

spedeSTEM tutorial

Bryan Carstens Matthew Demarest Maxim Kim Tara Pelletier Jordan Satler

Acknowledgements



Development of spedeSTEM was funded via a grant from the National Science Foundation (DEB-0918212).

Thanks to the OSU College of Arts and Sciences Technology Services, particularly Sanford Shew and Timothy Smith.

Thanks to Dan Ence, Sarah Hird, John McVay and Noah Reid for discussions about spedeSTEM.

Thanks to all users who have analyzed their data using our program.





spedeSTEM: Introduction

spedeSTEM is a program that delimits species using maximum likelihood and information theory. Specifically, the probabilities of multiple permutations of putative evolutionary lineages are calculated using STEM 2.0 (Kubatko et al. 2009) and ranked by model probability (see Anderson 2004). **spedeSTEM** takes as input ultrametric gene trees from multiple loci and a user-supplied estimate of theta, and returns a table of models ranked by model probability. The web-based software here conducts both discovery and validation analyses, and also generates the set up files and allows the users to subsample alleles from large nexus files. **spedeSTEM** does not estimate gene trees; for this we suggest PAUP (Swofford 2002) or Garli (Zwickl 2006).

spedeSTEM also includes modules for subsampling, simulation testing. Currently these are limited to the python version, but we hope to implement them into the web-based application at some point in the future.

Species delimitation using species trees operates by comparing the probability of models where putative lineages are separate to the probability of models where they are the same. For example, consider one motivating example, the carnivorous plant *Sarracenia alata* that has a disjunct distribution in eastern and western regions of its range (Koopman & Carstens 2010; Zellmer *et al.* 2012; Carstens & Satler 2013). If we computed the probabilities of each of the models below, we might use a likelihood ratio test as suggested by Knowles and Carstens (2007).



Species delimitation using species trees: We could compute the probability of each of these models using STEM (Kubatko et al. 2009). Stem computes the maximum likelihood species tree using a coalescent model that accounts for the loss of ancestral polymorphism due to genetic drift. The probability is calculated using the gene tree density $f(g_j | S, \tau)$ given by Rannala & Yang (2003).

BIOINFORMATICS APPLICATIONS NOTE 2009, pages 1-3 doi:10.1093/bioinformatics/btp079

Phylogenetics

STEM: species tree estimation using maximum likelihood for

gene trees under coalescence

Laura S. Kubatko^{1,*}, Bryan C. Carstens² and L. Lacey Knowles³

¹Departments of Statistics and Evolution, Ecology, and Organismal Biology, The Ohio State University, Columbus, OH 43210, ²Department of Biological Sciences, Louisiana State University, Baton Rouge, LA 70803 and ³Ecology and Evolutionary Biology, University of Michigan, Ann Arbor, MI 48109, USA

Received on November 28, 2008; revised and accepted February 04, 2009

Associate Editor: Martin Bishop





If we compare our two *a priori* models using a Likelihood ratio test, we find that we can not reject the model where all of *S. alata* constitutes a single evolutionary lineage ($\chi^2 = 3.512$, p = 0.0609)



However, what if we consider a third model, as suggested by the optimal results from a Structure analysis? Here, we find that we can reject the null ($\chi^2 = 24.054$, p = < 0.0001)

Optimal partitioning of the data suggested by STRUCTURE results . . .

8/22/13



What is the appropriate null model for species delimitation? If we assume, geography and try to delimit 2 lineages, we do not reject the null. If we treat the results of the Structure analysis as the alternative hypothesis, then we do reject the null. Which of these is appropriate? It could be that neither is; there are dozens of possible alternative models, up to an including a model that treats all sampling localities as putative lineages.

This question is inherently related to your philosophy of statistics. We argue that null hypothesis testing is not particularly useful for a historical discipline such as phylogeography that lacks experimental replication or controls.

- spedeSTEM does not rely on hypothesis testing, and can consider any number of possible models (in theory). In practice, we have successfully considered >20000 for 10 putative lineages.
- Information theory is used to evaluate the probability of models (see Anderson 2008 for a readable introduction to information theory).
- All models are ranked using information theory, and the model probabilities are computed.
- STEM computes the probability of the model given the data in seconds, so thousands of models can be compared.



- SpedeSTEM 2.0 takes as input ultrametric gene trees from multiple loci.
- SpedeSTEM 2.0 includes validation and discovery approaches to species delimitation.
- SpedeSTEM 2.0 contains commands to assist the user in data formatting and allows for the subsampling of data sets.
- This tutorial covers both the **command line** and **web-based** versions of spedeSTEM 2.0. Both are written in Python, but the former has more functionality.

spedeSTEM: web-based



http://spedestem.asc.ohio-state.edu/

spedeSTEM: command line

• Download the python version of spedeSTEM 2.0 from the following URL.

http://carstenslab.org.ohio-state.edu/software.html

• It contains all example files used here.

Step 1: Clean your genetrees.

SpedeSTEM Online Run SpedeSTEM Resources Contac	t Us	~ ~ ~
Species delimitation using Maxin Likelihood	num	Sign up Login About us Besearch in the Carstens Jab
Clean Trees	Help	seeks to understand how biological diversity is generated
Generate Settings File	Help	using computational approaches.
Discovery	Help	by identifying the limits of
Validation	Help	evaluate the relative contributions of evolutionary processes and infer the ecological and environmental forces that have contributed to the formation of population genetic structure.
Department of Ecology Evolution and Organismal Biology		

8/22/13

Step 1: Clean your genetrees.

OOO SpedeSTEM Web Interface ×	R ^M
← → C ☆ _ spedestem.asc.ohio-state.edu/runspedestem/clea	ntrees Q 값 트
SpedeSTEM Online Run SpedeSTEM Resources Conta	ct Us
Setup: Clean Trees	Sign up Login
Email Scaling file Choose File No file chosen Tree file 1 Choose File No file chosen Add another treefile Run Back Department of Ecology Evolution and Organismal Biology	Format gene trees. An ultrametric gene tree is taken as input, and polytomies are randomly resolved by randomly adding a very very short branch. The scaling file contains a column of scalars in square brackets that standardize the substitution rate across loci (see next slide).
National Science Foundation WHERE DISCOVERIES BEGIN Funded by NSF DEB 0918212	

Scaling Factors File



Scaling file. A scaling file is required to run spedeSTEM – this is simply a text file with scaling factors for each locus. Loci should be in the same order as entered. Scaling factors are used to normalize the length of the gene trees. Let the scaling factor of each locus be a number that normalizes the (# segregating sites across loci / length of the locus). So, if you have two loci, locus A with 10 snps / 1000 bp and locus B with 10 snps / 500 bp, your scaling factor of locus A would be [2.0] (because 2*(10/1000) = (10/500)). It does not matter which of the loci are chosen as the locus with a scaling factor of [1.0] (we used locus B in this example). Place these in a simple text file as shown to the left.

Step 2: Generate settings file.

SpedeSTEM Online Run SpedeSTEM Resources	Contact Us	
Species delimitation using N Likelihood	laximum	Sign up Login About us Research in the Carstens lab
Clean Trees	Help	seeks to understand how biological diversity is generated
Generate Settings File	Help	using computational approaches
Discovery	Help	by identifying the limits of
Validation	Help	evaluate the relative contributions of evolutionary processes and infer the ecological and environmental forces that have contributed to the formation of population genetic structure.
Department of Ecology Evolution and Organismal B	liology	
National Science Foundation		

Setup 2: Generate a settings file.

OOO SpedeSTEM Web Interface ×	R _M
← → C ☆ C spedestem.asc.ohio-state.edu/runspedestem/generatesettings	s Q☆ Ξ
SpedeSTEM Online Run SpedeSTEM Resources Contact Us	
Satur: Ganarata Sattings Filo	Sign up Login
Traits file Imail Imail Back	Generate settings file . A settings file is required to run spedeSTEM, and the setup will generate this file. Two simple files are taken as input – a traits file (a tab delimited with columns equal to trails, species, groups, as in the *Beast and BPP traits file) and a value for $\theta = 4N_e\mu$. The later is input into the "BT value" field.
National Science Foundation WHERE DISCOVERIES BEGIN	

Cmd line setup

• navigate to spedeSTEM directory

cd [path to directory]

- access help by entering:
- ./SpedeSTEM_2.py setup -h
- generate settings file by entering:

```
./SpedeSTEM_2.py setup -c
cpDNA.tre Sa4.tre Sa135.tre
Sa163.tre Sa297.tre Sa302.tre
Sa323.tre Sa405.tre -s
scaling.txt -t traits.txt -bt
0.123
```

• generate genetrees file by entering:

cat cleaned.*.tre > genetrees.tre

SpedeSTEM-master — bash — 80×48 Last login: Mon Apr 22 10:54:31 on ttys000 dhcp-254-235-238:~ bcarstens\$ cd /Users/bcarstens/spedeSTEM/SpedeSTEM-master dhcp-254-235-238:SpedeSTEM-master bcarstens\$./SpedeSTEM 2.py setup -h usage: SpedeSTEM 2.py setup [-h] [-c treeFile [treeFile ...]] [-s scalingFile] [-t traitsFile] [-bt thetaValue] optional arguments: -h, --help show this help message and exit -c treeFile [treeFile ...], --clean treeFile [treeFile ...] prepare tree file(s) for analysis -s scalingFile, --scalingFile scalingFile provide a text file with each scaling factor on a new line in the format: [scalingFactor]. Values prefixed in the order that trees appear after the clean command, DEFAULT: [1.0] for each -t traitsFile, --traits traitsFile read in traits file as Beast format -bt thetaValue, --theta thetaValue set theta value for Beast formatted traits files, DEFAULT: 1.0 dhcp-254-235-238:SpedeSTEM-master bcarstens\$./SpedeSTEM 2.py setup -c cpDNA.tre Sa4.tre Sa135.tre Sa163.tre Sa297.tre Sa302.tre Sa323.tre Sa405.tre -s scaling. txt -t traits.txt -bt 0.123 Namespace(clean=['cpDNA.tre', 'Sa4.tre', 'Sa135.tre', 'Sa163.tre', 'Sa297.tre', 'Sa302.tre', 'Sa323.tre', 'Sa405.tre'], command='setup', scalingFile=['scaling.t xt'], theta=[0.123], traits=['traits.txt']) *********** ####### Performing Setup ####### ****** Cleaning Phylip tree... Parsing Beast traits file for Species, Traits, and Grouping dhcp-254-235-238:SpedeSTEM-master bcarstens\$ cat cleaned.*.tre > genetrees.tre dhcp-254-235-238:SpedeSTEM-master bcarstens\$ 🛽

Discovery: available on web based server

OOO SpedeSTEM Web Interface ×	R
← → C f Spedestem.asc.ohio-state.edu/runstem/discovery	☆ =
SpedeSTEM Online Run SpedeSTEM Resources Contact Us	
Sign up Log	jin
Discovery Abou	t us
Settings file Choose File Settings.disco Gene tree file Choose File genetrees.tre seeks to un	the Carstens lab
Email bryan.c.carstens@gmail.com biological dusing comp	iversity is generated utational approaches.
Save We investig by identifyi	ate empirical systems ng the limits of
evolutionary evaluate the	 lineages, in order to relative contributions
of evolution infer the eco	ary processes and ological and
environmen contributed	tal forces that have to the formation of
Deputation of the population o	genetic structure.
Department of Ecology Evolution and Organismal Biology	
• out (1).zip	Show All ×

Discovery Required files

Settings file. Discovery analyses require a settings file formatted in the setup step.

Genetrees file. A gene trees files is also required. This file includes a cleaned ultrametric tree for each locus preceded by a scaling factor. This file is the same as the one used in the validation step.



Discovery

- navigate to spedeSTEM directory
- cd [path to directory]
- access help by entering:
- ./SpedeSTEM_2.py discovery -h
- generate settings file by entering:

./SpedeSTEM_2.py discovery _t
genetrees.tre _s settings.disco

• Output: species trees and likelihoods in *results.txt* and information theoretic table in *itTable.txt*.

SpedeSTEM-master — bash — 80×48 Last login: Mon Apr 22 10:55:52 on ttys001 dhcp-254-235-238:~ bcarstens\$ cd /Users/bcarstens/spedeSTEM/SpedeSTEM-master dhcp-254-235-238:SpedeSTEM-master bcarstens\$./SpedeSTEM 2.py discovery -h usage: SpedeSTEM 2.py discovery [-h] -t treeFile -s settingsFile [-v] optional arguments: -h, --help show this help message and exit -t treeFile, --tree treeFile specify tree file -s settingsFile, --settings settingsFile specify settings file in STEM format -v, --verbose execute in verbose mode, DEFAULT: off dhcp-254-235-238:SpedeSTEM-master bcarstens\$./SpedeSTEM 2.py discovery -t genet rees.tre -s settings.disco Namespace(command='discovery', settings=['settings.disco'], tree=['genetrees.tre verbose=False) ****** ###### DISCOVERY ANALYSIS ###### ******* CAUTION ----stemOut.txt is about to be deleted. If you would like to preserve it, remove it from this directory. Are you ready to continue? (y/n): y ---- SETTINGS ------In Varification Mode: False Tree File: genetrees.tre.save Settings File: settings.disco Associations File: associations.txt Number of loci sampled each replicate: 8 Number of replicates: 1 In Verbose Mode: False ----- END SETTINGS ----------- EXECUTION -----Sampling 8 loci from master tree file genetrees.tre.save... Completed 1 of 1 replicates... Completing Analysis... See 'results.txt' and 'itTable.txt' files dhcp-254-235-238:SpedeSTEM-master bcarstens\$

Validation: available on web based server

OOO SpedeSTEM Web Interface ×		R _M
← → C ☆ D spedestem.asc.ohio-state.edu/runstem/validation	☆	≡
SpedeSTEM Online Run SpedeSTEM Resources Contact Us		
Subscription Subscription <td></td> <td></td>		
Out (2).zip I out (1).zip I Show the second seco	w All) ×

Validation Required files

Settings file. Discovery analyses require a settings file formatted in the setup step.

Genetrees file. A gene trees files is also required. This file includes a cleaned ultrametric tree for each locus preceded by a scaling factor – this file is the same as the one used in the Discovery step.



validation

navigate to spedeSTEM directory

cd [path to directory]

- access help by entering:
- ./SpedeSTEM 2.py validation -h
- generate settings file by entering:

./SpedeSTEM 2.py validation -t genetrees.tre -s settings.stem -a associations.stem

Output: species trees and • likelihoods in *results.txt* and information theoretic table in itTable.txt.

SpedeSTEM-master — bash — 80×52 dhcp-254-235-238:~ bcarstens\$ cd /Users/bcarstens/spedeSTEM/SpedeSTEM-master dhcp-254-235-238:SpedeSTEM-master bcarstens\$./SpedeSTEM 2.py validation -h usage: SpedeSTEM 2.py validation [-h] -t treeFile -s settingsFile -a associationsFile [-v] optional arguments: -h. --help show this help message and exit -t treeFile, --tree treeFile specify tree file -s settingsFile, --settings settingsFile specify settings file in STEM format -a associationsFile, --associations associationsFile specify associations file in STEM format -v, --verbose execute in verbose mode, DEFAULT: off dhcp-254-235-238:SpedeSTEM-master bcarstens\$./SpedeSTEM 2.py validation -t gene trees.tre -s settings.stem -a associations.stem Namespace(associations=['associations.stem'], command='validation', settings=['s ettings.stem'], tree=['genetrees.tre'], verbose=False) ********* ###### VALIDATION ANALYSIS ###### ***** ----- CAUTION -----<u>stemOut.txt is about</u> to be deleted. If you would like to preserve it, remove it from this directory. re you ready to continue? (y/n): y SETTINGS -----In Varification Mode: True Free File: genetrees.tre.save Settings File: settings.stem \ssociations File: associations.stem Number of loci sampled each replicate: 8 Number of replicates: 1 In Verbose Mode: False ----- END SETTINGS ----------- EXECUTION -----Sampling 8 loci from master tree file genetrees.tre.save... 25 permutations to run... Completed 10 of 25 permutations... Completed 20 of 25 permutations... Completed 1 of 1 replicates... Completing Analysis... See 'results.txt' and 'itTable.txt' files dhcp-254-235-238:SpedeSTEM-master bcarstens\$ 🗌 Carstens, Demarest, Kim, Pelletier, Satler 24

spedeSTEM: subsampling

• Takes as input a settings and genetrees file.