PRK12726 PRK14974 PTKc_Chk
--	--	---	--	--	---

PTKc_IGF-1R PTZ00031 rplB SMC_prok_A SRP54 STKc_CDK_like STKc_OSR1_SPAK STKc_PAK_I HLH superfamily PLN02479 PLN03102 PRK08162 COG4642 DUF936 Glyco_transf_8 nadB Peptidase_C19 STKc_CDK7 Ycf1 superfamily PTK_Jak1_rpt1 PTK_Ryk PTK_Tyk2_rpt1 PTKc_Fes PTKc_FGFR2 PTKc_FGFR4 _alpha	TCP1_beta TCP1_epsilon Fasciclin GRAS COG2112 COG4886 DUF1785 Glyco_transf_GTA_t ype superfamily Peptidase_C19D Peptidase_C19E Peptidase_C19F Peptidase_C19K Peptidase_C19O Peptidase_C19R PKc_MAPKK_plant _like PLN03150 PRK00911 PRK07560 PTZ00416 sdhA UBP12 AAA chap_CCT_theta	Fer4_NifH Fer4_NifH superfamily gag_pre-integr gag_pre-integr superfamily ParA PHA02517 PKc_MAPKK PLN02423 PTKc_Aatyk2 PTKc_PDGFR_beta PTZ00174 Ribosomal_L2 rpl2 rplB_bact SRP SRP54_euk STKc_CDK2_3 STKc_MEKK4 STKc_SLK_like TCP1_theta PRK14709 PTKc_Hck PTKc_Met_Ron	PTKc_Tyro3 STKc_CDK8_like STKc_CNK2-like STKc_PAK1 STKc_PAK3 STKc_PAK5 COG5078 UBCc AMP-binding CaiC DUF4106 PLN02330 PLN02501 PRK06187 PRK07529 PRK08275 PRK08276 PRK08316 ttLC_FACS_AEE21 _like COG4231 AdoMet_MTases superfamily APH_ChoK_like superfamily	EF2 FHY3 FusA GTP_EFTU IlvD LRR_8 Methyltransf_29 Ras_like_GTPase superfamily STKc_MAPK STKc_MST3_like DSPc PMM superfamily STKc_AGC STKc_CDK1_euk STKc_CDK10 Tra5 Transposase_22 chap_CCT_beta chap_CCT_zeta STKc_MAP4K3_like STKc_MAPKKK_Bc k1_like TCP1_zeta Branch	UBCc superfamily UQ_con 4CL DUF788 EntE PTKc_FAK PTKc_Fyn_Yrk PTKc_Src STKc_CDK12 STKc_FA2-like STKc_MEKK2 STKc_MST1_2 STKc_Nek8 ttLC_FACS_AIkK_li ke ttLC_FACS_like APP APP_MetAP superfamily IOR_alpha MORN TAXi_C TPP_enzymes superfamily TPP_IOR
--	--	--	---	--	--

