

Transposons: a blessing curse

Scott Teresi

Department of Horticulture
Genetics & Genome Sciences Program
IMPACTS

Monday 28th October, 2024



Table of Contents:

Personal Introduction

- B.S in Biology from The College of William & Mary (2019)
 - Minor in Computational & Applied Math & Statistics



Personal Introduction

- B.S in Biology from The College of William & Mary (2019)
 - Minor in Computational & Applied Math & Statistics
- Joined the Puzey Lab
 - Was able to collaborate with Pat's lab on the strawberry genome project



Article | [Open access](#) | Published: 25 February 2019

Origin and evolution of the octoploid strawberry genome

[Patrick P. Edger](#) , [Thomas J. Poorten](#), [Robert VanBuren](#), [Michael A. Hardigan](#), [Marivi Colle](#), [Michael R. McKain](#), [Ronald D. Smith](#), [Scott J. Teresi](#), [Andrew D. L. Nelson](#), [Ching Man Wai](#), [Elizabeth I. Alger](#), [Kevin A. Bird](#), [Alan E. Yocca](#), [Nathan Pumplin](#), [Shujun Ou](#), [Gil Ben-Zvi](#), [Avital Brodt](#), [Kobi Baruch](#), [Thomas Swale](#), [Lily Shiu](#), [Charlotte B. Acharya](#), [Glenn S. Cole](#), [Jeffrey P. Mower](#), [Kevin L. Childs](#), ... [Steven J. Knapp](#) 

[+ Show authors](#)

Nature Genetics 51, 541–547 (2019) | [Cite this article](#)

87k Accesses | 419 Citations | 396 Altmetric | [Metrics](#)

Personal Introduction

- B.S in Biology from The College of William & Mary (2019)
 - Minor in Computational & Applied Math & Statistics
- Joined the Puzey Lab
 - Was able to collaborate with Pat's lab on the strawberry genome project
- Joined Pat's Lab for the 2018 Plant Genomics REU
 - More on this later!



Personal Introduction

- B.S in Biology from The College of William & Mary (2019)
 - Minor in Computational & Applied Math & Statistics
- Joined the Puzey Lab
 - Was able to collaborate with Pat's lab on the strawberry genome project
- Joined Pat's Lab for the 2018 Plant Genomics REU
 - More on this later!
- Returned to Pat's Lab for my PhD in Fall 2019



What are TEs?

- Transposable elements (TEs, transposons) are mobile, repetitive, genetic entities.

What are TEs?

- Transposable elements (TEs, transposons) are mobile, repetitive, genetic entities.
- TEs are ubiquitous components of most genomes, and can often make up a large portion of the genome.

What are TEs?

- Transposable elements (TEs, transposons) are mobile, repetitive, genetic entities.
- TEs are ubiquitous components of most genomes, and can often make up a large portion of the genome.
- TEs are a major force of mutation and genome evolution.

TEs were discovered by Barbara McClintock

- Barbara McClintock discovered TEs in maize in the 1940s.



1

TEs were discovered by Barbara McClintock

- Barbara McClintock discovered TEs in maize in the 1940s.
- She observed that a certain region of DNA was consistently breaking and moving, resulting in kernel color changes.



2

TEs were discovered by Barbara McClintock

- Barbara McClintock discovered TEs in maize in the 1940s.
- She observed that a certain region of DNA was consistently breaking and moving, resulting in kernel color changes.
- She called these oddities “controlling elements”, due to their ability to create phenotypic variation.



2

TEs were discovered by Barbara McClintock

- Barbara McClintock discovered TEs in maize in the 1940s.
- She observed that a certain region of DNA was consistently breaking and moving, resulting in kernel color changes.
- She called these oddities “controlling elements”, due to their ability to create phenotypic variation.
- However, her work was largely met with skepticism because it was so far ahead of its time and limited to maize.



2

So why should we care about TEs?

So why should we care about TEs?

- TEs create genetic variation
and they are everywhere!

So why should we care about TEs?

- TEs create genetic variation and they are everywhere!
- TEs create cancer

Review Article | Published: 09 June 2017

Transposable elements in cancer

[Kathleen H. Burns](#) 

Nature Reviews Cancer 17, 415–424 (2017) | [Cite this article](#)

23k Accesses | 59 Altmetric | [Metrics](#)

So why should we care about TEs?

- TEs create genetic variation and they are everywhere!
- TEs create cancer
- TEs contribute to disease resistance in a variety of organisms

RESEARCH

Open Access



New reference genome sequences of hot pepper reveal the massive evolution of plant disease-resistance genes by retroduplication

Seungill Kim¹, Jeun Park^{1,2}, Seon-In Yeom³, Yong-Min Kim⁴, Eunyoung Seo¹, Ki-Tae Kim⁵, Myung-Shin Kim¹, Je Min Lee⁶, Kyeongchae Cheong^{2,3}, Ho-Sub Shin¹, Saet-Byul Kim¹, Koeun Han^{1,7}, Jundae Lee⁸, Minkyu Park², Hyun-Ah Lee¹, Hye-Young Lee¹, Youngsill Lee¹, Soohyun Oh¹, Joo Hyun Lee¹, Eunhye Choi¹, Eunbi Choi¹, So Eui Lee¹, Jongbum Jeon¹, Hyunbin Kim¹, Gobong Choi¹, Heunyeong Song⁹, Junil Lee¹, Sang-Choon Lee¹, Jin-Kyung Kwon^{1,7}, Hea-Young Lee^{1,7}, Namjin Koo¹, Yunji Hong¹, Ryan W. Kim⁴, Won-Hee Kang¹, Jin Hoe Huh¹, Byoung-Cheol Kang^{1,7}, Tae-Jin Yang¹, Yong-Hwan Lee^{2,3}, Jeffrey L. Bennetzen⁹ and Doli Choi^{1*}

So why should we care about TEs?

- TEs create genetic variation and they are everywhere!
- TEs create cancer
- TEs contribute to disease resistance in a variety of organisms
- TEs are associated with fruit color variation
 - Apples, grapes, peppers, blood oranges, strawberry and more

ARTICLE

<https://doi.org/10.1104/plc.114.09518>

OPEN

A high-quality apple genome assembly reveals the association of a retrotransposon and red fruit colour

Liyi Zhang¹, Jiang Hu², Xiaolei Han¹, Jingjing Li¹, Yuan Gao¹, Christopher M. Richards³, Caixia Zhang¹, Yi Tian¹, Guiming Liu⁴, Hera Gul¹, Dajiang Wang¹, Yu Tian², Chuanxin Yang², Minghui Meng², Gaopeng Yuan¹, Guodong Kang¹, Yonglong Wu¹, Kun Wang¹, Hengtao Zhang⁵, Depeng Wang² & Peihua Cong¹

The Plant Cell, Vol. 24: 1242–1255, March 2012, www.plantcell.org © 2012 American Society of Plant Biologists. All rights reserved.

Retrotransposons Control Fruit-Specific, Cold-Dependent Accumulation of Anthocyanins in Blood Oranges

Eugenio Butelli,^a Concetta Licciardello,^b Yang Zhang,^a Jianjun Liu,^c Steve Mackay,^a Paul Bailey,^a Giuseppe Reforgiato-Recupero,^b and Cathie Martin^{a,1}

^aJohn Innes Centre, Norwich NR4 7UH, United Kingdom

^bCentro di Ricerca per l'Agrumicoltura e le Colture Mediterranee, 95024 Acireale, Italy

^cSichuan Academy of Agricultural Sciences, Chengdu City, Sichuan 610066, China

Retrotransposon-Induced Mutations in Grape Skin Color

Shozo Kobayashi,^{1*} Nami Goto-Yamamoto,² Hirohiko Hirochika³

So why should we care about TEs?

- TEs create genetic variation and they are everywhere!
- TEs create cancer
- TEs contribute to disease resistance in a variety of organisms
- TEs are associated with fruit color variation
 - Apples, grapes, peppers, blood oranges, strawberry and more
- TEs created the peppered moth phenotype

LETTER

doi:10.1038/nature17951

The industrial melanism mutation in British peppered moths is a transposable element

Arjen E. van't Hof¹*, Pascal Campagne¹*, Daniel J. Rigden¹, Carl J. Yung¹, Jessica Lingley¹, Michael A. Quail², Neil Hall¹, Alistair C. Darby¹ & Ilik J. Saccheri¹



1

¹Image: BBC

So why should we care about TEs?

- TEs create genetic variation and they are everywhere!
- TEs create cancer
- TEs contribute to disease resistance in a variety of organisms
- TEs are associated with fruit color variation
 - Apples, grapes, peppers, blood oranges, strawberry and more
- TEs created the peppered moth phenotype
- Metabolite diversity and morphological variation in tomato

A Retrotransposon-Mediated Gene Duplication Underlies Morphological Variation of Tomato Fruit

Han Xiao,¹ Ning Jiang,² Erin Schaffner,^{1,3*} Eric J. Stockinger,³ Esther van der Knaap^{1†}

ARTICLE

<https://doi.org/10.1101/s41467-020-17874-2>

OPEN



The impact of transposable elements on tomato diversity

Marisol Domínguez¹, Elise Dugas¹, Médine Benchouaia², Basile Leducque¹, José M Jiménez-Gómez³, Vincent Colot^{1,5*} & Leandro Quadrana^{1,6}

So why should we care about TEs?

- TEs create genetic variation and they are everywhere!
- TEs create cancer
- TEs contribute to disease resistance in a variety of organisms
- TEs are associated with fruit color variation
 - Apples, grapes, peppers, blood oranges, strawberry and more
- TEs created the peppered moth phenotype
- Metabolite diversity and morphological variation in tomato
- Maize domestication, and photoperiod sensitivity

CACTA-like transposable element in *ZmCCT* attenuated photoperiod sensitivity and accelerated the postdomestication spread of maize

Qin Yang^{a,1}, Zhi Li^{a,1}, Wengiang Li^{b,1}, Lixia Ku^{b,1}, Chao Wang^a, Jianrong Ye^a, Kun Li^a, Ning Yang^b, Yipu Li^a, Tao Zhong^a, Jiansheng Li^a, Yanhui Chen^{c,2}, Jianbing Yan^{a,2}, Xiaohong Yang^{a,2}, and Mingliang Xu^{a,2}

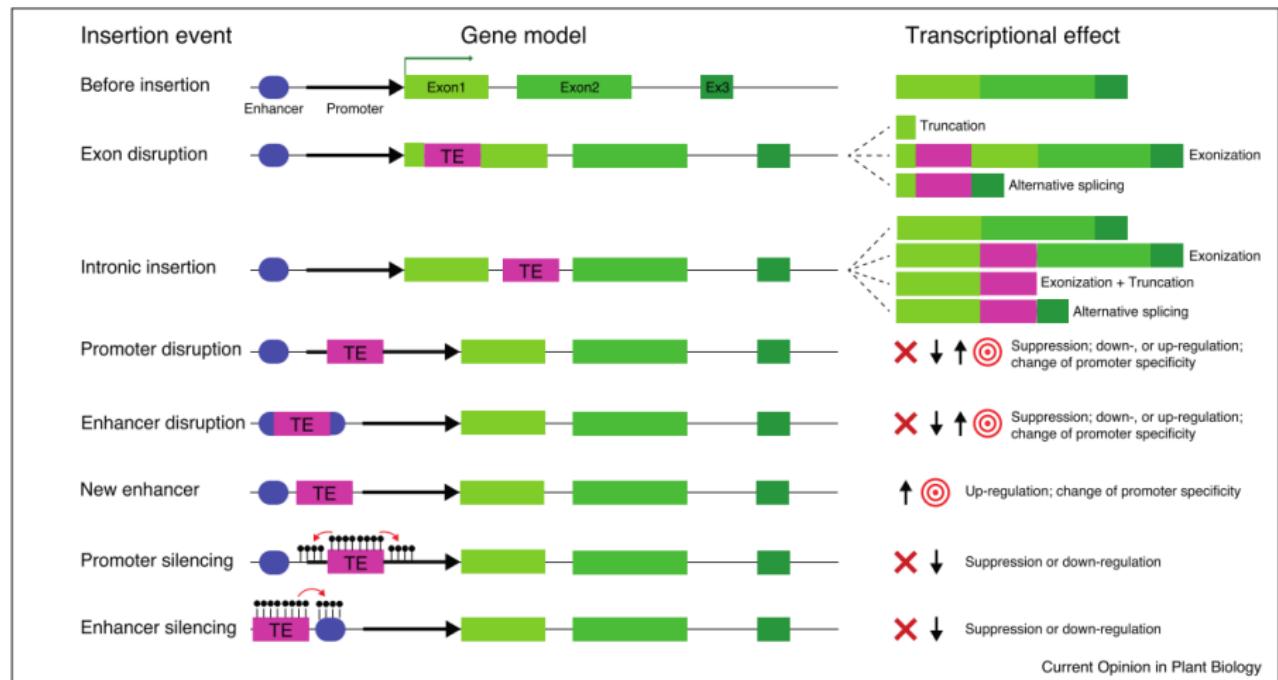
Identification of a functional transposon insertion in the maize domestication gene *tb1*

Anthony Studer¹, Qiong Zhao¹, Jeffrey Ross-Ibarra^{2,3} & John Doebley¹

So why should we care about TEs?

And many more (undiscovered) examples!

How does a TE create variation?



Current Opinion in Plant Biology

TEs have historically been maligned as “junk” DNA

- TEs were initially thought to be “junk” DNA, due to their apparent “selfish” nature. All they do is make more copies of themselves!



1

¹Image: [Forbes.com](https://www.forbes.com/sites/forbestechcouncil/2018/05/15/transposable-elements-the-selfish-gene/)

TEs have historically been maligned as “junk” DNA

- TEs were initially thought to be “junk” DNA, due to their apparent “selfish” nature. All they do is make more copies of themselves!
- Similarities between TEs and viruses cast them as genomic invaders.



1

¹Image: [Forbes.com](https://www.forbes.com/sites/forbestechcouncil/2018/03/06/transposable-elements-are-not-junk-dna/#:~:text=Transposable%20elements%20(TEs)%20have%20historically,more%20copies%20of%20themselves!)

TEs have historically been maligned as “junk” DNA

- TEs were initially thought to be “junk” DNA, due to their apparent “selfish” nature. All they do is make more copies of themselves!
- Similarities between TEs and viruses cast them as genomic invaders.
- Many of the new, emerging genomes in the 1990s and early 2000s such as *Homo sapiens* appeared to be filled with “dead” TEs.



1

¹Image: [Forbes.com](https://www.forbes.com/sites/forbestechcouncil/2018/05/08/the-truth-about-junk-dna/#:~:text=Junk%20DNA%20is%20the%20common,of%20the%20genome%20that%20has%20no%20function.)

TEs have historically been maligned as “junk” DNA

- TEs were initially thought to be “junk” DNA, due to their apparent “selfish” nature. All they do is make more copies of themselves!
- Similarities between TEs and viruses cast them as genomic invaders.
- Many of the new, emerging genomes in the 1990s and early 2000s such as *Homo sapiens* appeared to be filled with “dead” TEs.
- Their high-copy number inflates genome size and sequencing costs, as well as complicates sequence alignment



1

¹Image: [Forbes.com](https://www.forbes.com/sites/forbestechcouncil/2014/05/01/the-truth-about-junk-dna/#:~:text=Junk%20DNA%20is%20the%20name,of%20the%20genome%20that%20we%20don't%20use.)

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features
 - Location — close to, or away from, genes

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features
 - Location — close to, or away from, genes
 - Insertion preference

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features
 - Location — close to, or away from, genes
 - Insertion preference
 - Selective forces to remove

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features
 - Location — close to, or away from, genes
 - Insertion preference
 - Selective forces to remove
 - Methylation status

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features
 - Location — close to, or away from, genes
 - Insertion preference
 - Selective forces to remove
 - Methylation status
 - Likelihood to capture, duplicate, and shuffle gene sequences

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features
 - Location — close to, or away from, genes
 - Insertion preference
 - Selective forces to remove
 - Methylation status
 - Likelihood to capture, duplicate, and shuffle gene sequences
 - Activity — dead or alive

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features
 - Location — close to, or away from, genes
 - Insertion preference
 - Selective forces to remove
 - Methylation status
 - Likelihood to capture, duplicate, and shuffle gene sequences
 - Activity — dead or alive
 - Intact — full-length or fragmented

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features
 - Location — close to, or away from, genes
 - Insertion preference
 - Selective forces to remove
 - Methylation status
 - Likelihood to capture, duplicate, and shuffle gene sequences
 - Activity — dead or alive
 - Intact — full-length or fragmented
 - Copy number — low copy or abundant

TEs are difficult to study

- There are a multitude of different TE types, each with varying characteristics and distributions!
 - Size
 - Replication mechanism
 - Autonomous vs. non-autonomous
 - Structural features
 - Location — close to, or away from, genes
 - Insertion preference
 - Selective forces to remove
 - Methylation status
 - Likelihood to capture, duplicate, and shuffle gene sequences
 - Activity — dead or alive
 - Intact — full-length or fragmented
 - Copy number — low copy or abundant
 - TE insertions into other TEs

TEs are difficult to study

Classification		Structure	TSD	Code	Occurrence
Order	Superfamily				
Class I (retrotransposons)					
LTR	Copia	→ GAG AP INT RT RH →	4–6	RLC	P, M, F, O
	Gypsy	→ GAG AP RT RH INT →	4–6	RLG	P, M, F, O
	Bel-Pao	→ GAG AP RT RH INT →	4–6	RLB	M
	Retrovirus	→ GAG AP RT RH INT ENV →	4–6	RLR	M
	ERV	→ GAG AP RT RH INT ENV →	4–6	RLE	M
DIRS	DIRS	← GAG AP RT RH YR ←	0	RYD	P, M, F, O
	Ngaro	→ GAG AP RT RH YR ← → →	0	RYN	M, F
	VIPER	→ GAG AP RT RH YR → → →	0	RYV	O
PLE	Penelope	← → RT EN → →	Variable	RPP	P, M, F, O
LINE	R2	— RT EN —	Variable	RIR	M
	RTE	— APE RT —	Variable	RIT	M
	Jockey	— ORFI — APE RT —	Variable	RIJ	M
	L1	— ORFI — APE RT —	Variable	RIL	P, M, F, O
	I	— ORFI — APE RT RH —	Variable	RII	P, M, F
SINE	tRNA	—	Variable	RST	P, M, F
	7SL	—	Variable	RSL	P, M, F
	5S	—	Variable	RSS	M, O

Structural features



Protein coding domains

AP, Aspartic protease	APE, Apurinic endonuclease	ATP, Packaging ATPase	C-INT, C-integrase	CYP, Cysteine protease	EN, Endonuclease
ENV, Envelope protein	GAG, Capsid protein	HEL, Helicase	INT, Integrase	ORF, Open reading frame of unknown function	
POL B, DNA polymerase B	RH, RNase H	RPA, Replication protein A (found only in plants)	RT, Reverse transcriptase		
Tase, Transposase (* with DDE motif)		YR, Tyrosine recombinase	Y2, YR with YY motif		

TEs are difficult to study

Class II (DNA transposons) - Subclass 1					
TIR	Tc1-Mariner		TA	DTT	P, M, F, O
	hAT		8	DTA	P, M, F, O
	Mutator		9-11	DTM	P, M, F, O
	Merlin		8-9	DTE	M, O
	Transib		5	DTR	M, F
	P		8	DTP	P, M
	PiggyBac		TTAA	DTB	M, O
	PIF-Harbinger		3	DTH	P, M, F, O
	CACTA		2-3	DTC	P, M, F
Crypton	Crypton		0	DYC	F
Class II (DNA transposons) - Subclass 2					
Helitron	Helitron		0	DHH	P, M, F
Maverick	Maverick		6	DMM	M, F, O
Structural features					
Long terminal repeats			Terminal inverted repeats		Coding region
Diagnostic feature in non-coding region					Non-coding region
					Region that can contain one or more additional ORFs
Protein coding domains					
AP, Aspartic proteinase	APE, Apurinic endonuclease	ATP, Packaging ATPase	C-INT, C-integrase	CYP, Cysteine protease	EN, Endonuclease
ENV, Envelope protein	GAG, Capsid protein	HEL, Helicase	INT, Integrase	ORF, Open reading frame of unknown function	
POL B, DNA polymerase B	RH, RNase H	RPA, Replication protein A (found only in plants)	RT, Reverse transcriptase		
Tase, Transposase (* with DDE motif)		YR, Tyrosine recombinase	Y2, YR with YY motif		

Our understanding of TEs is limited by the tools we have to study them.

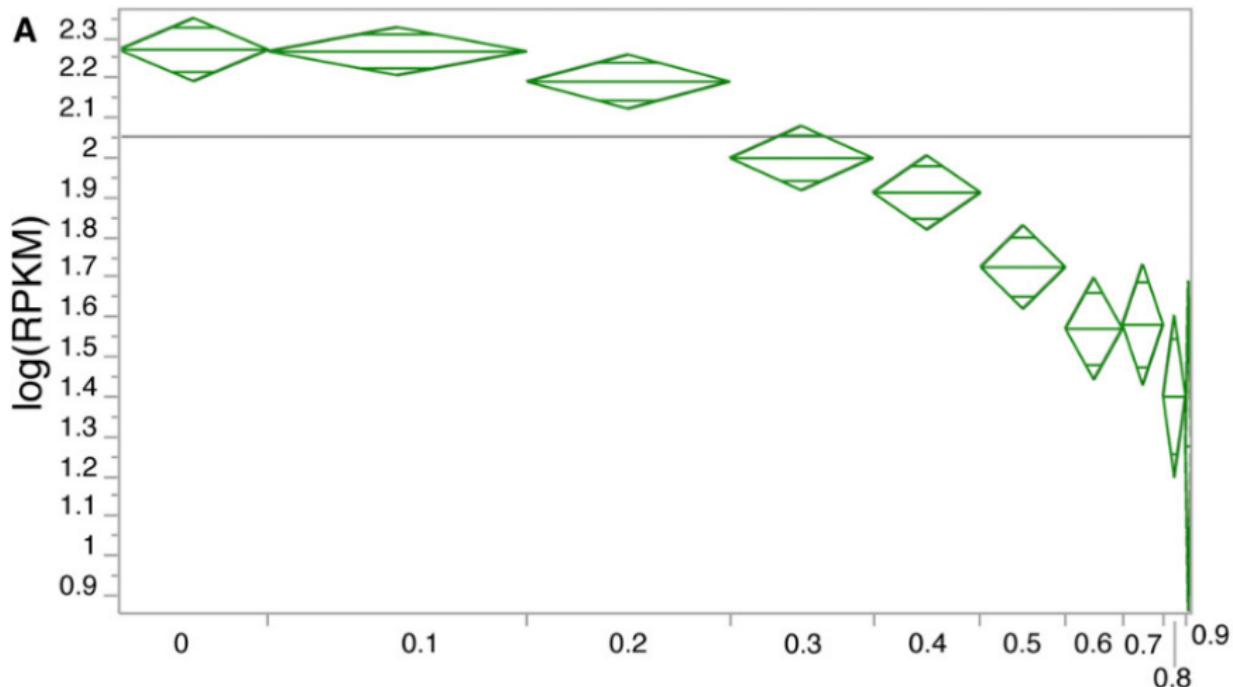
Letter

Epigenetic silencing of transposable elements: A trade-off between reduced transposition and deleterious effects on neighboring gene expression

Jesse D. Hollister and Brandon S. Gaut¹

Department of Ecology and Evolutionary Biology, University of California, Irvine, Irvine, California 92697-2525, USA

TE Density Software — Inspiration



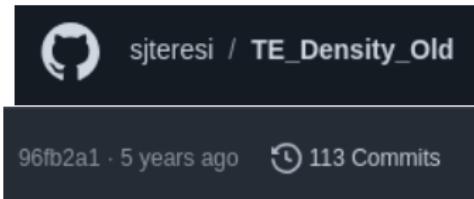
TE Density Software — REU Edition

- I joined Pat's lab for the summer of 2018 as part of the Plant Genomics REU.
- I wanted to develop software to calculate TE density around genes, and better understand the different associations between TEs and genes.



TE Density Software — REU Edition

- I joined Pat's lab for the summer of 2018 as part of the Plant Genomics REU.
- I wanted to develop software to calculate TE density around genes, and better understand the different associations between TEs and genes.



TE Density Software — REU Edition

I managed to get a working prototype by the end of the summer and had some cool results

KEGG Pathways

<i>pathway ID</i>	<i>pathway description</i>	<i>count in gene set</i>	<i>false discovery rate</i>
03010	Ribosome	30	1.5e-13
00500	Starch and sucrose metabolism	14	0.000142

PFAM Protein Domains

<i>pathway ID</i>	<i>pathway description</i>	<i>count in gene set</i>	<i>false discovery rate</i>
PF00982	Glycosyltransferase family 20	6	2.23e-08
PF02358	Trehalose-phosphatase	6	2.23e-08

INTERPRO Protein Domains and Features

<i>pathway ID</i>	<i>pathway description</i>	<i>count in gene set</i>	<i>false discovery rate</i>
IPR001830	Glycosyl transferase, family 20	6	3.81e-08
IPR003337	Trehalose-phosphatase	6	3.81e-08
IPR006379	HAD-superfamily hydrolase, subfamily IIB	4	0.000485

pace

Teresi *et al.* *Mobile DNA* (2022) 13:11
<https://doi.org/10.1186/s13100-022-00264-4>

Mobile DNA

SOFTWARE

Open Access

TE Density: a tool to investigate the biology of transposable elements



Scott J. Teresi^{1,2}, Michael B. Teresi³ and Patrick P. Edger^{1,2*}

TE Density Software — Development

 sjteresi Update text and column names for Intragenic values, when parsed with ... 

219dd8a · 3 weeks ago  289 Commits

- For every gene, report the amount of TE-occupied base-pairs normalized by a given search space. I.e Gene X has 50% of the 5000 BP upstream area occupied by TEs, and then that number can further be broken down into the various TE types

Algorithm 1: Calculate Density

CalculateDensity

inputs : set of windows W

output : density ρ

dataset: pseudomolecule C

get *genes* from C

get *transposons* from C

initialize overlaps O

foreach *geneName* $g_j \in genes$ **do**

foreach *window* $w_j \in W$ **do**

 // see (1) (2) (3)

$o = Overlap(g_j, w_j, transposons)$

$O[g_j, w_j] = o$

initialize densities ρ

foreach *geneName* $g_j \in genes$ **do**

foreach *TE identity* $t_j \in transposons$ **do**

foreach *window* $w_j \in W$ **do**

 // see (4)

$d = Density(g_j, t_j, w_j, O_j)$

$\rho[g_j, w_j, t_j] = d$

TE Density Software — Development

- For every gene, report the amount of TE-occupied base-pairs normalized by a given search space. I.e Gene X has 50% of the 5000 BP upstream area occupied by TEs, and then that number can further be broken down into the various TE types
- Software really isn't that complex (just a lot of iteration), it was mostly just a matter of program architecture and tying everything together.

Algorithm 1: Calculate Density

```
CalculateDensity
  inputs : set of windows  $W$ 
  output : density  $\rho$ 
  dataset: pseudomolecule  $C$ 
  get genes from  $C$ 
  get transposons from  $C$ 
  initialize overlaps  $O$ 
  foreach geneName  $g_j \in genes$  do
    foreach window  $w_j \in W$  do
      // see (1) (2) (3)
       $o = Overlap(g_j, w_j, transposons)$ 
       $O[g_j, w_j] = o$ 
  initialize densities  $\rho$ 
  foreach geneName  $g_j \in genes$  do
    foreach TE identity  $t_j \in transposons$  do
      foreach window  $w_j \in W$  do
        // see (4)
         $d = Density(g_j, t_j, w_j, O_j)$ 
         $\rho[g_j, w_j, t_j] = d$ 
```

TE Density Software — Development

- For every gene, report the amount of TE-occupied base-pairs normalized by a given search space. I.e Gene X has 50% of the 5000 BP upstream area occupied by TEs, and then that number can further be broken down into the various TE types
- Software really isn't that complex (just a lot of iteration), it was mostly just a matter of program architecture and tying everything together.
- Insane 3D matrix to store all of the data (e.g $20 \times 15 \times 2$ for each gene of 108,000 in strawberry)

Algorithm 1: Calculate Density

CalculateDensity

inputs : set of windows W

output : density ρ

dataset: pseudomolecule C

get *genes* from C

get *transposons* from C

initialize overlaps O

foreach *geneName* $g_j \in genes$ **do**

foreach *window* $w_j \in W$ **do**

 // see (1) (2) (3)

$o = Overlap(g_j, w_j, transposons)$

$O[g_j, w_j] = o$

initialize densities ρ

foreach *geneName* $g_j \in genes$ **do**

foreach *TE identity* $t_j \in transposons$ **do**

foreach *window* $w_j \in W$ **do**

 // see (4)

$d = Density(g_j, t_j, w_j, O_j)$

$\rho[g_j, w_j, t_j] = d$

TE Density Software — Notable Improvements

- Speed — From a matter of weeks to 1-3 days

TE Density Software — Notable Improvements

- Speed — From a matter of weeks to 1-3 days
- Flexibility — Can be run on any genome

TE Density Software — Notable Improvements

- Speed — From a matter of weeks to 1-3 days
- Flexibility — Can be run on any genome
- Testing — Unit tests and integration tests

TE Density Software — Notable Improvements

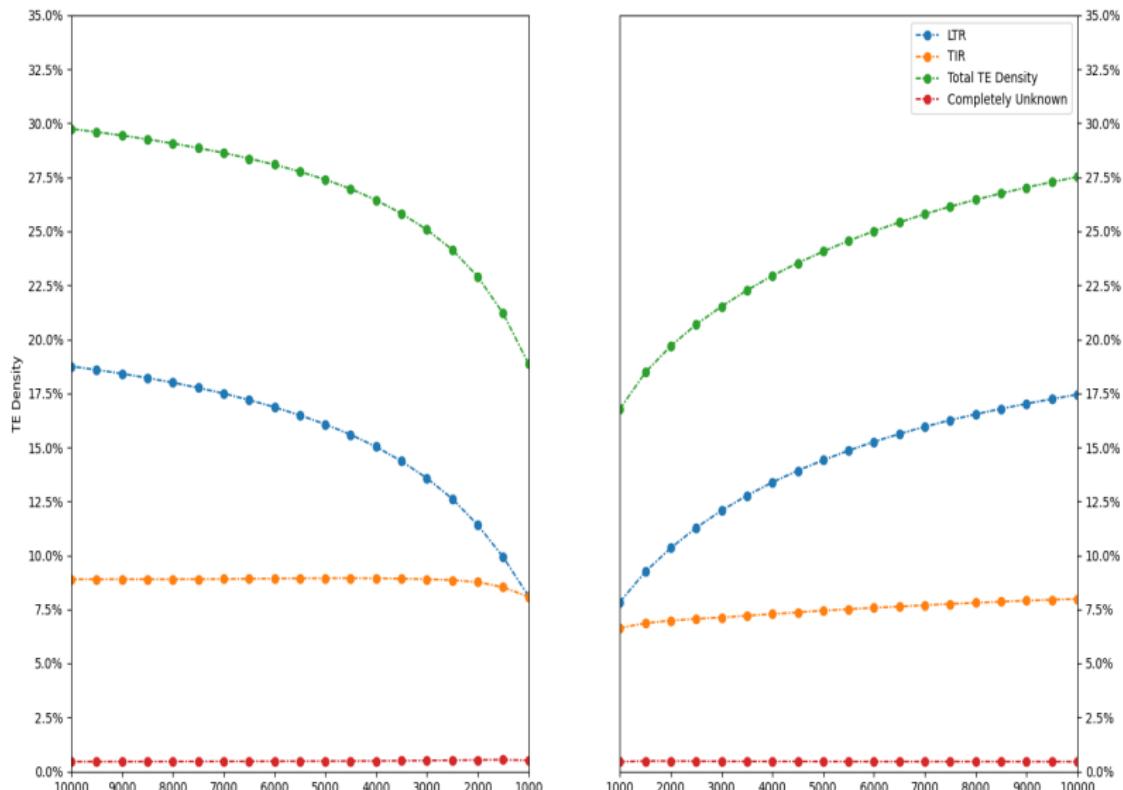
- Speed — From a matter of weeks to 1-3 days
- Flexibility — Can be run on any genome
- Testing — Unit tests and integration tests
- Documentation — README and example figures

TE Density Software — Notable Improvements

- Speed — From a matter of weeks to 1-3 days
- Flexibility — Can be run on any genome
- Testing — Unit tests and integration tests
- Documentation — README and example figures
- Reproducibility — Python virtual environment

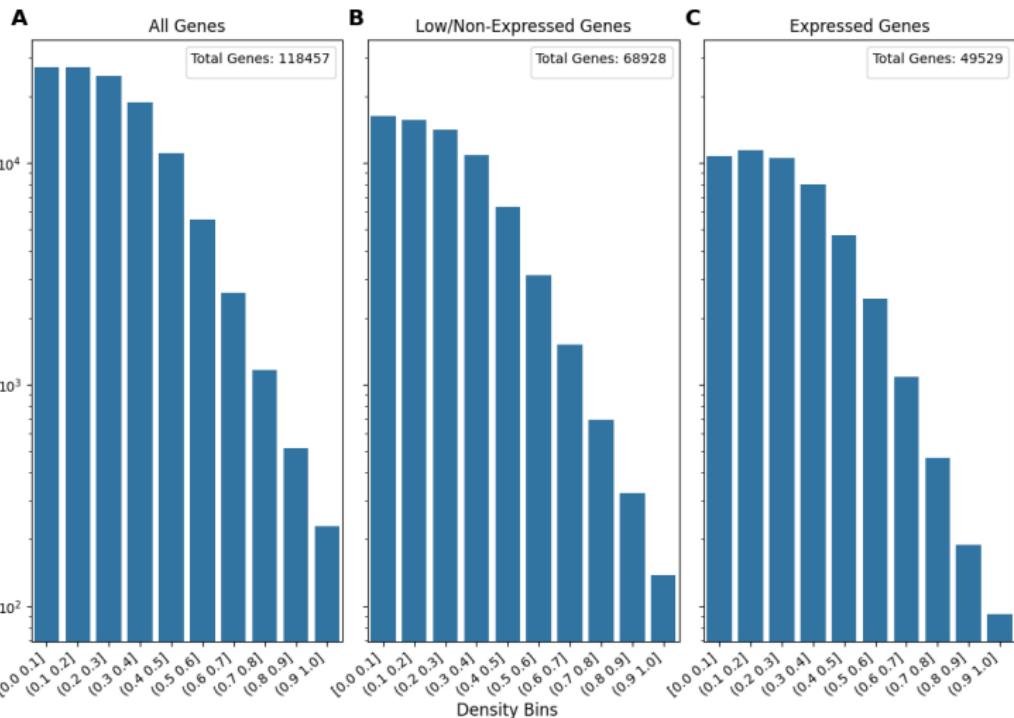
You can examine genome-wide trends of TE-gene relationships

Average TE Density of All Genes as a Function of Window Size and Location



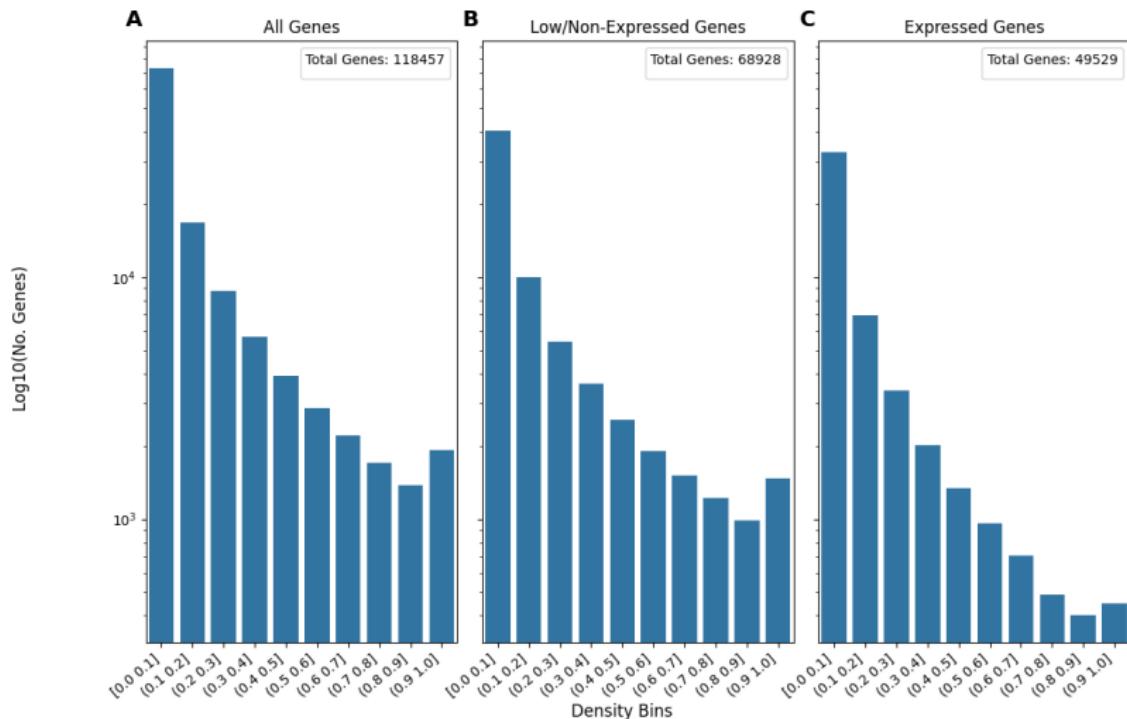
Re-examine the relationship with gene expression

No. Genes Binned by TIR TE 5KB Upstream Density According to Expression Profile

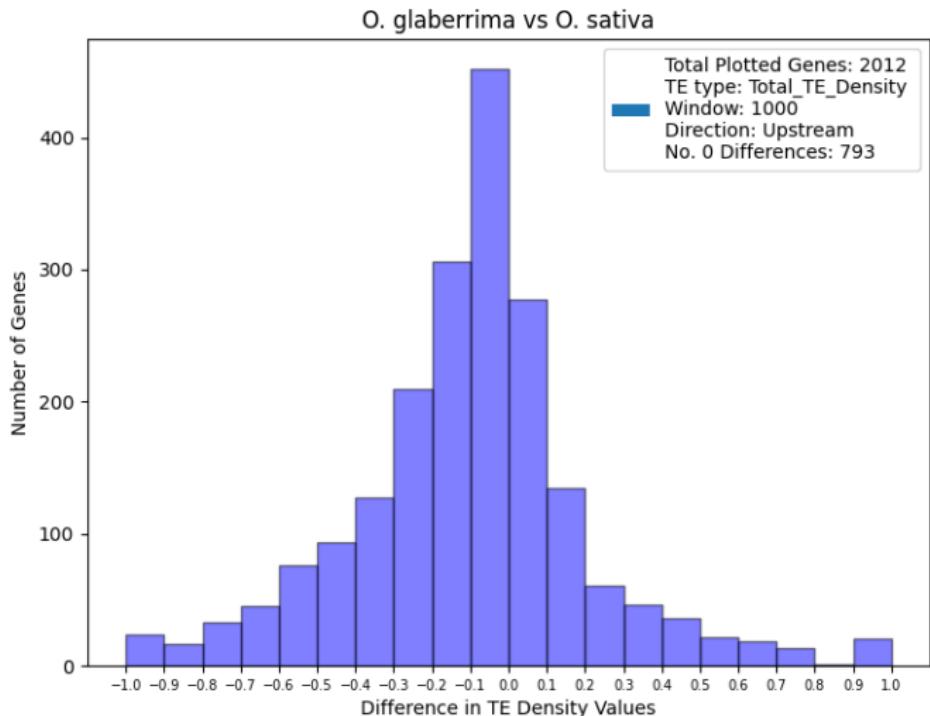


Re-examine the relationship with gene expression

No. Genes Binned by LTR TE 5KB Upstream Density According to Expression Profile



Identify TE differences between closely-related genes



What kinds of genes are associated with TEs?

More on this later!

PANTHER GO-Slim Biological Process	Fold Enrichment	FDR
Unclassified (UNCLASSIFIED)	1.11	0.00552
biological process (GO:0008150)	0.49	0.00276
cellular process (GO:0009987)	0.48	0.00687
organic substance metabolic process (GO:0071704)	0.46	0.0279
metabolic process (GO:0008152)	0.46	0.0175
primary metabolic process (GO:0044238)	0.44	0.0283
nitrogen compound metabolic process (GO:0006807)	0.42	0.0282
cellular metabolic process (GO:0044237)	0.42	0.0157

Strawberry TE Project — strawberry as a system

- Strawberries are octoploid (8 copies of each gene), and have a complex genome.

Strawberry TE Project — strawberry as a system

- Strawberries are octoploid (8 copies of each gene), and have a complex genome.
- Strawberries have a short history of domestication.

Strawberry TE Project — strawberry as a system

- Strawberries are octoploid (8 copies of each gene), and have a complex genome.
- Strawberries have a short history of domestication.
- Strawberries have a rich set of genomics resources, making it good system for studying polyploidy, domestication, and TE-gene relationships.

Strawberry is an extremely recent domesticate

- Cultivated strawberry (*Fragaria x ananassa*) is a very recent domesticate, with the first cultivated strawberry appearing in the 18th century in the Royal Gardens of Versailles.



Sources: Edger et al (2019) - Nature Genetics & Darrow, G. M. (1966) - The Strawberry

Image: [NC State Extension](#)

Strawberry is an extremely recent domesticate

- Cultivated strawberry (*Fragaria x ananassa*) is a very recent domesticate, with the first cultivated strawberry appearing in the 18th century in the Royal Gardens of Versailles.
- French Botanist Antoine Nicolas Duchesne is credited with its genesis; he created it by crossing the two New World species *Fragaria virginiana* and *Fragaria chiloensis*.



Sources: Edger et al (2019) - Nature Genetics & Darrow, G. M. (1966) - The Strawberry

Image: [NC State Extension](#)

Strawberry TE Project

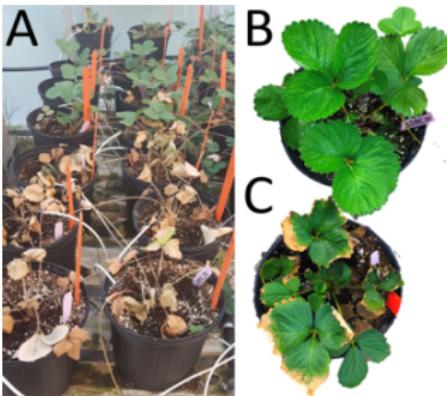
Given my findings from my REU, I am interested in understanding the role of TEs in strawberry domestication.

I do this by leveraging the cultivated strawberry genome as well as one of its wild progenitor species to compare and contrast TEs and genes.

A Comparable Wild Progenitor — *Fragaria chiloensis* “Del Norte”



- *Fragaria chiloensis* is a wild strawberry species that is a progenitor of the cultivated strawberry.
- It is also an octoploid species, and is quite hardy.
- It is found along the Pacific coast of North America.
- Our lab has sequenced and assembled the genome of this species, “Del Norte”



Strawberry TE Project

- What were my questions?

Strawberry TE Project

- What were my questions?
 - Do TEs shape phenotypes in strawberry?

Strawberry TE Project

- **What were my questions?**

- Do TEs shape phenotypes in strawberry?
- What kinds of genes are enriched with TEs? Are any of these genes associated with domestication phenotypes?

Strawberry TE Project

- **What were my questions?**

- Do TEs shape phenotypes in strawberry?
- What kinds of genes are enriched with TEs? Are any of these genes associated with domestication phenotypes?
- What TE Density differences exist between wild and cultivated strawberries?

Strawberry TE Project

- **What were my questions?**

- Do TEs shape phenotypes in strawberry?
- What kinds of genes are enriched with TEs? Are any of these genes associated with domestication phenotypes?
- What TE Density differences exist between wild and cultivated strawberries?
- Can I generate a list of candidate TE-impacted genes for future study?

What datasets did I generate?

- TE pangenome with EDTA — What's that?
- Ortholog predictions
- TE Density calculations for each gene
- GO enrichments for TE-dense genes

Pangenomes — What and Why?

- Pangenomes are useful for understanding the genetic diversity of a species.

Pangenomes — What and Why?

- Pangenomes are useful for understanding the genetic diversity of a species.
- A pangenome is a collection of all the genes or TEs in a group of species, and it delineates the shared and unshared portions.

Pangenomes — What and Why?

- Pangenomes are useful for understanding the genetic diversity of a species.
- A pangenome is a collection of all the genes or TEs in a group of species, and it delineates the shared and unshared portions.
- TE pangenomes can be better than individual genome annotations because you have an enhanced ability to filter out false positives and negatives, particularly low-copy or fragmented TEs.

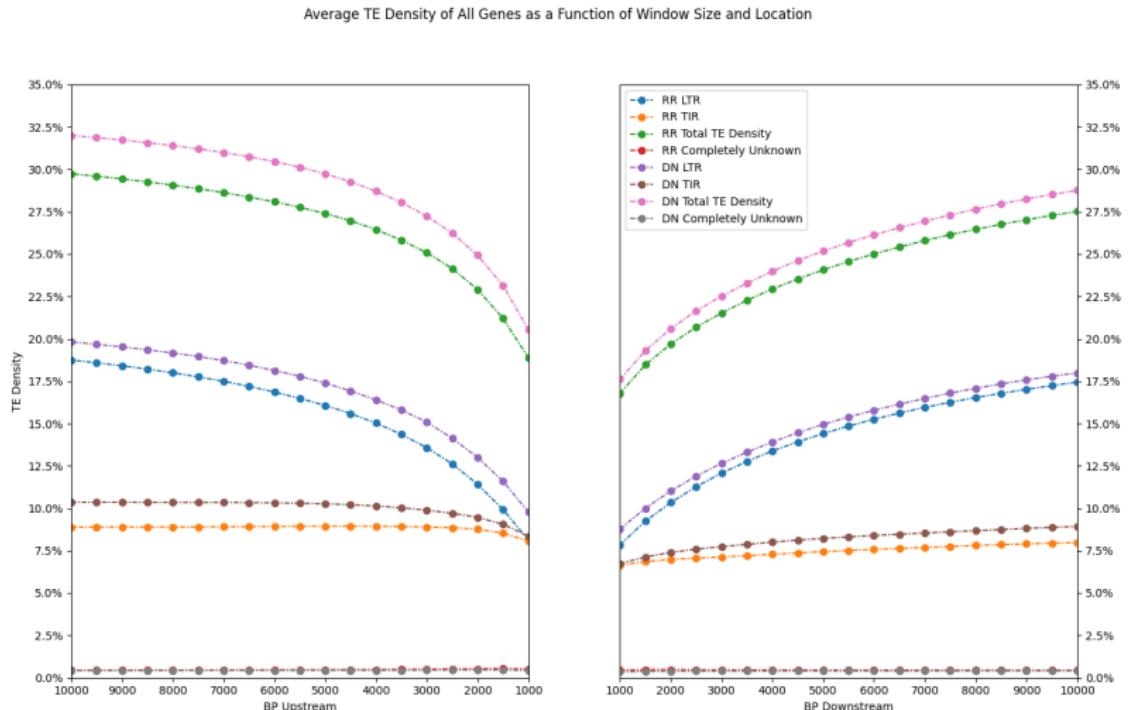
Pangenomes — What and Why?

- Pangenomes are useful for understanding the genetic diversity of a species.
- A pangenome is a collection of all the genes or TEs in a group of species, and it delineates the shared and unshared portions.
- TE pangenomes can be better than individual genome annotations because you have an enhanced ability to filter out false positives and negatives, particularly low-copy or fragmented TEs.
- I wanted a consistent and comparable TE annotation for both species, where I could compare and contrast differing TE families (which are defined by sequence similarity).

TE pangenome reveals subtle differences between wild and cultivated strawberry

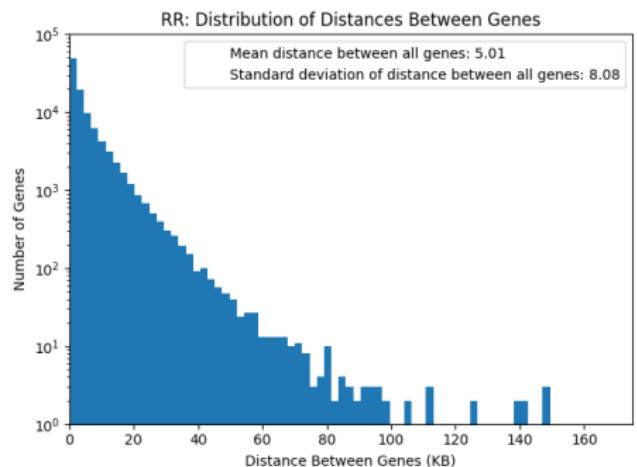
Class	Subtype	H4	DN	RR
LINE	L1	0.57%	0.46%	0.47%
LTR	Copia	3.46%	3.67%	3.51%
	Gypsy	5.05%	7.75%	7.90%
	unknown	7.62%	13.72%	12.57%
SINE	tRNA	0.00%	0.00%	0.00%
	unknown	0.37%	0.48%	0.56%
TIR	CACTA (En/Spm)	3.71%	4.88%	3.86%
	Mutator	2.43%	2.35%	2.16%
	PIF Harbinger	0.97%	0.91%	0.91%
	Tc1 Mariner	0.04%	0.02%	0.03%
	hAT	1.85%	2.16%	1.93%
Low Complexity	None	0.01%	0.07%	0.02%
nonLTR	pararetrovirus	0.13%	0.17%	0.16%
nonTIR	helitron	1.87%	1.17%	1.24%
rDNA	45S	0.01%	0.01%	0.01%
Repeat Fragment	None	0.54%	0.45%	0.56%
Total		28.62%	38.28%	35.88%

Cultivated strawberry is less TE-dense near genes

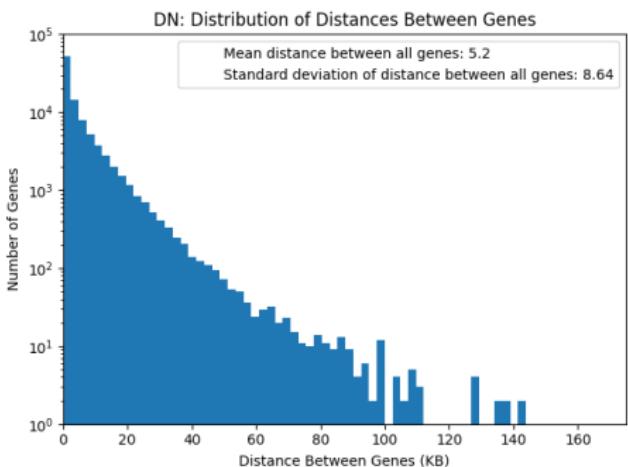


Cultivated strawberry is more gene-dense

Cultivated



Wild



What are the functional enrichments of TE-dense genes? How much do they differ?

- **What were my questions?**

- Do TEs shape phenotypes in strawberry?
- What TE Density differences exist between wild and cultivated strawberries?
- Can I generate a list of candidate TE-impacted genes for future study?
- What kinds of genes are enriched with TEs? Are any of these genes associated with domestication phenotypes?

What are the functional enrichments of TE-dense genes? How much do they differ?

GO Term	P-Value
Unique to Wild Strawberry	
Pectic galactan metabolic process	0.04842
Starch catabolic process	0.00785
Peptidoglycan biosynthetic process	0.04842
Tricarboxylic acid cycle	0.03893
Leaf development	0.00014
Embryo development ending in seed dormancy	0.00697
Regulation of developmental process	0.01241
Maltose biosynthetic process	0.04842
Malate metabolic process	0.00971
Circadian regulation of gene expression	0.01794
Histone H3-K9 demethylation	0.00342
Plant-type secondary cell wall biogenesis	0.00051
Cellular water homeostasis	0.006800
Pigment metabolic process	0.02453
Plant ovule development	0.02538
Negative regulation of flower development	0.04217

What are the functional enrichments of TE-dense genes? How much do they differ?

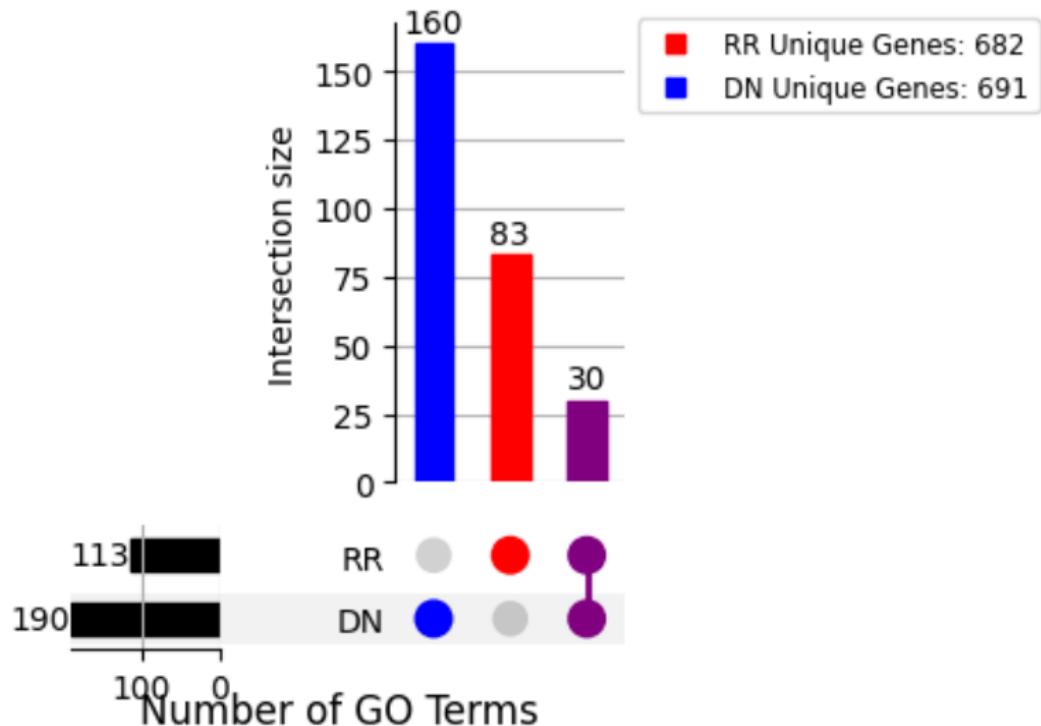
GO Term	P-Value
Unique to Domesticated Strawberry	
Fructose 6-phosphate metabolic process	0.00776
Brassinosteroid homeostasis	0.03174
Response to flooding	0.02728
Regulation of jasmonic acid biosynthetic process	0.03944
Abscisic acid homeostasis	0.00882
Reponse to microbial phytotoxin	0.03944
Flower development	0.01750
Seed development	0.01059
Shoot apical meristem development	0.02723
Glycine betain biosynthetic process from choline	0.00882
UDP-glucose metabolic process	0.02725

What are the functional enrichments of TE-dense genes? How much do they differ?

GO Term	P-Value	
	Wild vs. Domesticated	
Cellular defense response	0.00234	0.00305
Response to osmotic stress	0.02834	0.02382
Response to heat	0.01600	0.01904
Gibberellin biosynthetic process	0.00623	0.04642
Regulation of gene expression	0.00440	0.00440
Vegetative to Reproductive Phase Transition of Meristem	0.00883	0.00592

What are the functional enrichments of TE-dense genes? How much do they differ?

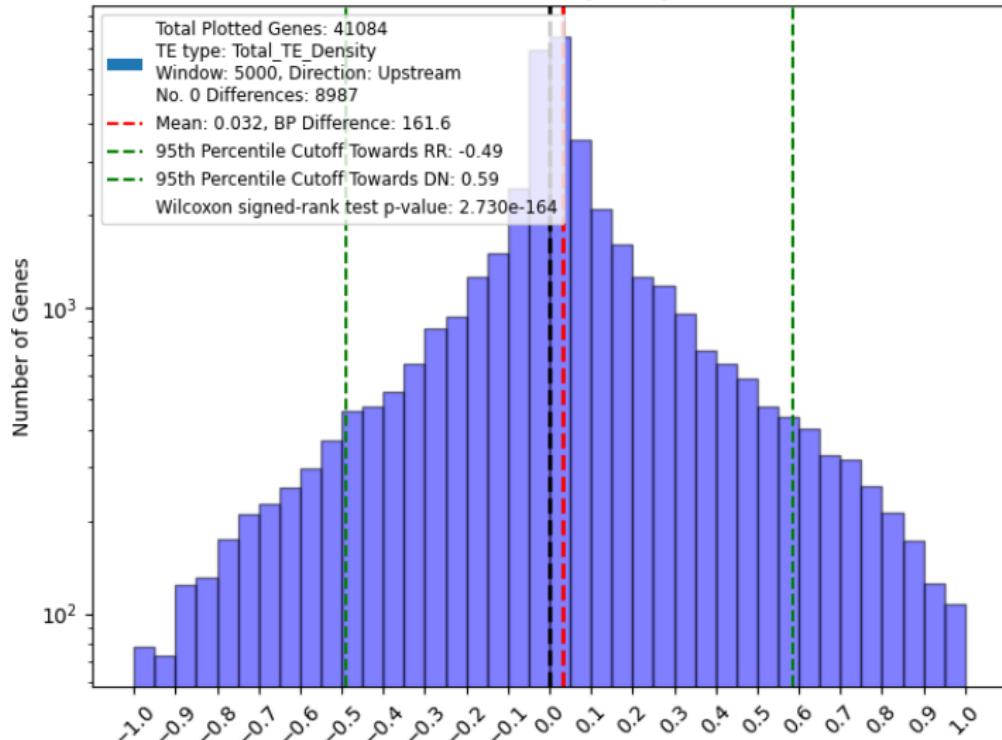
Total TE Density 5000 Upstream



Syntelogs show major differences in TE density

All TEs

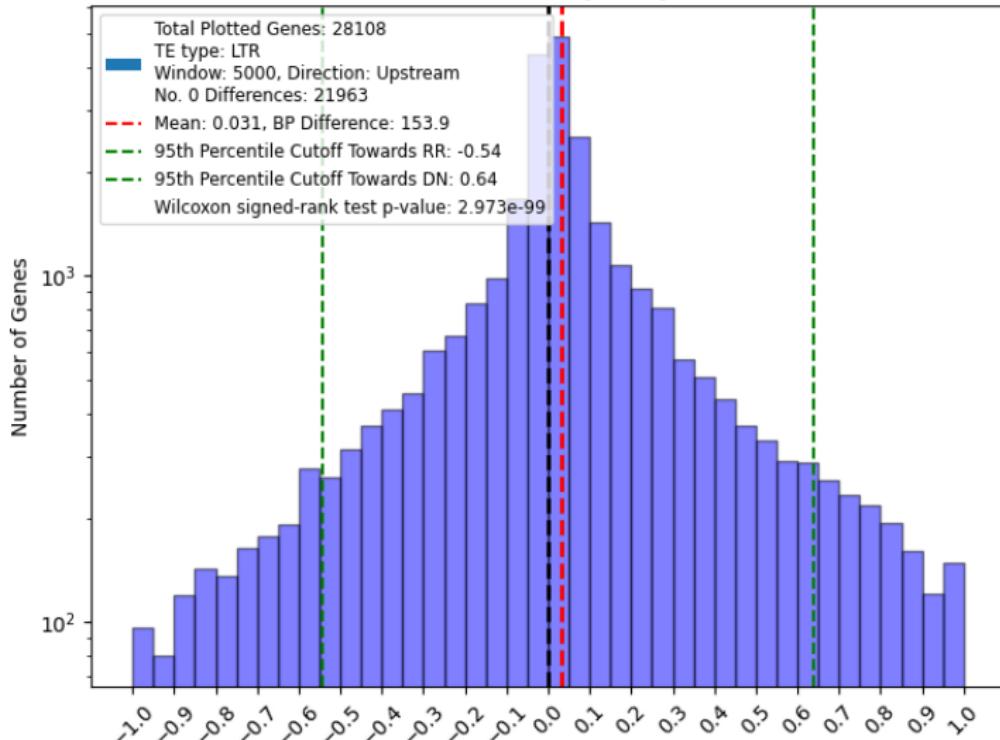
Del Norte vs Royal Royce



Syntelogs show major differences in TE density

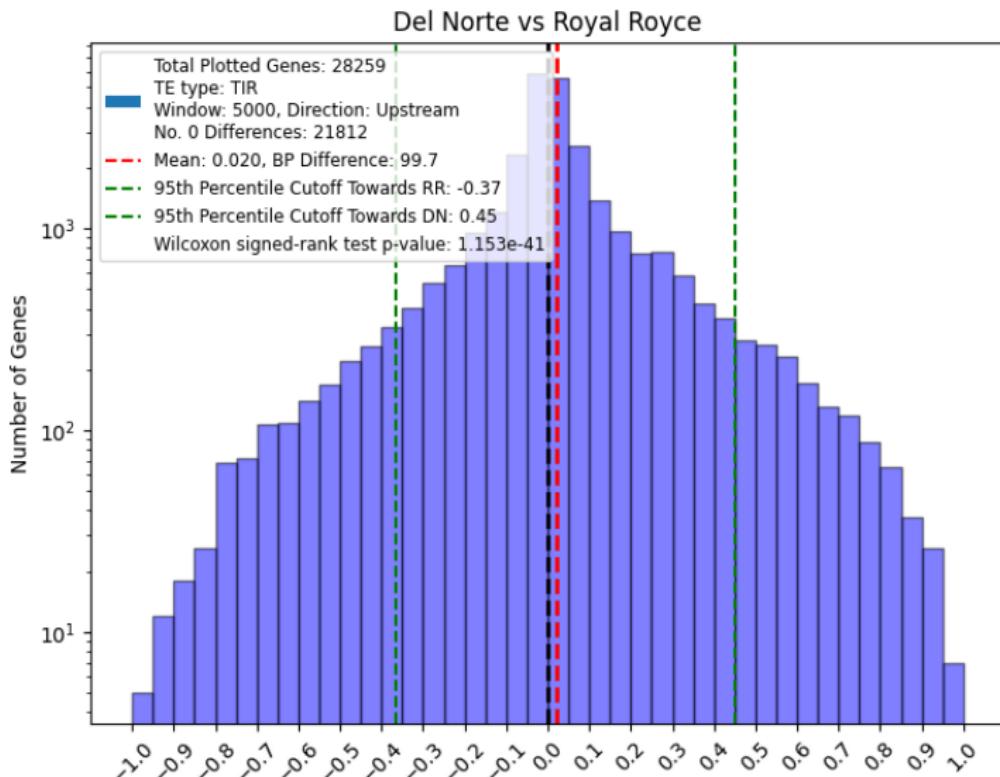
LTR TEs

Del Norte vs Royal Royce



Syntelogs show major differences in TE density

TIR TEs



Syntelogs show major differences in TE density

Cultivated strawberry has a less TEs around positionally conserved genes.

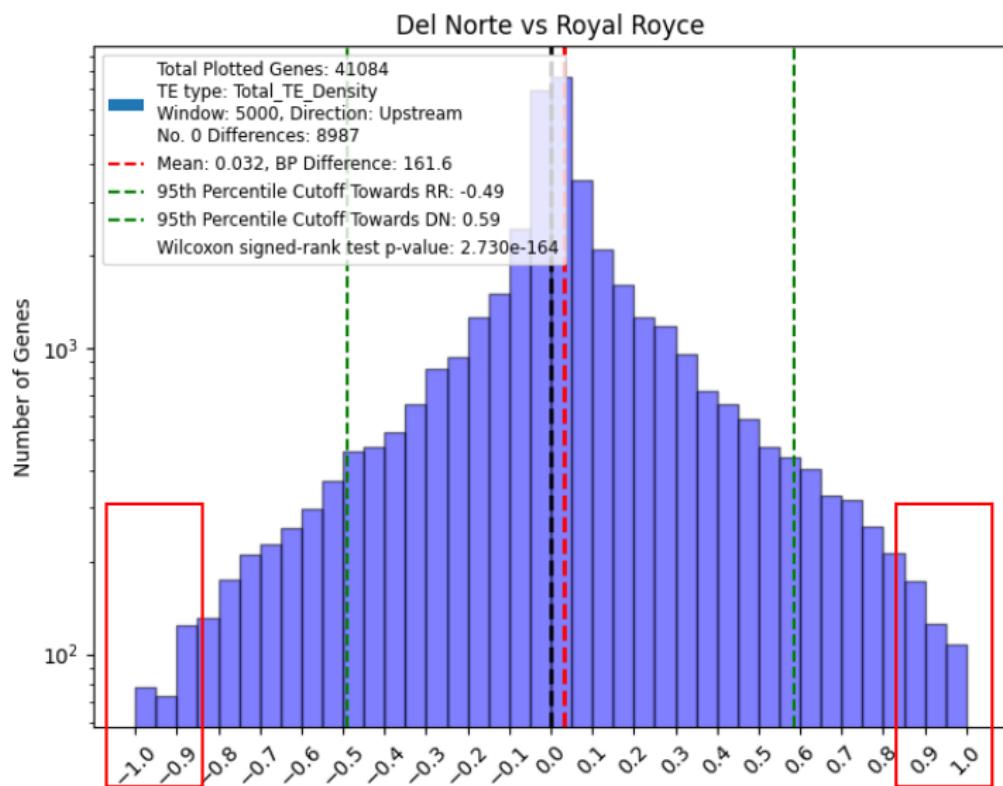
- For every TE type
- For every measurement window
- For every direction – upstream/downstream

Syntelogs show major differences in TE density

TE Category	Window	BP Bias	P-Value
Total TE	2500	113.7	1.603e-228
	5000	161.6	2.730e-164
	7500	202.3	4.871e-137
LTR	2500	125.6	5.776e-141
	5000	153.9	2.973e-99
	7500	171.9	7.172e-79
TIR	2500	73.2	4.140e-51
	5000	99.7	1.153e-41
	7500	126.7	5.753e-49
TIR/CACTA	2500	148.7	3.844e-33
	5000	218.0	1.813e-41
	7500	250.0	1.368e-48

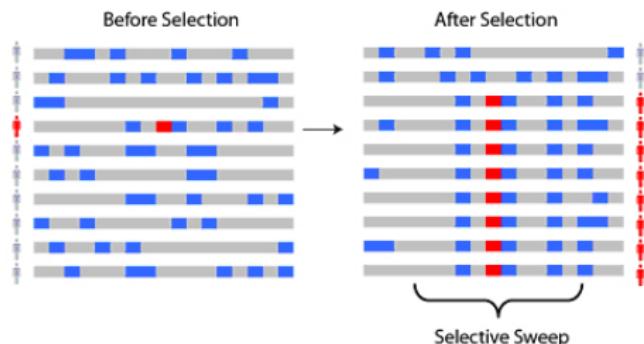
Table: Summary table of TE density differences between Del Norte and Royal Royce syntelogs. The BP Bias column reflects the average base-pair bias towards Del Norte, and the P-Value column reflects the significance of the Wilcoxon signed-rank test.

What kinds of genes differ in their TE presence between wild and cultivated strawberry?



Unique TE insertions in selective sweeps in cultivated strawberry

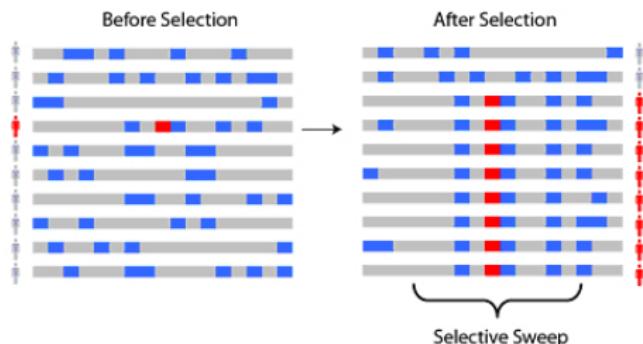
- The process of domestication entails the selection and preferential propagation of specific phenotypes, which leaves signatures of selection in the genome¹.



¹Hardigan et al. (2021) - Molecular Biology and Evolution

Unique TE insertions in selective sweeps in cultivated strawberry

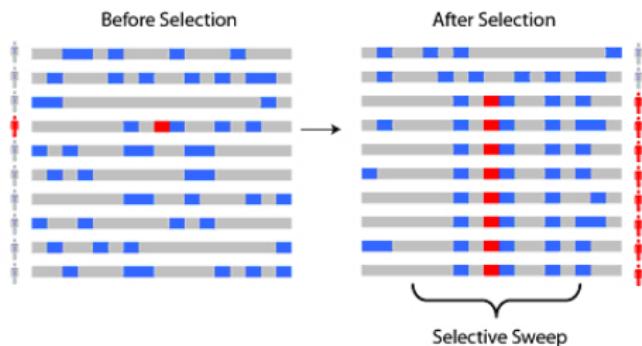
- The process of domestication entails the selection and preferential propagation of specific phenotypes, which leaves signatures of selection in the genome¹.
 - Selective sweeps are created when selection decreases genetic variation among neutral loci that are linked to the selected locus.



¹Hardigan et al. (2021) - Molecular Biology and Evolution

Unique TE insertions in selective sweeps in cultivated strawberry

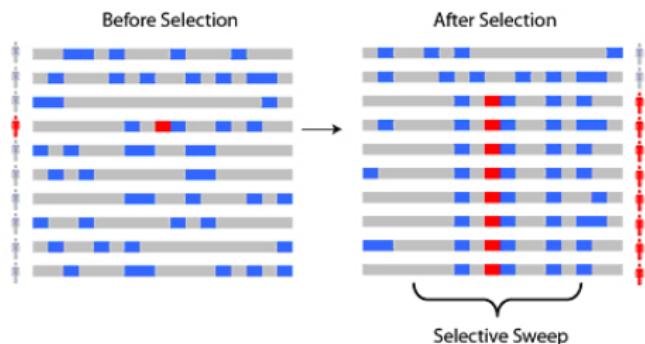
- The process of domestication entails the selection and preferential propagation of specific phenotypes, which leaves signatures of selection in the genome¹.
- Selective sweeps are created when selection decreases genetic variation among neutral loci that are linked to the selected locus.
- Selective sweeps are regions of the genome where a beneficial allele has rapidly increased in frequency, and the region has reduced genetic diversity.



¹Hardigan et al. (2021) - Molecular Biology and Evolution

Unique TE insertions in selective sweeps in cultivated strawberry

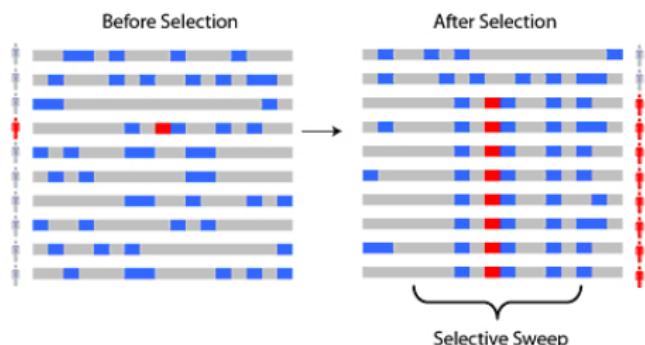
- I used data from previous publications¹ to intersect my differing TE-dense genes with regions that we know are associated with early and modern domestication.



¹Fan and Whitaker (2024) - Plant Cell

Unique TE insertions in selective sweeps in cultivated strawberry

- I used data from previous publications¹ to intersect my differing TE-dense genes with regions that we know are associated with early and modern domestication.
- I identified 60 genes residing within selective sweep regions that had significant functional enrichments and a TE density difference greater than or equal to 0.75.



What notable genes did I find?

- Fructose 6 phosphate gene — sugar metabolism
- A SWEET gene — sugar transport
- COP1 Supressor 2 — light signaling and development

What notable genes did I find?

- Fructose 6 phosphate gene — sugar metabolism
- A SWEET gene — sugar transport
- COP1 Supressor 2 — light signaling and development
- **Why do I think TEs are involved?**

What notable genes did I find?

- Fructose 6 phosphate gene — sugar metabolism
- A SWEET gene — sugar transport
- COP1 Supressor 2 — light signaling and development
- **Why do I think TEs are involved?**
 - They are all in selective sweep regions, early and/or modern domestication zones

What notable genes did I find?

- Fructose 6 phosphate gene — sugar metabolism
- A SWEET gene — sugar transport
- COP1 Supressor 2 — light signaling and development
- **Why do I think TEs are involved?**
 - They are all in selective sweep regions, early and/or modern domestication zones
 - There is no analogous TE in the wild strawberry genome, TE is unique to cultivated strawberry

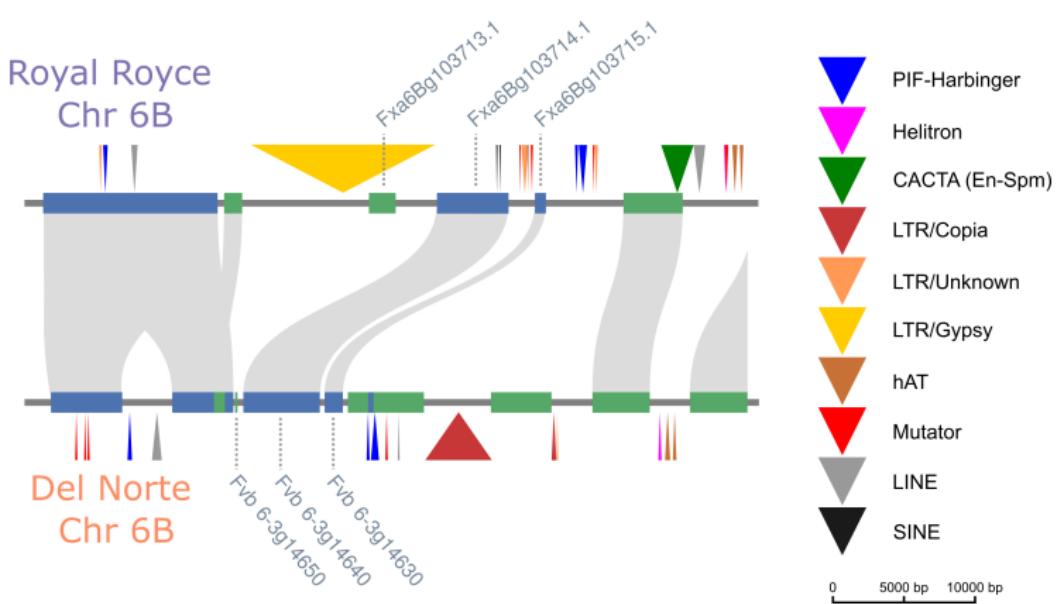
What notable genes did I find?

- Fructose 6 phosphate gene — sugar metabolism
- A SWEET gene — sugar transport
- COP1 Supressor 2 — light signaling and development
- **Why do I think TEs are involved?**
 - They are all in selective sweep regions, early and/or modern domestication zones
 - There is no analogous TE in the wild strawberry genome, TE is unique to cultivated strawberry
 - The TE is immediately upstream of the gene, and is not decayed at all

What notable genes did I find?

- Fructose 6 phosphate gene — sugar metabolism
- A SWEET gene — sugar transport
- COP1 Supressor 2 — light signaling and development
- **Why do I think TEs are involved?**
 - They are all in selective sweep regions, early and/or modern domestication zones
 - There is no analogous TE in the wild strawberry genome, TE is unique to cultivated strawberry
 - The TE is immediately upstream of the gene, and is not decayed at all
 - The TE is in a family with relatively low copy number, suggesting it is rather new

Fructose



Summary

- I applied TE Density software to strawberries, and found that cultivated strawberry has overall less TEs, and that it has less TEs near genes for every TE type and measurement parameter.

Summary

- I applied TE Density software to strawberries, and found that cultivated strawberry has overall less TEs, and that it has less TEs near genes for every TE type and measurement parameter.
- Functional enrichment analyses point to distinct patterns of TE-associated genes between the two genomes, particularly in pathways relevant to strawberry breeding.

Summary

- I applied TE Density software to strawberries, and found that cultivated strawberry has overall less TEs, and that it has less TEs near genes for every TE type and measurement parameter.
- Functional enrichment analyses point to distinct patterns of TE-associated genes between the two genomes, particularly in pathways relevant to strawberry breeding.
- I identified a number of fruit, defense, transcription factor, and development genes enriched for novel TE presence in domesticated strawberry, demonstrating the potential for TEs to influence a variety of important domestication related traits.

Summary

- I applied TE Density software to strawberries, and found that cultivated strawberry has overall less TEs, and that it has less TEs near genes for every TE type and measurement parameter.
- Functional enrichment analyses point to distinct patterns of TE-associated genes between the two genomes, particularly in pathways relevant to strawberry breeding.
- I identified a number of fruit, defense, transcription factor, and development genes enriched for novel TE presence in domesticated strawberry, demonstrating the potential for TEs to influence a variety of important domestication related traits.
- **The TE Density software offers a powerful method to identify genes that are potentially impacted by TEs**

Summary

- I applied TE Density software to strawberries, and found that cultivated strawberry has overall less TEs, and that it has less TEs near genes for every TE type and measurement parameter.
- Functional enrichment analyses point to distinct patterns of TE-associated genes between the two genomes, particularly in pathways relevant to strawberry breeding.
- I identified a number of fruit, defense, transcription factor, and development genes enriched for novel TE presence in domesticated strawberry, demonstrating the potential for TEs to influence a variety of important domestication related traits.
- The TE Density software offers a powerful method to identify genes that are potentially impacted by TEs
- Most of the time, when someone identifies a TE-influenced gene, it is because they first identified an unusual phenotype, *and then* they identified the gene, *and then* saw a TE was sitting right next to it. My methodology can be used to start with the genotype and work forwards.

Future Directions

- We know that TEs can influence gene expression, how many of the genes are being up vs. downregulated by a TE?

Future Directions

- We know that TEs can influence gene expression, how many of the genes are being up vs. downregulated by a TE?
- If a TE is altering the expression of something like a hormone gene, that can have far reaching effects. How can we analyze the downstream effects?

Future Directions

- We know that TEs can influence gene expression, how many of the genes are being up vs. downregulated by a TE?
- If a TE is altering the expression of something like a hormone gene, that can have far reaching effects. How can we analyze the downstream effects?
- What are the tissue-specific or temporal expression contexts of my candidate genes?

Future Directions

- We know that TEs can influence gene expression, how many of the genes are being up vs. downregulated by a TE?
- If a TE is altering the expression of something like a hormone gene, that can have far reaching effects. How can we analyze the downstream effects?
- What are the tissue-specific or temporal expression contexts of my candidate genes?
- Why do I have a fructose gene *inside* a TE (an LTR element no less!)? Evidence suggests it is a real gene, TE seems real too. Where did they both come from?

Future Directions

- We know that TEs can influence gene expression, how many of the genes are being up vs. downregulated by a TE?
- If a TE is altering the expression of something like a hormone gene, that can have far reaching effects. How can we analyze the downstream effects?
- What are the tissue-specific or temporal expression contexts of my candidate genes?
- Why do I have a fructose gene *inside* a TE (an LTR element no less!)? Evidence suggests it is a real gene, TE seems real too. Where did they both come from?
- What are the subgenome-specific patterns of TE insertions, what does the expression and density of the other homeologs look like?

Collaborations

Collaborations

I am teaming up with the Batman and the Antwoman

Collaborations

David Ray — The Batman



- Dr. David Ray of Texas Tech
- Bats have unusually long lifespans, given their body size (Longest living wild bat was 41 years old).
- Comparative analysis of bats and rodents' genomes suggests a relation between non-LTR retrotransposons, cancer incidence, and ageing¹

¹Ricci et al. (2023) - Scientific Reports

Collaborations

Janina Rinke — The Antwoman



- The Global Ant Genomics Alliance!
- Janina Rinke and Dr. Lukas Schrader of the University of Münster
- Ants are social insects and have evolved complex mechanisms for communication and chemical sensing

Collaborations

Janina Rinke — The Antwoman

GAGA_TEs_Analysis

main / +

History Find file Edit Code

Project information

66 Commits 4 Branches 0 Tags

README Add Wiki

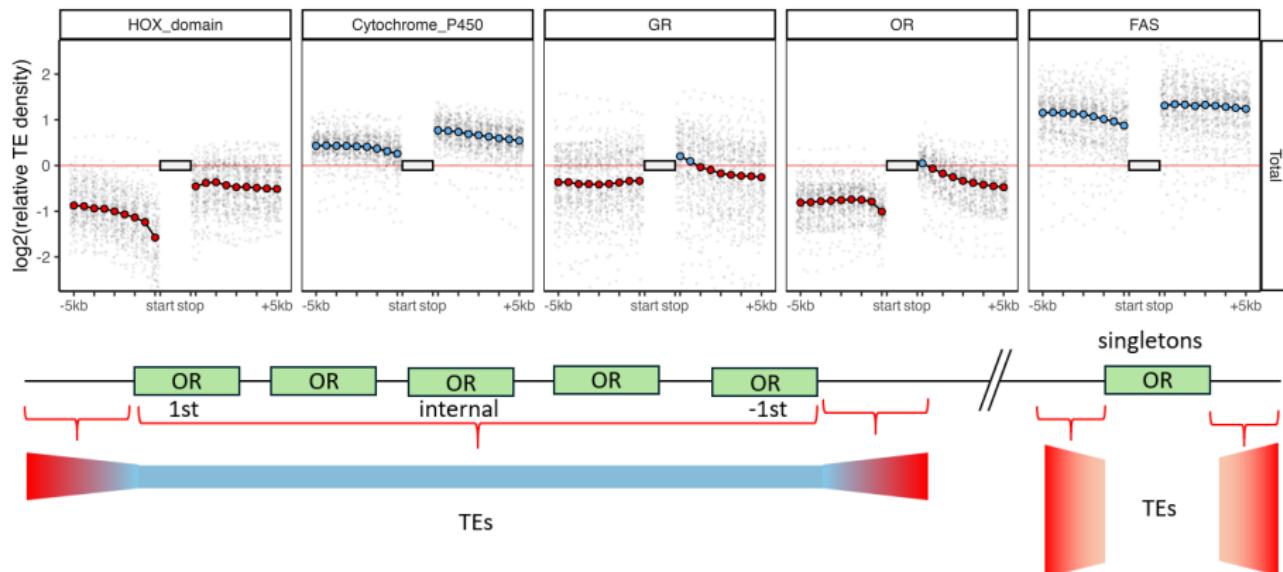
Created on September 12, 2023

Name	Last commit	Last update
.gitlab/merge_request_temp...	add simple merge request template	4 months ago
1_code	Update TE Density submodule to handl...	1 month ago
TE_Density_Submodule @ d064e0a0	Update TE Density submodule to handl...	1 month ago
.bash_history	changed working trees	5 months ago
.gitignore	changed working trees	5 months ago
.gitmodules	Add TE Density as submodule	5 months ago
Makefile	Create dotplots, address error handling...	1 month ago
README.md	changed working trees	5 months ago

Collaborations

Janina Rinke — The Antwoman

Preliminary results show depletion/enrichment of TEs near certain gene families, as well as an association between TE abundance and gene family expansions/contractions.



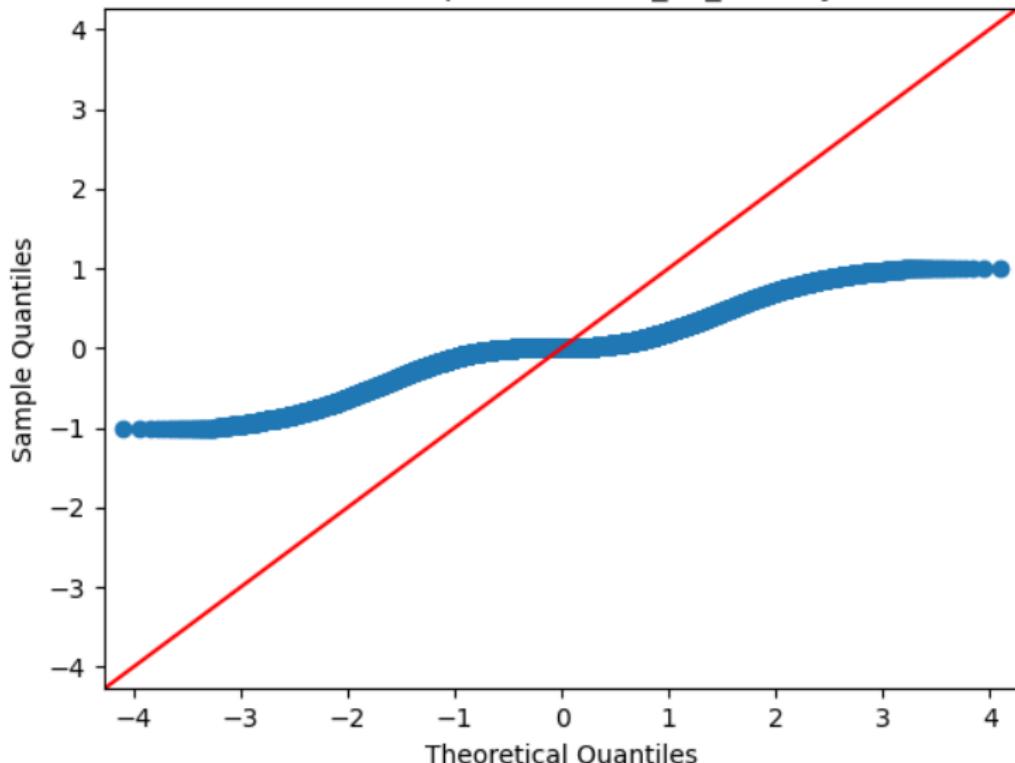
Acknowledgements

- Jay & Pam Teresi, Michael, Steven, Jane
- Erin Kramer
- Pat Edger & the Edger Lab
- Ning Jiang, Jiming Jiang, Jianrong Wang, Shin-Han Shiu
- Scott Stelpflug, Sarah Jensen, Harley Durbin-Rowan
- Shujun Ou
- Ryan Williams
- Jyothi Kumar
- Meghan Hill
- Alder Fulton
- Plant Genomics REU
- Ross Hatlen, Mallory St. Clair, Jordan Brock, Michael Gasdick
- James Bowers, Brandon Hinton, Drew Reibel
- Dan Wyrembelski
- Nicholas Panchy
- Joshua Puzey

Questions?



Del Norte vs Royal Royce
5000 BP Upstream Total_TE_Density



Extras

