2020 WORKSHOP ON POPULATION AND SPECIATION GENOMICS, CESKY KRUMLOV

The Multi-Species Coalescent (MSC) and its Application in Phylogenetics and Species Delimitation

L. Lacey Knowles

Dept. of Ecology and Evolutionary Biology University of Michigan

Software: Delineate Jeet Sukumaran Dept. of Biology, Evolutionary Biology Program San Diego State University

https://github.com/jeetsukumaran/delineate

Software: Decrypt

Arnaud Becheler Dept. of Ecology and Evolutionary Biology University of Michigan <u>https://becheler.github.io/pages/applications.html</u>

Delimitation models that bring speciation to the multispecies coalescent

Inference of species boundaries (beyond the MSC)

Software: *Delineate* <u>https://github.com/jeetsukumaran/delineate</u>

• Phylogenetic modeling approach that delineates species versus population lineages under a protracted speciation model

Software: *Decrypt* <u>https://becheler.github.io/pages/applications.html</u>

• Model of the geography of genetic divergence under a spatially explicit coalescent to evaluate competing hypotheses about cryptic diversity (inferred under the MSC)

• approach integrates an explicit model of speciation into the "censored" or "multispecies" coalescent model to organize a set of population lineages sampled from one or more species into mutually-exclusive and jointly-comprehensive subsets, where each subset of population lineages represents a distinct species



(Sukumaran, Holder, and Knowles, 2020)

Seemingly fractal nature of species diversity

With more geographic sampling, get more and more species inferred under the MSC



Seemingly fractal nature of species diversity



Data Requirements

DELINEATE requires two items of data:

- A population tree
- A species assignment table

File: Dynastes_calib_GeoPop.tre File: constraints1.txt

Input Data: Population Tree

This is rooted ultrametric tree where each tip lineage represents a population or deme. This tree is typically obtained through a classical "censored" or "multispecies" coalescent analysis, such as results from BP&P (either mode A01 or A10) or StarBeast. The tree can be be specified either in NEXUS or Newick format.

Data Requirements

- A population tree
- A species assignment table

Input Data: Species Assignment Table

This is a tab-delimited plain text file with at least three columns:

- "lineage"
- "species"
- "status"

lineage	spe	cies	status
DGRP1	gr	1	
DGRR1	gr	1	
Dhy6	hy	1	
• •	•	•	
• •		•	
• •		•	

There may be more than these three columns, but these columns are mandatory (and all other columns will be ignored). The order of columns does not matter, as the DELINEATE programs will use the labels specified in the header row (see below) to identify the columns.

The first row is the header row: i.e., column labels. Subsequent rows will map each tip in the population tree (the "lineage" column) to a species label ("species") as well as an indication of whether this species assignment is known or not ("status"). *Every* tip in the population tree must be represented by a row (and no more than one row) in this table. If species assignments are not known for some population lineages (as would be expected in species discovery type analyses), then the "species" field can be left blank or populated with an arbitrary value, such as a "?" or "NewSp.?" etc., but the "status" field should be set to "0". Species assignments that *are* known will have the status field set to "1".

To do a quick check of your files from your terminal:

(base) m-c02wt02ehh24:~ knowlesl\$ delineate-check

delineate-check -c Downloads/constraints1.txt -t Downloads/Dynastes_calib_GeoPop.tre -f newick

Given a population lineage tree file, "population-tree.nex", and a species assignment table file "species-mappings.tsv", then the following command will run a DELINEATE analysis on the data::

```
delineate-estimate partitions --tree-file population-tree.nex --config-file data1.tsv"
```

or, using the short-form options:

delineate-estimate partitions -t population-tree.nex --t data1.tsv"

This command has the following components:

- delineate-estimate : This is the name of the program to be run.
- estimate : This is the command or operation that the program will be running.
- --tree-file population-tree.nex or -t population-tree.nex : The --tree-file flag, or its short-form synonym, -t, specifies that the next element will be the path to tree file with data on the population lineage tree. In this example, the file is located in the current working directory, i.e., population-tree.nex. If it was in another directory, then the path could be /home/bilbo/projects /orc-species-delimitation/data1/population-tree.nex, for example.
- --config-file data1.tsv or -c data1.tsv The --config-file flag, or its short-form synonym, -c, specifies that the next element will be the path to species assignment configuration file. In this example, the file is located in the current working directory, i.e., *delineate-species.tsv*. Again, if it was in another directory, then the path could be */home/bilbo* /*projects/orc-species-delimitation/data1/data1.tsv*, for example.

(base) m-c02wt02ehh24:~ knowlesl\$ delineate-estimate partitions -c Downloads/constraints1.txt -t Downloads/Dynastes_calib_GeoPop.tre -f newick

Basic Run Output

Executing this command will run the DELINEATE analysis and will produce the following output files:

- "data1.delimitation-results.json" File: constraints1.delimitation-results.json
- "data1.delimitation-results.trees" File: constraints1.delimitation-results.trees

More generally, unless the <u>-o</u> or <u>--output-prefix</u> flag (see below) is used to explicitly specify an alternate output prefix for all results generated, the files will take on a prefix given by the file stemname of the configuration file, "*data1*" in the this example.

Basic Run Output

Executing this command will run the DELINEATE analysis and will produce the following output files:

- "data1.delimitation-results.json" File: constraints1.delimitation-results.json
- "data1.delimitation-results.trees"

The first file, with the general name of "*<output-prefix>.delimitation-results.json*", is the primary results file. As can be inferred from its extension, it is a JSON format text file, and it consists of a single a dictionary. The dictionary provides information on the estimated speciation completion rate as well as the probabilities of all the possible partitions of the population lineage leafset into species sets, given the species assignment constraints, ranked by the probability of each partition.

Basic Run Output

Executing this command will run the DELINEATE analysis and will produce the following output files:

- "data1.delimitation-results.json"
- "data1.delimitation-results.trees" File: constraints1.delimitation-results.trees

The first file, with the general name of "*<output-prefix>.delimitation-results.json*", is the primary results file. As can be inferred from its extension, it is a JSON format text file, and it consists of a single a dictionary. The dictionary provides information on the estimated speciation completion rate as well as the probabilities of all the possible partitions of the population lineage leafset into species sets, given the species assignment constraints, ranked by the probability of each partition.

Basic Run Output

Executing this command will run the DELINEATE analysis and will produce the following output files:

- "data1.delimitation-results.json"
- "data1.delimitation-results.trees" File: constraints1.delimitation-results.trees

The second file, with the general name of "*<output-prefix>.delimitation-results.trees*", provides supporting results. Basically this is a collection of trees, with one tree for each partition considered. The topology of the trees are identical, corresponding to the topology of the input tree (i.e., the population lineage tree), as are the tip labels. However, the tips have associated with them some extra metadata that will be available for viewing in a program like FigTree. Most important of this is species, i.e., the label corresponding to the identity of the species assignment in that partition. In the case of species assignments that are constrained (i.e., status indicated by "1"), these will be identical to the assignment and invariant across all partitions, of course. However, in the case of population lineages of *unknown* species affinities (i.e., status indicated by "0"), this may be an existing species label (if the population lineage was assigned to an existing species in the partition under consideration) or a new, arbitrary species label (if the population lineage was assigned to a new distinct species in the partition under consideration). In addition, in FigTree you can also choose to have the branches colored by "status", and this will highlight population lineages of (a priori) known vs unknown species affinities, and thus quickly identify the assigned species identities of the lineages of interest.

Basic Run Output

Executing this command will run the DELINEATE analysis and will produce the following output files:

- "data1.delimitation-results.json"
- "data1.delimitation-results.trees" File: constraints1.delimitation-results.trees

The second file, with the general name of "*<output-prefix>.delimitation-results.trees*", provides supporting results. Basically this is a collection of trees, with one tree for each partition considered. The topology of the trees are identical, corresponding to the topology of the input tree (i.e., the population lineage tree), as are the tip labels. However, the tips have associated with them some extra metadata that will be available for viewing in a program like FigTree. Most important of this is species, i.e., the label corresponding to the identity of the species assignment in that partition. In the case of species assignments that are constrained (i.e., status indicated by "1"), these will be identical to the assignment and invariant across all partitions, of course. However, in the case of population lineages of *unknown* species affinities (i.e., status indicated by "0"), this may be an existing species label (if the population lineage was assigned to an existing species in the partition under consideration) or a new, arbitrary species label (if the population lineage was assigned to a new distinct species in the partition under consideration). In addition, in FigTree you can also choose to have the branches colored by "status", and this will highlight population lineages of (a priori) known vs unknown species affinities, and thus quickly identify the assigned species identities of the lineages of interest.

Open in Text editor: constraints1.delimitation-results.trees

🕨 🗢 💿					Evaluation (25					
Currently Open Documents	~	🌣 ~/Dow	nlo	oads/co	nstraints1.delimitation-results.json 💲	(no function selected)	۰ ،	# v	•	# •
constraints1.delimitation-results.json	$\mathbf{\otimes}$	1 -		{						
🔥 constraints1.delimitation-results.json	\otimes	2		1	<pre>'speciation_completion_rate": 0.5298355815850312,</pre>					
constraints1.txt	×	3			<pre>'speciation_completion_rate_source": "estimated",</pre>					
DelineateCommands	\mathbf{x}	4			<pre>'speciation_completion_rate_estimate_lnl": -9.477</pre>	65239836891,				
DecryptCommanfd.txt		5			'num_partitions": 925,					
DelineateCommands.txt		6			<pre>'num_partitions_in_confidence_interval": 119, 'report mle enly": false</pre>					
lumper1.conf		8			'report_mice_only Talse. 'report_constrained_probability_threshold"null.					
Coreanabet Conf		9			report constrained cumulative probability thresh	old": null,				
Screenshot		10	•		'partitions": [
splitter1 conf	$11 \checkmark \{$									
		12 -	•		"species_leafsets" [
		13								
		14			"DGRP1", "DGPP1"					
		16	_							
		17	-		[
		18			"DhsH3",					
		19			"DhsG1",					
		20			"DhsP1"					
		21 •	-		1,					
		22 •							05 1	1.60
		23	_	10 m	Constraints Locimitation-results.json	anacias losfasta ^	Evaluation (25 days left)			
		25 -	-	∿r ~/D0		species_learsets V	* v		# •	
		26		45 46	"YSNE1"					
		27 🖬	-	47	⊾],					
		28	•	48						
		29		49 50	L 1.					
		30		51	• [
		32 -	-	52	"DtTN1"					
		33		53 54	►, ▼ [
		34		55	"DhlC01",					
		35		56	"ERVL2",					
		36		57	"Dh LP /"					
		37		59	▼ [
		38		60	"DhpB1Br"					
+ 0-		39		61						
		L:	i C	63	"constrained_probability": 0.154816662647	89598,				
				64	<pre>"constrained_cumulative_probability": 0.1</pre>	.5481666264789598,				

FigTree v1.4.4 - constraints1.delimitation-results.trees Node Clade Taxa 4 Q~ Selection Mode Prev/Next Collapse Reroot Rotate Annotate Colour Find Layout Д 3 -1 Zoom: Screenshot Fish Eye: ma Root Length: occ Curvature: sep Align Tip Labels sep sep Current Tree: 1 / 925 R Appearance P lic R Trees Time Scale R lic **Tip Labels** \$ mor |Display: species pas User selection 0 Colour by: ecu ecu Font Size: 8 🗘 ecu Setup: Colour Font tri Decimal Format: tri Sig. Digits: 4 🗘 tri **Tip Shapes** tri R First tree reflects constraints in configuration file Node Labels R tri **Node Shapes** R (yellow are unknowns) tri Node Bars \$ 0.7 **Branch Labels** R

Open in FigTree; File: constraints1.delimitation-results.trees

Open in FigTree; File: constraints1.delimitation-results.trees



Open in FigTree; File: constraints1.delimitation-results.trees



Original Alpha Taxonomy

lineage	spec	ies	status
DGRP1	gr	1	
DGRR1	gr	1	
Dhy6	hy	1	
Dhy3Br	hy	1	
DtTN1	ty	1	
Dhym5	mo	1	
Dma2	ma	1	
CCV01	000	1	
DhsP1	sep	1	
DhsH3	sep	1	
DhsG1	sep	1	
DhlC01	lic	1	
ERVL2	lic	1	
DhlP7	lic	1	
DhmB2Br	mor	0	
DhpB1Br	pas	0	
DheP8	ecu	0	
DheC06	ecu	0	
YSNE1	ecu	0	
DhtT9	tri	0	
DhbV2	blu	0	
DhrSL5	rei	0	
DhrM1	rei	0	
DHHG2	her	0	
DhhD1	her	0	

File: constraints1.txt

Data Requirements

- A population tree
- A species assignment table







D. h. reidi (rei)

D. h. hercules (her)

Hypotheses based on different species criteria:

Original Alpha Taxonomy

Data Requirements

- A population tree
- A species assignment table





Could test different hypotheses about species boundaries based on different species criteria

Huang & Knowles (2016) Syst. Biol.

Hypotheses based on different species criteria:

Original Alpha Taxonomy

"Lumper"Taxonomy

lineage	spec	cies	status	lineage	spe	cies st	atus	linea
DGRP1	gr	1		DGRP1	gr	grhy	1	DGRP1
DGRR1	gr	1		DGRR1	gr	grhy	1	DGRR1
Dhy6	hy	1		Dhy6	hy	grhy	1	Dhy6
Dhy3Br	hy	1		Dhy3Br	hy	grhy	1	Dhy3E
DtTN1	ty	1		DtTN1	ty	motyma	1	DtTN1
Dhym5	mo	1		Dhym5	mo	motyma	1	Dhym5
Dma2	ma	1		Dma2	ma	motyma	1	Dma2
CCV01	occ	1		CCV01	occ	sepocc	1	CCV01
DhsP1	sep	1		DhsP1	sep	sepocc	1	DhsP1
DhsH3	sep	1		DhsH3	sep	sepocc	1	DhsH3
DhsG1	sep	1		DhsG1	sep	seppcc	1	DhsG1
DhlC01	lic	1		DhlC01	lic	lic <mark>1</mark>		DhlCC
ERVL2	lic	1		ERVL2	lic	lic <mark>1</mark>		ERVL2
DhlP7	lic	1		DhlP7	lic	lic <mark>1</mark>		DhlP7
DhmB2Br	mor	0		DhmB2Br	mor	mor 0		DhmB2
DhpB1Br	pas	0		DhpB1Br	pas	pas 0		DhpB1
DheP8	ecu	0		DheP8	ecu	ecu 0		DheP8
DheC06	ecu	0		DheC06	ecu	ecu 0		DheCO
YSNE1	ecu	0		YSNE1	ecu	ecu 0		YSNE1
DhtT9	tri	0		DhtT9	tri	tri 0		DhtTS
DhbV2	blu	0		DhbV2	blu	blu 0		DhbV2
DhrSL5	rei	0		DhrSL5	rei	rei 0		DhrSL
DhrM1	rei	0		DhrM1	rei	rei 0		DhrM1
DHHG2	her	0		DHHG2	her	her 0		DHHG2
DhhD1	her	0		DhhD1	her	her 0		DhhD1

Data Requirements

- A population tree
- A species assignment table

"Splitter" Taxonomy

lineage	spec	ies	status
DGRP1	gr1	1	
DGRR1	gr2	1	
Dhy6	hy1	1	
Dhy3Br	hy2	1	
DtTN1	ty	1	
Dhym5	mo	1	
Dma2	ma	1	
CCV01	occ	1	
DhsP1	sep	1	
DhsH3	sep	1	
DhsG1	sep	1	
DhlC01	lic	1	
ERVL2	lic	1	
DhlP7	lic	1	
DhmB2Br	mor	0	
DhpB1Br	pas (0	
DheP8	ecu	0	
DheC06	ecu	0	
YSNE1	ecu	0	
DhtT9	tri (0	
DhbV2	blu (0	
DhrSL5	rei	0	
DhrM1	rei	0	
DHHG2	her	0	
DhhD1	her	0	

Data Requirements

- A population tree
- A species assignment table

lineage	spec	ies	status
DGRP1	gr	1	
DGRR1	gr	1	
Dhy6	hy	1	
Dhy3Br	hy	1	
DtTN1	ty	1	
Dhym5	mo	1	
Dma2	ma	1	
CCV01	000	1	
DhsP1	sep	1	
DhsH3	sep	1	
DhsG1	sep	1	
DhlC01	lic	1	
ERVL2	lic	1	
DhlP7	lic	1	
DhmB2Br	mor	0	
DhpB1Br	pas	0	
DheP8	ecu	0	
DheC06	ecu	0	
YSNE1	ecu	0	
DhtT9	tri	0	
DhbV2	blu	0	
DhrSL5	rei	0	
DhrM1	rei	0	
DHHG2	her	0	
DhhD1	her	0	

- Number of inferred species naturally increases as the known species assignments are made under more of a splitter, rather than a lumper, than perspective.

- However, naturally as well, the relationship is not a simple numerical one but is rather dependent on the branch lengths leading to the assigned species subtrees.

++		+	+	+	
Regime	Const	rained	Inferr	ed New	
++		+	+	+	
lumper2	3	1	0	I	
lumper1	4	1	0		
original	8	2	2		
splitter1	12	2	2		
splitter2	14	10	10		
++		+	+	+	

A new era of species delimitation models that brings speciation models to the multispecies coalescent

- Erroneous species boundaries are inferred from current model-based genetic approaches based on MSC Delimitation under the MSC:
 - genetic structure = species
- Relying on heuristics to interpret results from current genetic methods (e.g., bpp) is not the answer
- 5 cm (≈2")
- Future of genetic-based species delimitation is with speciation-based delimitation models in which species criteria is explicitly incorporated into test of species status of unknowns (e.g., "splitter" vs "lumper")

Incorporating the Speciation Process into Species Delimitation

Jeet Sukumaran^{a,c,1,2}, Mark T. Holder^{b,1}, and L. Lacey Knowles^a

*Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor MI, USA 41809-1079; Department of Ecology and Evolutionary Biology, University of Kansas, Lawrence KS, USA 66045; Department of Biology, San Diego State University, San Diego CA, USA 92182-4614

This manuscript was compiled on October 10, 2019

Processes contribute to the generation of this structure apart from speciation, including population isolation or divergence. This results in "population" units being conflated with "species" units when these approaches are applied without corroborating data. This confusion between populations and species has become more and more problematical recently, ironically due to the increasing resolution and scope of the data

Our work pioneers a new age of species delimitation approaches that deal with the problem directly by modeling the issue instead of ignoring it! That is, by actually incorporating an explicit model of the speciation process — in particular, an extended or protracted speciation process — into species delimitation, we are able to discriminate between species and population (or other) boundaries in genomic data.

knowlesl@umich.edu

QUESTIONS?



Jeet Sukumaran San Diego State Univ.



Mark Holder Univ. of Kansas



support NSF & the UM



Software: *Decrypt* <u>https://becheler.github.io/pages/applications.html</u>

• Model of the geography of genetic divergence under a spatially explicit coalescent to evaluate competing hypotheses about cryptic diversity (inferred under the MSC)