Demographic inference based on Site frequency spectrum (SFS) – Part II

Vitor Sousa CE3C – center for ecology, evolution and environmental changes

> 2018 WSPG Cesky Krumlov 22 Jan 2020

> > vmsousa@fc.ul.pt





Outline part II

Example of Applications:

- Human dispersal out of Africa (high quality whole-genome) lessons on choice of models
- Human colonization of Siberia and America (ancient wholegenome data) - lessons on dealing with sequencing errors
- Deer mice colonization of Nebraska Sand Hills (targeted recapture data) – lessons on effects of filtering
- Inferring divergence times and gene flow in sawflies (ddRADseq data) – lessons from comparing models



Nourlangie, Kakadu National Park, NT, Australia

A genomic history of Aboriginal Australia

Anna-Sapfo Malaspinas^{1,2,3}*, Michael C. Westaway⁴*, Craig Muller¹*, Vitor C. Sousa^{2,3}*, Oscar Lao^{5,6}*, Isabel Alves^{2,3,7}*, Anders Bergström⁸*, Georgios Athanasiadis⁹, Jade Y. Cheng^{9,10}, Jacob E. Crawford^{10,11}, Tim H. Heupink⁴, Enrico Macholdt¹², Stephan Peischl^{3,13}, Simon Rasmussen¹⁴, Stephan Schiffels¹⁵, Sankar Subramanian⁴, Joanne L. Wright⁴, Anders Albrechtsen¹⁶, Chiara Barbieri^{12,17}, Isabelle Dupanloup^{2,3}, Anders Eriksson^{18,19}, Ashot Margaryan¹, Ida Moltke¹⁶, Irina Pugach¹², Thorfinn S. Korneliussen¹, Ivan P. Levkivskyi²⁰, J. Víctor Moreno-Mayar¹, Shengyu Ni¹², Fernando Racimo¹⁰, Martin Sikora¹, Yali Xue⁸, Farhang A. Aghakhanian²¹, Nicolas Brucato²², Søren Brunak²³, Paula F. Campos^{1,24}, Warren Clark²⁵, Sturla Ellingvåg²⁶, Gudjugudju Fourmile²⁷, Pascale Gerbault^{28,29}, Darren Injie³⁰, George Koki³¹, Matthew Leavesley³², Betty Logan³³, Aubrey Lynch³⁴, Elizabeth A. Matisoo-Smith³⁵, Peter J. McAllister³⁶, Alexander J. Mentzer³⁷, Mait Metspalu³⁸, Andrea B. Migliano²⁹, Les Murgha³⁹, Maude E. Phipps²¹, William Pomat³¹, Doc Reynolds⁴⁰, Francois-Xavier Ricaut²², Peter Siba³¹, Mark G. Thomas²⁸, Thomas Wales⁴¹, Colleen Ma'run Wall⁴², Stephen J. Oppenheimer⁴³, Chris Tyler-Smith⁸, Richard Durbin⁸, Joe Dortch⁴⁴, Andrea Manica¹⁸, Mikkel H. Schierup⁹, Robert A. Foley^{1,45}, Marta Mirazón Lahr^{1,45}, Claire Bowern⁴⁶, Jeffrey D. Wall⁴⁷, Thomas Mailund⁹, Mark Stoneking¹², Rasmus Nielsen^{1,48}, Manjinder S. Sandhu⁸, Laurent Excoffier^{2,3}, David M. Lambert⁴ & Eske Willerslev^{1,8,18}

Nature(2016)



Ewaninga Rock Carvings Conservation Reserve, NT, Australia

Australia harbors some of the oldest modern human remains outside Africa



Many sites and remains dated to be older than 40 kya, suggesting a human settlement 47.5-55 kya

One wave out of Africa vs Two waves out of Africa



83 high-coverage Aboriginal Australians genomes



Average depth of coverage: 65x Very good quality of genotype calls

Effect of depth of coverage on SFS



 Compared 2D SFS based on depth of coverage of observed data (mean larger than >20x), with a distribution 8 times smaller.

A note on recovering the SFS from genomic data a) Low depth of coverage, no GQ filter, allowing missing data

-1.0

-2.0

2.5

-2,5

-1,5

-2,0

True SFS

-1.0

Sample SFS -1.5

- Simulation study
- Low depth of coverage and missing data lead to biased SFS towards rare variants

Singletons

0.20

0.15

0.10

0.05

0.00

Pop1

Pop2

Relative SFS





True SFS

83 high-coverage Aboriginal Australians genomes





★ Archaic human genomes:

- 1 Neanderthal (~66 kya)
- 1 Denisovan (~52 kya)

Mutation rate assumed 1.25 x 10⁻⁸ /site/gen Scally and Durbin (2012) *Nat. Rev. Genet.*

Generation time

29 years/gen Fenner (2005) *Am. J. Phys. Anthropol.*

Since we want to infer demography we tried to minimize the number of sites affected by selection:

- 985 1Mb blocks outside genic regions and CpG islands (~4.3 Million SNPs)
- 5 dimensional SFS (16,875 entries)
- Confidence intervals obtained using block-bootstrap

Towards a model to test the hypotheses: One vs Two waves Out of Africa

- Data (SFS)
 - (Re-)Define model (hypotheses to test)
- Run fastsimcoal2
- Estimates!
 - Assess the fit to the data

Do you have an outgroup?

- Yes use the derived (unfolded) SFS
- No use the minor allele frequency spectrum (folded)

Do you have monomorphic sites?

- Yes then, given a mutation rate you can infer the absolute times and effective sizes
- No then all your estimates need to be relative to a fixed parameter (fixed Ne or fixed time)

We always get results...

Evidence of two waves Out of Africa:

- Old split leading to colonization of Australia (81kya)
- More recent split leading to colonization of Eurasia (67 kya)



Towards a model incorporating Neanderthal and Denisovan admixture



- Non-African populations: 1-4% estimated Neanderthal admixture
- Aboriginal Australians and New Guineans: 3-6% estimated Denisovan admixture
- Archaic admixture can affect times of split estimates

Evidence of archaic introgression



Total length (Mb) of:

- Putative Denisovan haplotype (PDH)
- Putative Neanderthal haplotypes (PNH)

Unadmixed Australo-Papuans

Accounting for shared ancestry of Neanderthal and Denisovan



Admixture occurs between modern humans and:

- Denisovan-related (D.R.) population
- Neanderthal-related (N.R.) population

Two-waves out of Africa



West Africans

Unsampled

East Africa

Europeans

East

Asians

Australians

Present

- Two different divergence times (∆t >> 0)
- Two independent bottlenecks associated with the two Out of Africa events

Two-waves out of Africa



Two-waves out of Africa



- Two different divergence times (Δt >> 0)
- Two independent bottlenecks associated with the two Out of Africa events



West ghost Eurasians Australians Africans

One wave out of Africa



- Similar divergence times (∆t close to zero)
- One single bottlenecks associated with the Out of Africa events
- A major admixture pulse with Neanderthal





- Similar divergence time (∆t close to zero)
- Bottleneck associated with the Out of Africa event



- Similar divergence time (∆t close to zero)
- Bottleneck associated with the Out of Africa event
- A major admixture pulse with Neanderthal in ancestors of all non-Africans



- Similar divergence time (∆t close to zero)
- Bottleneck associated with the Out of Africa event
- A major admixture pulse with Neanderthal in ancestors of all non-Africans



Model captures aspects about the observed data



What entries are not well fitted?



Pagani et al (2016) suggests two waves: Papuan genomes with signature of admixture with humans from first wave (at least 2% of their genome).

Model captures the higher derived allele sharing between Eurasians and Yoruba



Australia Europe Yoruba Chimp or East Asian

D-statistics suggest that Yoruba and Eurasians share more derived alleles than Yoruba and Australians



Summary

Aboriginal Australians genomes support a single major wave out of Africa

- Accounting for archaic admixture with Neanderthal and Denisovan was crucial to understand population divergence
- Genomic data consistent with a single major dispersal event out of Africa (60-104 kya)
- Two major dispersal waves into Asia: Aboriginal Australians diverged
 51-72 kya from Eurasians



ARTICLE

The population history of northeastern Siberia since the Pleistocene

Martin Sikora^{1,43}*, Vladimir V. Pitulko^{2,43}*, Vitor C. Sousa^{3,4,5,43}, Morten E. Allentoft^{1,43}, Lasse Vinner¹, Simon Rasmussen^{6,41}, Ashot Margaryan¹, Peter de Barros Damgaard¹, Constanza de la Fuente^{1,42}, Gabriel Renaud¹, Melinda A. Yang⁷, Qiaomei Fu⁷, Isabelle Dupanloup⁸, Konstantinos Giampoudakis⁹, David Nogués–Bravo⁹, Carsten Rahbek⁹, Guus Kroonen^{10,11}, Michaël Peyrot¹¹, Hugh McColl¹, Sergey V. Vasilyev¹², Elizaveta Veselovskaya^{12,13}, Margarita Gerasimova¹², Elena Y. Pavlova^{2,14}, Vyacheslav G. Chasnyk¹⁵, Pavel A. Nikolskiy^{2,16}, Andrei V. Gromov¹⁷, Valeriy I. Khartanovich¹⁷, Vyacheslav Moiseyev¹⁷, Pavel S. Grebenyuk^{18,19}, Alexander Yu. Fedorchenko²⁰, Alexander I. Lebedintsev¹⁸, Sergey B. Slobodin¹⁸, Boris A. Malyarchuk²¹, Rui Martiniano²², Morten Meldgaard^{1,23}, Laura Arppe²⁴, Jukka U. Palo^{25,26}, Tarja Sundell^{27,28}, Kristiina Mannermaa²⁷, Mikko Putkonen²⁵, Verner Alexandersen²⁹, Charlotte Primeau²⁹, Nurbol Baimukhanov³⁰, Ripan S. Malhi^{31,32}, Karl-Göran Sjögren³³, Kristian Kristiansen³³, Anna Wessman^{27,34}, Antti Sajantila²⁵, Marta Mirazon Lahr^{1,35}, Richard Durbin^{22,36}, Rasmus Nielsen^{1,37}, David J. Meltzer^{1,38}, Laurent Excoffier^{4,5*} & Eske Willerslev^{1,36,39,40}*

Nature (2019)





Colonization of Siberia



Yana RHS (31,600 years ago) Whole-genome depth of coverage 25x



Kolyma (9,800 years ago) Whole-genome depth of coverage 14x



Hypothesis: Continuity vs **Replacement of populations**

Data: Ancient and presentday samples; 625 blocks of 1Mb (~1.5 Million SNP), far from genic regions and CpG islands

Method: Composite likelihood - fastsimcoal2 (Excoffier et al, 2013 Plos Genetics)

Europe Ancient Ancient Paelo-(Sardinia) North siberian siberian Siberians (Yana) (Kolyma)







Neo-

(Even)



Fast

Asia

(Han)

Hypothesis: Continuity vs Replacement of populations

For instance:

 $\beta = 1$ indicates continuity: Kolyma descends from Yana

 $\beta = 0$ indicates replacement of Yana by Kolyma











Site frequency spectrum is affected by damage patterns in ancient DNA

- High proportion of singletons in Kolyma probably reflect errors
- Thus, all analyses were performed discarding the singletons



Model comparison and likelihood profiles consistent with replacement with gene flow



Model comparison and likelihood profiles consistent with replacement with gene flow













Estimates of best nested model indicate replacement with gene flow



Siberia and colonization of the Americas



Yana RHS (31,600 years ago) Whole-genome depth of coverage 25x



USR1 (11,500 years ago) Alaska

Kolyma (9,800 years ago) Whole-genome depth of coverage 14x



Estimates consistent with replacement with gene flow



- Kolyma is the closest population to Native Americans (USR1 and Karitiana)
- Native Americans with a contribution of up to 20% from Yana

Summary: 3 migration waves



• Ancient North Siberians (Yana) reached Siberia before 30 ka (thousand-years ago)



Summary: 3 migration waves

- Ancient North Siberians (Yana) reached Siberia before 30 kya
- Paleo-Siberians (Kolyma) migrated after Last Glacial Maximum (26.5 ka)
- Native-Americans are closer to Kolyma, with 20% of Yana contribution







Summary: 3 migration waves

- Ancient North Siberians (Yana) reached Siberia before 30 ka
- Paleo-Siberians (Kolyma) likely migrated after Last Glacial Maxima
- Native-Americans are closer to Kolyma, with 20% of Yana contribution
- Paleo-Siberians (Kolyma) were replaced by Neo-Siberians, likely associated with the cooler period "Younger Dryas" (12.8-11.5 ka)





3rd migration wave

Deer mice from Nebraska Sand Hills



S. Pfeifer, S. Laurent, V. Sousa, C. Linnen, H. Hoekstra, L. Excoffier, J. Jensen

Coat color adaptation in deer mice *Peromyscus maniculatus*

- Habitat (soil color) correlated with coat phenotype
- Field experiments suggest that light color confers selective advantage against visually hunting predators
- Nebraska Sand Hills were formed 8000 to 15,000 years ago



Linnen et al (2013) Science

Pfeifer*, Laurent*, Sousa* et al (in press) MBE

A transect across the Sand Hills (ON and OFF)

Sample locations "off" and "on" the Sand Hills

- 11 populations
- 330 individuals



- Genomic data (NGS) data
 - Target 10,000 random 1.5kb regions
 - 185kbp region comprising the *Agouti* gene
- Phenotypic data for each individual



Evidence for isolation by distance but three groups



43.5

43.0

42.5

Latitude 42.0

41.5

41.0

40.5

Model-based inference

Is there evidence of gene flow between Off and On the Sand Hills?



Estimates based on the joint **3D site frequency spectrum** (SFS): - folded SFS with 140,358 SNPs

Deer mice: Pairwise marginal 2D SFS Since we did not have an outgroup we used the folded SFS



Estimates support south colonization and high gene flow levels

- Recent time of colonization of Sand Hills ~3-5 kya, younger than formation of Sand Hills 8-15 kya
- High migration rates across all populations, inferred for all models

Migration rates above/below arrows in units of 2Nm, i.e. average number of immigrants per generation.



Deer mice: Model fit to marginal SFS



Some lessons I learned working with the deer mice data

- Be carefull when applying Hardy-Weinberg filters to your data
- Be carefull when filtering on depth of coverage applying the same thresholds for all individuals

The depth of coverage varied considerably across individuals



- Applying the same threshold for all individuals can lead to biases
- Apply a filter on DP for each individual

Effect of DP filters on the SFS Simulation study



Effect of HW filtering on demographic estimates Removing sites with HWE excess and deficit leads to different estimates



Sawflies and RAD data

MOLECULAR ECOLOGY

Molecular Ecology (2016)

doi: 10.1111/mec.13972

History, geography and host use shape genomewide patterns of genetic variation in the redheaded pine sawfly (*Neodiprion lecontei*)

ROBIN K. BAGLEY,* VITOR C. SOUSA,† MATTHEW L. NIEMILLER‡ and CATHERINE R. LINNEN*

*Department of Biology, University of Kentucky, Lexington, KY 40506, USA, †cE3c - Centre for Ecology, Evolution and Environmental Changes, Faculdade de Ciências, Universidade de Lisboa, 1749-016 Lisboa, Portugal, ‡Illinois Natural History Survey, Prairie Research Institute, University of Illinois Urbana-Champaign, Champaign, IL 61820, USA



Sawflies Neodiprion lecontei

- Hymenoptera
- Plant-feeding insects
- Pine tree specialists



ddRAD seq data

- 80 individuals from 77 localities and 13 host species
- 100 bp paired-end reads, mapped to reference genome of *N. lencontei*
- Depth of coverage filter DP>10



Given the detected three groups (North, Central, South):

- What is the the population tree topology?
- What are the split times?
- What are the migration levels among groups?

Comparing models with composite likelihoods

- Fastsimcoal2
 likelihood is "correct"
 if all SNPs are
 independent
- We can then compare the model likelihoods using Akaike Information Criterion (AIC)



Effective size (Ne)

Composite likelihood provide unbiased maximum likelihood parameter estimates, but the likelihoods are inflated

A strategy to compare models



- 2. Create a dataset with all SNPs (including linked SNPs)
- For each model, obtain the parameters that maximize the likelihood (this is ok even with linked sites!) and the corresponding expected SFS
- Create a dataset with "independent" SNPs (1 SNP per RAD tag)
- Given the expected SFS of each model, compute the "correct" likelihood for each model with the dataset with independent SNPs
- 6. Compare models with AIC



"Correct" likelihood for each model

Comparing alternative models

Table 2 Summary of the likelihoods for the sixteen demographic models tested. Lhood (ALL SNPs) and Lhood (1 SNP) correspond to the mean likelihood computed with the data sets containing 'all SNPs' (including monomorphic sites) and a 'single SNP' (without monomorphic sites) per RAD locus, respectively. Mean likelihoods were computed based on 100 expected site frequency spectra simulated according to the parameters that maximized the likelihood of each model. Topology names for each model are as indicated in Fig. S1 (Supporting information). AIC scores and relative likelihoods (Akaike's weight of evidence) were calculated based on the 'single SNP' data set following Excoffier *et al.* 2013.

Topology	Migration allowed?	Exponential growth?	North bottleneck?	log ₁₀ (Lhood) ALL SNPs	log ₁₀ (Lhood) 1 SNP	# Parameters	AIC	ΔΑΙϹ	Relative likelihood
North–South	No	No	No	-46502.02	-7381.4	7	34006.70	75.69	0.000
North-Central	No	No	No	-46475.82	-7369.0	7	33949.44	18.43	0.000
South-Central	No	No	No	-46502.18	-7381.6	7	34007.60	76.59	0.000
Trifurcation	No	No	No	-46501.54	-7380.4	5	33998.07	67.06	0.000
North-South	Yes	No	No	-46470.49	-7365.0	15	33947.25	16.24	~0.000
North–Central	Yes	No	No	-46462.24	-7361.5	15	33931.01	0.00	0.851
South-Central	Yes	No	No	-46467.69	-7363.8	15	33941.57	10.56	0.004
Trifurcation	Yes	No	No	-46470.28	-7364.7	11	33937.93	6.91	0.027
North–South	Yes	Yes	No	-46469.48	-7362.8	18	33942.91	11.90	0.002
North–Central	Yes	Yes	No	-46461.17	-7361.7	18	33937.82	6.80	0.028
South-Central	Yes	Yes	No	-46463.73	-7363.9	18	33948.15	17.13	~0.000
Trifurcation	Yes	Yes	No	-46467.72	-7363.3	14	33937.39	6.37	0.035
North–South	Yes	Yes	Yes	-46467.45	-7361.5	20	33940.86	9.85	0.006
North–Central	Yes	Yes	Yes	-46461.25	-7362.1	20	33943.82	12.81	0.001
South-Central	Yes	Yes	Yes	-46463.58	-7364.1	20	33953.08	22.07	0.000
Trifurcation	Yes	Yes	Yes	-46466.06	-7362.4	16	33936.93	5.92	0.044

Estimates favors a scenario where North and Central diverged more recently with asymmetric gene flow



The inferred population tree topology and divergence times are consistent with divergence and range expansion from different refugia after LGM

Summary

- Fastsimcoal2 can be applied to RAD seq data
- We used a strategy to obtain (as close as possible) the "correct" likelihood by dividing the data into blocks, inferring the expected SFS for each model with ALL SNPs, and then re-computing the "true" likelihood with independent SNPs (1 SNP per block)
- Despite the reduced number of SNPs we were able to discriminate models based on their likelihoods

Protocol for model comparison based on AIC when we have independent SNPs

- Get the observed SFS
- Define the alternative models
- Perform 50-100 runs under each model
- Select the runs with maximum likelihood under each model
- Compute the AIC (Akaike information critera) for each model
- Select the model with minimum AIC

Estimating SFS from observed data

- The sample size can vary across SNPs due to missing data
- How to deal with missing data?

	Freq. derived	Sample size	Rel. freq
SNP1	1	16	1/16
SNP2	6	12	1/2
SNP3	1	12	1/12
SNP4	6	16	3/8



Estimating SFS from observed data

- The sample size can vary across SNPs due to missing data
- How to deal with missing data?

	Freq. derived	Sample size	Rel. freq
SNP1	1	16	1/16
SNP2	6	12	1/2
SNP3	1	8	1/12
SNP4	6	16	3/8



Estimating SFS from observed data

- The sample size can vary across SNPs due to missing data
- How to deal with missing data?
- Solution:
 - Find minimimum sample size
 - Resample without replacement

	Freq. derived	Sample size	Rel. freq
SNP1	1	16	1/16
SNP2	6	12	1/2
SNP3	1	8	1/12
SNP4	6	16	3/8



Gavel et al. (2011) PNAS

Acknowledgements

Martin Sikora Laurent Excoffier Isabelle Dupanloup Stephan Peischl Eske Willerslev



Thank you!



FONDO NAZIONALE SVIZZERO SWISS NATIONAL SCIENCE FOUNDATION Danmarks

Grundforskningsfond Danish National Research Foundation

FONDS NATIONAL SUISSE SCHWEIZERISCHER NATIONALFONDS

NF



Museum of Comparative Zoology

Catherine R. Linnen Stefan Laurent Jeffrey D. Jensen Susanne Pfeifer Hopi E. Hoekstra Laurent Excoffier













EC Fundação para a Ciência e a Tecnologia MINISTÉRIO DA CIÊNCIA, TECNOLOGIA E ENSINO SUPERIOR

UID/BIA/00329/2015-2018 UID/BIA/00329/2019 CEECIND/02391/2017



MCSA 2018-2020: MAPgenome (N.799729)