Structural variation (A story of surprise \rightarrow frustration \rightarrow hope)





Alexander Sang-Jae Suh 서상재 徐商在

Senior lecturer University of East Anglia, UK Uppsala University, Sweden



@alexander_sul

What is structural variation?

Structural variant (SV): genomic

variation between individuals affecting the presence, abundance, position, and/ or direction of a nucleotide sequence

Mérot et al. 2020, Trends Genet.

Some key mutation types



Berdan et al. 2021, Mol. Ecol.

Part 1: Surprise



A) It's not a SNP!

Delicious effects of SVs



Discovery of gene regulation in 1940s



Barbara McClintock (Nobel Prize in Physiology or Medicine 1983)

Transposon Day: Barbara McClintock's 120th birthday is this week (June 16), check out: https://twitter.com/search?q=%23TransposonDay2022

TE-induced rapid adaptation

The industrial melanism of the peppered moth is probably the most famous textbook example for adaptation (in only a few decades)!



Van't Hof et al. 2016, Nature





High diversity of possible effects



Part 1: Surprise



B) Covariation

Two key mechanisms of structural change

Non-homologous end joining (**NHEJ**) (requires double-strand DNA breaks)



Non-allelic homologous recombination (NAHR) (requires sequence homology)

NHEJ correlates with frequency of DNA damage, NAHR correlates with frequency of (identical, large) repeats

Genome shrinking despite more TEs





Accordion model



Consider not only host popgen, but also TE popgen!

Genome size and life history traits



Dynamic genome (more TEs, fast shrinking)



Static genome (fewer TEs, slow shrinking)

Adaptive processes are often invoked but remain difficult to prove (few high-quality genome assemblies and lack of popgen data)!



Non-adaptive processes likely contribute to a large or very large degree!

Genomes: whack-a-transposon



More context in Suh 2021 TE lecture 6

Covariation between (epi)mutation types



Gene upregulation

Spillover of DNA methylation and/or histone modifications from new TE insertions to nearby genes!

Weissensteiner & Suh 2019 in Avian Genomics book

Host–TE conflict and reproductive isolation



Weissensteiner & Suh 2019 in Avian Genomics book



If an inversion or duplication leads to gene truncations, a toxin/antitoxin system can evolve to distort its transmission!

Centromere drive of some SVs



If a pericentric inversion or a centromere shift leads to a stronger centromere, it can distort its own transmission!

Kursel & Malik 2018, Curr. Opin. Cell Biol.

Questions?



Part 2: Frustration



A) Concepts and methods

What this lecture will not cover

- 1. Genome assembly: What is (not) assembled? Primers: Peona et al. 2018, Peona et al. 2021, Rhie et al. 2021, Nurk et al. 2022
- 2. Gene and repeat annotation: What is (not) annotated? Primers: Yandell & Ence 2012, Suh 2021 TE lecture 4, Goubert et al. 2022
- 3. Within-individual or germline/soma genome differences Primers: <u>Smith et al. 2021</u>, <u>Suh & Dion-Côté 2021</u>, <u>Borodin et al. 2022</u>
- 4. All SVs, all processes, all effects, all methods, all limitations. Talk to Valentina and me until 10 pm today!



Valentina Peona

Awareness of biology and technology



How can we make sure that what we see in our data is what we think it is?

Did we account for biological patterns/processes and technological limitations?

Terminology

Synteny vs. collinearity

Hs2 Pt12/13



Dot plot



Pattern vs. process





Feuk et al. 2005, PLoS Genet.

Beware of waves

My SNP explains everything! My inversion explains everything! My TE explains everything!



Each of these statements can be true, but what if there is covariation with other mutation types?

Taxon X is not known to have mutation type Y

We did not look for mutation type Y in taxon X

Reflection on biases

Confirmation bias

Survivorship bias





My own biases: I like transposable elements, centromere shifts, and simple answers to complicated questions!

Ultimate vs. proximate causes

Proximate: This TE is beneficial for the host

<u>Ultimate</u>: TEs jump to be beneficial for the host TEs jump because they can

<u>Proximate</u>: This asteroid caused diversification

<u>Ultimate</u>: Asteroids land to cause diversification Asteroids land eventually



What is the null hypothis?

Guilty until proven innocent Innocent until proven guilty

> Absence of evidence Evidence of absence



Theory applies to SNPs and to SVs



Background variation: What SNPs and SVs are there?

http://evomics.org/2022-workshop-on-population-and-speciation-genomics-cesky-krumlov/

SVs are nowhere as established as SNPs



Problem: Reliable SV genotyping (cf. SNP activies in this workshop) + accounting for covariation with other SVs (cf. this lecture) is essential but the SV field is not there yet.

One approach to find them all?



Trends in Ecology & Evolution

Mérot et al. 2020, Trends Genet.

How to pick a tool for finding SVs?

Repeat tools

Description

This page compiles a list of software for the detection, annotation, analysis, simulation and visualization of repetitive, mobile and selfish DNA and related entities.

It is maintained by <u>Tyler A. Elliott</u> \square and a more metadata rich form of the data can be found <u>here</u> \square . It was initiated with the help of Elizabeth Smikle and Miduna Rahulan, formerly and currently at the <u>Centre for Biodiversity Genomics</u> \square at the <u>University of Guelph</u> \square . Suggestions, updates and error corrections are welcome. Please feel free to add missing tools into the table, that would help a lot!

We encourage the authors of these tools to create pages for them on TE Hub, so that they can provide more information about their work, and link it back to this table. Please find a <u>template software sheet here</u>.

HUB

Overview of tools for repeat analysis

Tool≎Find	DOI‡Find	Alternate URL [↑] Find	Keywords I Polymorphism
AluMine 🖸	<u>https://doi.org/10.1101</u> /588434 ⊠		Alu, SINE, Genotype, Polymorphism, NGS/HTS
alu-detect	<u>https://doi.org/10.1093</u> /nar/gkt612 ☑		Alu, SINE, Genotype, Polymorphism, NGS/HTS, Paired- End

88 tools listed for TE insertion polymorphism analysis!

https://tehub.org/en/resources/repeat_tools; The TE Hub Consortium 2021, Mobile DNA

Read-based SV detection



Reliable read mapping and SV scoring is difficult near (other) repeats, near gaps, at misassemblies ...



It could all be so easy (if it wasn't for technological limitations) 11 ((7) ((K - ?) Х - <u>|</u>[_ Н 51 88 1




SV mapping with longer and longer reads













Sedlazeck et al. 2018 Nat. Rev. Genet.



Tandem duplications are (usually) collapsed in assemblies!

Sedlazeck et al. 2018 Nat. Rev. Genet.

Not all gaps are equal

Chromosome 18 of hooded/carrion crow





~ 3.5 Mb

Rule of thumb: centromeres are not *in* assemblies but in gaps within or between scaffolds!

Miga 2015, Chromosome Res.

Questions?



Coffee break (10 minutes)



<u>Task</u>: Form random groups of 3 and discuss what resources (assembly quality, gene/repeat annotation) there are for your respective study system.

Part 2: Frustration



B) Biology and more concepts

Transposable elements are very diverse



	Superfamily				
Class I (reti	rotransposons)				
LTR	Copia	GAG AP INT RT RH	4-6	RLC	P, M, F, O
LIK	Gypsy		4-6	RLG	P, M, F, O
	Bel-Pao	GAG AP RT RH INT	4-6	RLB	M
	Retrovirus		4-6	RLR	M
	ERV		4-6	RLE	м
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	0
PLE	Penelope	RT EN	Variable	RPP	P, M, F, O
LINE	R2	RT EN	Variable	RIR	м
	RTE	APE RT	Variable	RIT	м
	Jockey	ORFI APE RT	Variable	RIJ	м
	L1	ORFI APE RT	Variable	RIL	P, M, F, O
	1	ORFI APE RT RH	Variable	RII	P, M, F
SINE	tRNA		Variable	RST	P, M, F
	7SL		Variable	RSL	P, M, F
	55	<u></u>	Variable	RSS	M, O
Class II (DN	A transposons) - Su	bclass 1			
TIR	Tc1-Mariner	Tase*	TA	DTT	P, M, F, O
	hAT	Tase*	8	DTA	P, M, F, O
	Mutator	Tase*	9-11	DTM	P, M, F, O
	Merlin	Tase*	8-9	DTE	M, O
	Transib	Tase*	5	DTR	M, F
	Ρ	Tase Tase	8	DTP	P, M
	PiggyBac	Tase Tase	TTAA	DTB	M, O
	PIF- Harbinger	Tase* ORF2	3	DTH	P, M, F, O
	CACTA	► ↔ ← Tase – ORF2 → ↔ ← <	2-3	DTC	P, M, F
Crypton	Crypton	YR	0	DYC	F
Class II (DN	A transposons) - Su	bclass 2			
Helitron	Helitron	RPA Y2 HEL	0	DHH	P, M, F
Maverick	Maverick				

<u>Today's</u> <u>focus</u>: LINE, SINE, LTR, TIR

Weirder TEs in Suh 2021 <u>TE</u> <u>lecture 1</u>

Class I: LINE retrotransposons

Classification		Structure	TSD	Code	Occurrence	
Order	Superfamily					
Class I (retr	rotransposons)					
PLE	Penelope		Variable	RPP	P, M, F, O	
LINE	R2	RT EN	Variable	RIR	М	
	RTE	APE RT -	Variable	RIT	М	
	Jockey	- ORFI - APE RT -	Variable	RIJ	М	
	L1	- ORFI - APE RT -	Variable	RIL	P, M, F, O	
	1	- ORFI - APE RT RH	Variable	RII	P, M, F	



Target-primed reverse transcription (TPRT)





TPRT frequently undergoes premature termination (5' truncation)

Target site duplication (TSD)



TSDs are a hallmark of nearly all (retro)transposition mechanisms!

Kazazian 2004, Science

Class I: SINE retrotransposons

Classification		Structure	TSD	Code	Occurrence	
Order	Superfamily					
Class I (retr	otransposons)					
SINE	tRNA	<u></u>	Variable	RST	P, M, F	
	7SL	<u></u>	Variable	RSL	P. M. F	
	5S	<u></u>	Variable	RSS	M,O	

SINEs are parasites of LINEs! Trans-mobilization via LINE enzymes.



Class I: LTR retrotransposons

Classification		Structu	Structure					TSD	Code	Occurrence	
Order	Superfamily										
Class I (ret	rotransposons)										
LTR	Copia		GAG	AP	INT	RT	RH	→	46	RLC	P, M, F, O
	Gypsy	→ <u></u>	GAG	AP	RT	RH	INT	→	4-6	RLG	P, M, F, O
	Bel-Pao		GAG	AP	RT	RH	INT	→	4-6	RLB	м
	Retrovirus		GAG	AP	RT	RH	INT		4-6	RLR	м
	ERV	->	GAG	AP	RT	RH	INT		4-6	RLE	м
DIRS	DIRS	>		AP					0	RYD	P, M, F, O
	Ngaro	\rightarrow	GAG						0	RYN	M, F
	VIPER	\rightarrow	GAG	AP	RT	RH	YR		0	RYV	0



Replicative retrotransposition





Why LTR retrotransposons have LTRs

	LTR	GAG	AP	RT	RH	INT	(ENV)		LTR		
a	•	PBS (primer bind	ling site	e)						>]
ь.		trna	primer	www.ww	mannahan	mhannaha	uman man	•	•		
с .		~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	www.	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	*****	******		•	~AAAA	
d .	······*	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	www.	~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~~	kannakann	the when	walker walker wat		·		
e .	ann mar mar an a	vakaannakaannakaannakaann	Ano mano ma	\A.D.~~~\$\A.D.~~	whan whan	wannana	mpunne	harmhar	•	· AAAA•	22
• • f.	· · · · · · · · · · · · · · · · · · ·							•	•	· · ·	22
g							~~~	•	•	· · ·	25
h .	· · · ·						~~~	•	•	· · ·	25
• •								•	•	· · ·	
	• • •										-

Kazazian 2004, Science

Class II: DNA transposons

Classification		Structure	TSD	Code	Occurrence	
Order	Superfamily					
Class II (DN	A transposons) - Subcla	ss 1				
TIR	Tc1–Mariner	Tase*	TA	DTT	P, M, F, O	
h N N	hAT	Tase*	8	DTA	P, M, F, O	
	Mutator	Tase*	9–11	DTM	P, M, F, O	
	Merlin	Tase*	8-9	DTE	M,O	
	Transib	Tase*	5	DTR	M, F	
	Р	Tase Tase	8	DTP	P, M	
	PiggyBac	Tase Tase	ΤΤΑΑ	DTB	M,O	
	PIF– Harbinger	Tase* ORF2	3	DTH	P, M, F, O	
	CACTA	Tase ORF2	2-3	DTC	P, M, F	
Crypton	Crypton	YR	0	DYC	F	



Cut-and-paste transposition (TIR)

a DNA transposon

'Cut and paste' TE







How to increase in copy number?



I. DNA replication fork passes transposon



II. Newly replicated transposon is cut out...



III. ...and inserted into a not-yet replicated genomic site



IIII. DNA replication fork passes insertion site





I. Newly replicated transposon is cut out...



II. ...and transposed into a new locus



III. Following transposition, the double-stranded break is repaired by homology-dependent DNA repair



TE ≠ TE SINE

LINE



More context in <u>Suh 2021 TE lecture 2</u>

Inversion formation



"We found that inversion breakpoints frequently occur in centromeric and telomeric regions and are often flanked by long inverted repeats (0.5-50 kb)"



Assembling or mapping inversion breakpoints is difficult!



Inversions reduce recombination (2)



Rare recombination in (large) inversions Independents Satellite Independent

Faeder





More cases of NAHR



Fusions/fissions/translocations can decrease (new proximity to centromere) or increase (new proximity to telomere) recombination rates

Duplications can increase the chance of further nonallelic homologous recombination (NAHR)

Centromere shifts



Chromosome 18 of hooded/carrion crow





Centromere shifts across songbirds

. 8510	8520	8530	8540	8550 8	117090	117100	117110	11712	20 11	17130 1
Corvus cornix scaleguege	gtttagtatgtgat	aaggetttggg	gaatggetet	gcaacaggacgg	cctgctccagete	cagacctct	agccactcg	ggg <mark>c</mark> g	gacaga	gggggactto
Corvus hawaiiensiigutgu	gtttagtatgtgat		gaatggetet	gcaacaggacgg		cagacctct	agecacteg	ggg <mark>g</mark> gg	gacaga	gaggactto
Corvus nawallensiagttgt	gtttagtatgtgat		gaatggetet	gcaacaggacgg		agacetete	agecacteg	ggggeeg	gacaga	gaggactt
Corvus honeduloid agetge	atttagtatgtgat		gaatggetet	gcaacaggatgg		cagacctct	agecacteg		gacaga	gggggactto
Aphelocoma coerul aggigt	gtttagtatgtgat	aggetttgg	gaatggetet	gcgacaggatgg		cagacctct	agecactea	aaacta	gacaga	cagggggactto
Lanius ludivician aggt gt	gtttggtatgtgat	agcetttgg	aaatggetet	getecaggatg-			a	aaacta	glacaga	cacccattt
Ifrita kowaldi scaggtgt	gttcagtatgtgat	aaa			.		a	aaacta	gacaga	caggaactt
Paradisaea raggiaatgtgt	gtttagtatgtgat	ag					q	gggctgt	gacaga	cagggactt
Struthidea cinerelaggtgt	gtttagtatgtggc	aag			·		ď	gggctgt	gacaga	cagggactt
Dicrurus megarhyn <mark>aggtgt</mark>	atttagtatgtgat	aag					<mark>g</mark>	ggg <mark>c</mark> tgt	gacaga	cagggac <mark>tt</mark> g
Gymnorhina tibiceaggtgt	gttt <mark>agta</mark> tgtgat	aag					<mark>g</mark>	ggtctc	gacaga	<mark>ca</mark> ggg <mark>actt</mark> g
Rhagologus leucos <mark>aggt</mark> gt	.gttt <mark>agca</mark> tgtgat	aag					g	ggg <mark>c</mark> tgt	gacaga	<mark>ca</mark> ggg <mark>actt</mark> o
Pachycephala philaggtgt	gtttagtgtgtgat	gag					<mark>g</mark>	ggg <mark>c</mark> tgt	gacaga	<mark>cagggactt</mark> o
Oreocharis arfakiaggegt	gtttagtatgtgat	aag					<mark>g</mark>	ggg <mark>c</mark> tg	gacaga	<mark>eggactt</mark> e
Falcunculus frontagetgt	gtttagtatgtgat	aag					g	ggg <mark>c</mark> g	gacaga	cagggactto
Machaerirhynchus aggogg	c tttagtatgtgat	aag					g	ddd <mark>c</mark> rd	gacaga	cagggactto
Eulacestoma nigrolaggogt	ggttagtatgtaat	aag					g	dddord	gacaga	cagggactto
Edolisoma coerule aggugu	gtttaggatgtgat	Jag					9	gggccg	gacaga	agggactto
Dyaphorophyla casago go	gtttagtgtgtgtgat						9	ggggggggg	g <mark>cac</mark> aga	caggggactt
Mehana achroganha agttat	atttagtatgtgat	1ag 1ag					9 g		gacaga	cagggactt
Mujagra bebetior aggigt	gtttagtatgtgtgat	aag					d	aaacta	gacaga	caggggactt
Chastorhunchus na aggigt	gtttagtacgtgat	aag					d	aaacta	gacaga	caggggactt
Ptilorrhoa leucostaggtgt	gtgtagtatgtgat	aad			.		d	taacta	gacaga	cagggggt
Aleadryas rufinucaggtgt	gtttggcatgtgat	aag			.		a	aaacta	gacaga	cagggactt
Drvoscopus gamben aggtgg	ggttagtgtgtgat	aag					q	agactat	gacaga	cagggactt
Erpornis zantholeraggtgt	gtttagtgtgtgat	aag								<mark>t</mark> o
Vireo altiloquus aggtgt	gtttagtttgtgat	aag			·		<mark>g</mark>	ggg <mark>c</mark> tgt	gacaga	cagggactt
Pteruthius melanoaggtgt	gtttagtatgtgat	aag					<mark>g</mark>	ggg <mark>c</mark> tgt	gacacac	<mark>cagggactt</mark> g
Daphoenositta chr <mark>aggtgt</mark>	g ttcagcat gtgag	aag					g	ggg <mark>ct</mark> gt	gacaga	<mark>ca</mark> ggg <mark>actt</mark> a
Mystacornis crossagatgt	gttt <mark>aatac</mark> gtgat	aag					g	ggg <mark>cc</mark> gi	gacaga	<mark>cagggac<mark>tt</mark>g</mark>
Rhipidura dahli s <mark>agg</mark> tgt	gtttagtgtgtgat	aa <mark>c</mark>					<mark>g</mark>	ggg <mark>c</mark> tgt	gacaga	caggga <mark>c-t</mark> o
Taeniopygia gutta <mark>aggt</mark> gt	gtttagtatgtgat	gag					<mark>g</mark>	ggg <mark>c</mark> tg	gacaca	<mark>ca</mark> ggg <mark>attt</mark> g
Ficedula albicollaggatt	gtttagtatgtgat	gag					<mark>g</mark>	ggg <mark>c</mark> tg	gacaca	gagggatgt
Malurus_cyaneus_caggtot	gtttagtgtgggat	<u></u>					g	gggttgi	gacaga	cagggactt

>1 Mb satellite DNA array inserted in a formerly 5-kb intergenic region!



Not so stable chromosomes after all?



Am I stuck in confirmation bias?

Centromere shifts: frequency independent recombination reduction (unlike inversions)

Have centromere shifts been proposed as alternatives to inversions in speciation literature?

Are some variants that were interpreted as inversions actually centromere shifts?

Can centromere shifts fix more frequently than inversions because of meiotic drive?

Talk to me or email me (<u>alexander.suh@ebc.uu.se</u>) if you have references and/or criticism! PS: Manuscript in preparation about this model.

Coffee break (20 minutes)



<u>Task</u>: Gather in same groups of 3 and discuss 1) what SVs you <u>want</u> to study, 2) what SVs you <u>can</u> study, and 3) what data you need to be <u>less</u> frustrated.

Part 3: Hope



How frustrated are you?

- What types of SVs do you want to study?
- What types of SVs can you study?
- What data do you need to be less frustrated?



Genomes: ecosystems of selfish genes





Biodiversity inside each genome!

Cellular organisms

Phylum

Class



Family Genus

> Species Individual

Transposable elements

Class

Subclass Order

Superfamily

Family

Subfamily

Сору

More context in Suh 2021 <u>TE</u> lecture 3

Too much TE data, too few TEologists

Repetitive elements in the era of biodiversity genomics: insights from 600+ insect genomes

John S. Sproul, D Scott Hotaling, Jacqueline Heckenhauer, Ashlyn Powell, D Amanda M. Larracuente,
Joanna L. Kelley, Steffen U. Pauls, D Paul B. Frandsen
doi: https://doi.org/10.1101/2022.06.02.494618

Posted June 03, 2022.

Our

findings suggest this RE-annotation bottleneck, driven largely by uneven taxonomic representation in RE reference databases, is worsening. Although the diversity of available insect genomes has rapidly expanded, the rate of community contributions to RE databases (essential for RE annotation) has not kept pace, preventing high resolution study of REs in most groups. We highlight the tremendous opportunity and need for the field of biodiversity genomics to embrace REs and suggest collective steps for making progress towards this goal.

Complete human genome in April 2022

- Nearly 200 million bp more than the previous human reference (GRCh38) with 1956 new genes (99 proteincoding) and 0 assembly gaps!
- Homozygous cell line sequenced with: 120x coverage of Oxford Nanopore ultra-long reads, 70x PacBio CLR long reads, 30x PacBio HiFi long reads, 50x 10X Genomics linked reads, BioNano DLS optical maps, Arima Genomics Hi-C maps.





Money is less of a limitation now than sample amount + quality + repetitiveness!

Nurk et al. 2022, Science
What's next: Machine learning?

DeepTE: a computational method for de novo classification of transposons with convolutional neural network

https://github.com/LiLabAtVT/DeepTE, Yan et al. 2020 Bioinformatics

TransposonUltimate: software for transposon classification, annotation and detection

https://github.com/DerKevinRiehl/TransposonUltimate, Riehl et al. 2022 Nucl. Acids Res.



Prediction: AI training (cf. SV biology and curation) will be a key bottleneck for evaluating machine learning results!

More community initiatives needed



TE Worldwide Slack #te-hub channel

Let's keep in touch in the WPSG 2022 Slack channel #structural variation! (https://wpsg2022participants.slack.com)



Alexander Suh 4:38 PM Please click de for the SV types you are interested in:

Inversions et <mark>/</mark>1)



☺

Translocations

61)

☺

Insertions/deletions

1) G

Fusions/fissions

et 👍 1)

Centromere shifts

¢

¢ 1 Duplications

Conclusion: Genomics is no silver bullet



Time spent looking at genomics-y data

Genomics + cytogenetics = cytogenomics



What to take with a grain of salt

- 1. How can we declare something as absent in a genome (evidence of absence vs. absence of evidence)?
- 2. How can we study unassembled or underassembled regions (multicopy genes, GC-rich genes, TEs)?
- 3. How can we compare species with different assembly qualities, data types, or annotation efforts?
- 4. How can we account for unknown peculiarities (sex chromosomes, B chromosomes, ...)?



The transposition goes on!



UNIVERSITET



IB Leibniz Institute for the Analysis of Biodiversity Change

From 01 April 2023: Professor at University of Bonn and Head of the **Centre for Molecular Biodiversity Research** in Bonn/Hamburg, Germany

Talk to me or email me (<u>alexander.suh@ebc.uu.se</u>) if interested in developing cytogenomics and/or genome annotation/curation for scaling across animals

PhD/postdoc/researcher positions available!



Questions?



Thank you, Český Krumlov!

