Custom bioinformatic pipelines for community DNA barcoding Aug 2014

Douglas Chesters, 中国科学院动物研究所



Sections

- Building a reference library
- Taxonomic assignment to queries
 - Distance based
 - Phylogenetic based

Building a reference database ACGITCTGGCG TIGGTATGTAGO TATTCAAACTGG ACTGACGC TGTCATC CGTCCGTG CGTCCGTG

CGT

567

-GT

AC

~

T

G

Customized DNA barcoding



Reference / Library



Apis mellifera AAAAAATTTGGG



Apis cerana AAAAAAAAAGGG



Bombus terrestris TTTTTTTTTGGG



Bombus terrestris TTTTTTTTTGGA



Lasioglossum zephyrum GGGGGGGGGGGGG

Queries



Specimen 1 AAAAATTTTGGG



Specimen 2 AAACCAAAAGGG



Specimen 3 TTGGGGGGGGGGG

Where do we get data, Queries



Where do we get data, References (the data mining approach)

Index of ftp://ftp.ncbi.nl × +	
Itp://ftp.ncbi.nlm.nih.gov/genbank/	∇
ypmv1.seq.yz	23320 KD
📄 gbinv10.seq.gz	47124 KB
gbinv11.seq.gz	24195 KB
gbinv12.seq.gz	17628 KB
📄 gbinv13.seq.gz	22563 KB
gbinv14.seq.gz	21770 KB
📄 gbinv15.seq.gz	21065 KB
📄 gbinv16.seq.gz	16349 KB
gbinv17.seq.gz	60509 KB
📄 gbinv18.seq.gz	72508 KB
📄 gbinv19.seq.gz	50547 KB
☐ gbinv2.seq.gz	37811 KB
📄 gbinv20.seq.gz	17160 KB
gbinv21.seq.gz	18676 KB
📄 gbinv22.seq.gz	1696 KB
📄 gbinv23.seq.gz	53054 KB
binv24.seq.gz	56254 KB
📄 gbinv25.seq.gz	38932 KB
📄 gbinv26.seq.gz	14873 KB
📄 gbinv27.seq.gz	14547 KB
gbinv28.seq.gz	16347 KB
abinv29.sea.az	18203 KB

Building a reference library from mined data

Index of ftp://ftp.ncbi.nl × +	
ftp://ftp.ncbi.nlm. nih.gov /genbank/	~
gomvised.85	23320 KD
gbinv10.seq.gz	47124 KB
gbinv11.seq.gz	24195 KB
gbinv12.seq.gz	17628 KB
gbinv13.seq.gz	22563 KB
gbinv14.seq.gz	21770 KB
gbinv15.seq.gz	21065 KB
gbinv16.seq.gz	16349 KB
gbinv17.seq.gz	60509 KB
gbinv18.seq.gz	72508 KB
gbinv19.seq.gz	50547 KB
gbinv2.seq.gz	37811 KB
gbinv20.seq.gz	17160 KB
gbinv21.seq.gz	18676 KB
gbinv22.seq.gz	1696 KB
gbinv23.seq.gz	53054 KB
gbinv24.seq.gz	56254 KB
gbinv25.seq.gz	38932 KB
gbinv26.seq.gz	14873 KB
gbinv27.seq.gz	14547 KB
gbinv28.seq.gz	16347 KB
abinv29.sea.az	18203 KB



Reference / Library



Apis mellifera AAAAAATTTGGG



Apis cerana AAAAAAAAAGGG



Bombus terrestris TTTTTTTTGGG



Bombus terrestris TTTTTTTTGGA



Lasioglossum zephyrum GGGGGGGGGGGGGG

Building a reference library from mined data

MOLECULAR ECOLOGY RESOURCES

Molecular Ecology Resources (2014)

doi: 10.1111/1755-0998.12256

Automated DNA-based plant identification for large-scale biodiversity assessment

ANNA PAPADOPOULOU,* DOUGLAS CHESTERS,† INDIANA CORONADO,‡ GISSELA DE LA CADENA,* ANABELA CARDOSO,* JAZMINA C. REYES,‡ JEAN-MICHEL MAES,§ RICARDO M. RUEDA‡ and JESÚS GÓMEZ-ZURITA*

*Animal Biodiversity and Evolution, Institut de Biologia Evolutiva (CSIC-Univ. Pompeu Fabra), 08003 Barcelona, Spain, †Key Laboratory of Zoological Systematics and Evolution, Institute of Zoology, Chinese Academy of Sciences, Beijing, China, ‡Herbario y Jardín Botánico Ambiental, Universidad Nacional Autónoma de Nicaragua, León, Nicaragua, §Museo Entomológico de León, León, Nicaragua

Download data	 Automate using wget (linux), curl (OSX)
	• Format white facta
Process database	 Format, .go to fasta Sequence IDs, accession or GI number, Species name Remove sequences too short / long. Remove duplicates
Select taxon	 NCBI taxonomy IDs 2759=Eukaryota, 33208=Metazoa, 50557=Insecta
Find homologs in database	 Which queries? Software options, Blast or Uclust To consider, algorithm, strand
	. Trimming and Orientation
Process homologs	 Infining and Orientation De-replication: one per species or retain variation; exemplars, random, most complete, or most representative

Download data	 Automate using wget (linux), curl (OSX)
Process database	 Format, .gb to fasta Sequence IDs, accession or GI number, Species name Remove sequences too short / long. Remove duplicates
Select taxon	 NCBI taxonomy IDs 2759=Eukaryota, 33208=Metazoa, 50557=Insecta
Find homologs in database	 Software options, Blast or Uclust To consider, algorithm, strand
	• Trimming and Orientation
Process homologs	 Infining and Orientation De-replication: one per species or retain variation; exemplars, random, most complete, or most representative

Download data	 Automate using wget (linux), curl (OSX)
Process database	 Format, .gb to fasta Sequence IDs, accession or GI number, Species name Remove sequences too short / long. Remove duplicates
Select taxon	 NCBI taxonomy IDs 2759=Eukaryota, 33208=Metazoa, 50557=Insecta
Find homologs in database	 Software options, Blast or Uclust To consider, algorithm, strand
Process homologs	 Trimming and Orientation De-replication: one per species or retain variation; exemplars, random, most complete, or most representative

Download data	 Automate using wget (linux), curl (OSX)
Process database	 Format, .gb to fasta Sequence IDs, accession or GI number, Species name Remove sequences too short / long. Remove duplicates
Select taxon	 NCBI taxonomy IDs 2759=Eukaryota, 33208=Metazoa, 50557=Insecta
Find homologs in database	 Software options, Blast or Uclust To consider, algorithm, strand
Process homologs	 Trimming and Orientation De-replication: one per species or retain variation; exemplars, random, most complete, or most representative

Download data	 Automate using wget (linux), curl (OSX)
Process database	 Format, .gb to fasta Sequence IDs, accession or GI number, Species name Remove sequences too short / long. Remove duplicates
Select taxon	 NCBI taxonomy IDs 2759=Eukaryota, 33208=Metazoa, 50557=Insecta
Find homologs in database	 Which queries? Software options, Blast or Uclust To consider, algorithm, strand
Process homologs	 Trimming and Orientation De-replication: one per species or retain variation; exemplars, random, most complete, or most representative

Download data	 Automate using wget (linux), curl (OSX)
Process database	 Format, .gb to fasta Sequence IDs, accession or GI number, Species name Remove sequences too short / long. Remove duplicates
Select taxon	 NCBI taxonomy IDs 2759=Eukaryota, 33208=Metazoa, 50557=Insecta
Find homologs in database	 Software options, Blast or Uclust To consider, algorithm, strand
Process homologs	 Trimming and Orientation De-replication: one per species or retain variation; exemplars, random, most complete, or most representative

Reference / Library



Apis mellifera AAAAAATTTGGG



Apis cerana AAAAAAAAAGGG



Bombus terrestris TTTTTTTTGGG



Bombus terrestris TTTTTTTTGGA



Lasioglossum zephyrum GGGGGGGGGGGGGG

Taxonomic assignment: TACTGACGCGGCG TTGGTATGTAGO GC **Distance** based CGTCCGTC

AX

G

Taxonomic assignment: Distance based

- Advantage
 - MSA not required
 - Computational
- Disadvantage
 - accuracy
- Blast and Uclust
 - local versus global pairwise alignment
 - Distance measure and gaps
 - Gaps as single event, ignore terminal gaps, ambiguities

Improve scoring of a query to reference pair

REFERENCE-Apis cerana QUERY-Specimen 1





Papadopoulou et al. 2014. *Molecular Ecology Resources*

Jung based **Taxonomic assignment:** CGAGCGTCI ATTCAAACTGG

CTTATTA

GT

AGP

AGC

T

Retrieve 16S homologs	 Download mtgenomes for mammal refseq taxa Extract 16S rRNA (IrRNA) Use these as Blast queries, retrieve only full length entries (>1500)
Species filtering	 Random exemplar / longest / most representative sequence (Chesters and Zhu 2014)
Multiple Sequence Alignment	 References – Muscle (Edgar 2004) Queries to Reference profile – Pynast (Caporaso et al. 2010) MSA refinement
Infer Phylogenetic relationships	 RAxML (Stamatakis 2006) Branch swaps for references constrained according to high quality previously published tree (Bininda-Emonds et al. 2007)
Taxonomic assignment	 Requires complete taxonomies for reference members of the tree These are assigned to queries on the tree

Retrieve 16S homologs	 Download mtgenomes for mammal refseq taxa Extract 16S rRNA (IrRNA) Use these as Blast queries, retrieve only full length entries (>1500)
Species filtering	 Random exemplar / longest / most representative sequence (Chesters and Zhu 2014)
Multiple Sequence Alignment	 References – Muscle (Edgar 2004) Queries to Reference profile – Pynast (Caporaso et al. 2010) MSA refinement
Infer Phylogenetic relationships	 RAxML (Stamatakis 2006) Branch swaps for references constrained according to high quality previously published tree (Bininda-Emonds et al. 2007)
Taxonomic assignment	 Requires complete taxonomies for reference members of the tree These are assigned to queries on the tree

Retreive 16S homologs	 Download mtgenomes for mammal refseq taxa Extract 16S rRNA (IrRNA) Use these as Blast queries, retrieve only full length entries (>1500)
Species filtering	 Random exemplar / longest / most representative sequence (Chesters and Zhu 2014)
Multiple Sequence Alignment	 References – Muscle (Edgar 2004) Queries to Reference profile – Pynast (Caporaso et al. 2010) MSA refinement
Infer Phylogenetic relationships	 RAxML (Stamatakis 2006) Branch swaps for references constrained according to high quality previously published tree (Bininda-Emonds et al. 2007)
Taxonomic assignment	 Requires complete taxonomies for reference members of the tree These are assigned to queries on the tree

Retreive 16S homologs	 Download mtgenomes for mammal refseq taxa Extract 16S rRNA (IrRNA) Use these as Blast queries, retrieve only full length entries (>1500)
Species filtering	 Random exemplar / longest / most representative sequence (Chesters and Zhu 2014)
Multiple	 References – Muscle (Edgar 2004)
Sequence	• Queries to Reference profile – Pynast (Caporaso et al. 2010)
Alignment	MSA refinement
Alignment	MSA refinement
Alignment Infer Phylogenetic relationships	 MSA refinement RAxML (Stamatakis 2006) Branch swaps for references constrained according to high quality previously published tree (Bininda-Emonds et al. 2007)
Alignment Infer Phylogenetic relationships	 MSA refinement RAxML (Stamatakis 2006) Branch swaps for references constrained according to high quality previously published tree (Bininda-Emonds et al. 2007)

Retreive 16S homologs	 Download mtgenomes for mammal refseq taxa Extract 16S rRNA (IrRNA) Use these as Blast queries, retrieve only full length entries (>1500)
Species filtering	 Random exemplar / longest / most representative sequence (Chesters and Zhu 2014)
Multiple Sequence Alignment	 References – Muscle (Edgar 2004) Queries to Reference profile – Pynast (Caporaso et al. 2010) MSA refinement
Infer Phylogenetic relationships	 RAxML (Stamatakis 2006) Branch swaps for references constrained according to high quality previously published tree (Bininda-Emonds et al. 2007)
Taxonomic assignment	 Requires complete taxonomies for reference members of the tree These are assigned to queries on the tree

Retreive 16S homologs	 Download mtgenomes for mammal refseq taxa Extract 16S rRNA (IrRNA) Use these as Blast queries, retrieve only full length entries (>1500)
Species filtering	 Random exemplar / longest / most representative sequence (Chesters and Zhu 2014)
Multiple Sequence Alignment	 References – Muscle (Edgar 2004) Queries to Reference profile – Pynast (Caporaso et al. 2010) MSA refinement
Infer Phylogenetic relationships	 RAxML (Stamatakis 2006) Branch swaps for references constrained according to high quality previously published tree (Bininda-Emonds et al. 2007)
Taxonomic assignment	 Requires complete taxonomies for reference members of the tree These are assigned to queries on the tree

Overview:

Phylogeny for taxonomic assignment

- De novo tree of both references and queries
 bagpipe_phylo.pl
- New queries added to existing phylogeny
 - EPA
 - pplacer

Taxonomic assignment on de novo tree of references and queries

- bagpipe_phylo
 - New standalone implementation of a key component of BAGpipe (Papadopoulou et al. 2014)
 - Requires phylogeny of references and queries, and the NCBI taxonomy database
 - Will be available soon ...

unpublished

Example: Next Generation Sequencing of Leech Gut Contents

- Mammal 16S rDNA
- Degraded DNA
- Only short fragments retrievable.



unpublished, in collaboration with Douglas Yu (KIZ)



Cercartetus_cau Peta-tus_biarios Elephantulus_edwardii Chrysochloris_asiatica Episoriculus_fumidus Blarinella_griselda Hylomys_suillus Myzopoda_aurita Scotophilus_viridis Histiotus_magellanicus Hypsugo_anchietae Hypsugo_cadornae Nyctalus_leisleri Plecotus_auritus Pipistrellus_subflavus Lasiurus_atratus Myotis_capaccinii Myotis_ridleyi Myotis_fortidens Murina_ussuriensis Miniopterus_tristis Otomops_martiensseni Thyroptera_tricolor Chiroderma_trinitatum Artibeus_lituratus Dermanura_rava Sturnira tildae Anoura_caudifer Micronycteris_hirsuta Phylloderma_stenops Chrotopterus_auritus Pteronotus_quadridens Hipposideros_armiger Rhinolophus_paradoxolophus Cormura_brevirostris Taphozous_australis Pteralopex_atrata Rousettus_leschenaultii Mani 1100

> Melusys_ushus Mirounga_leonina

> > Pusa_siemca Eumetopias_jubatus

Maltes cane (Cana) We DEmigare Units Rhindo 193-55765 rea

Cercartetus_ca

Microtus_levis Castor capadonsis Coendou insidiosus Marmota himalayana

Lepus americanus

ademus j tropum

Rattus Juscipes

Miopithecus_talapoin

. UNIVERSE & 84454

Cercopithecus_denti

Chiorocebus_pygerythrus

Papio kindae I Uniques_498091 uniques_1535 uniques_77184

Presbytis_melalophos

augs_25863 o san i enso

Hylobates_agilis

Chiropotes_albinasus Callithrix jacchus

JURIQUES_6806

Daubentonia_madagascariensis Lepilemur_ruficaudatus Galago_senegalensis

Tupaia_tana Gazella_bennettii Nanger_dama Kobus_leche Cephalophus_nigrifrons

Oryx_gazella

Pamaliscus-lunatus Hemitragus_jayakari Biskudots, Atayaura Biskudots, Atayaur Costiluosinge Sumatraensis Union Constant Bubalus, depressicornis ranelanhus, eccintus Tragelaphus_scriptus

localegedsemionus

Ervus nippon untracus reevesi

្លាំនុក្តីតិ ត្រូង ក្តីដែលការទេ ហើង 1224 = Orinithoring incluis_anatings ដក់ទីទា

uniques_4232 uniques_27219 uniques_7707# uniques 76849 uniques 589 uniques 3042 uniques_34552 uniques_29120 uniques_52485 uniques 44622 uniques_99582 uniques_46466 uniques 2209 uniques_91978 uniques 84830 uniques_5918# uniques_95778 uniques_102340 HYAENIDAE_Hyaena_hyaena~ HYAENIDAE_Crocuta_crocutauniques_82158 uniques_3772* uniques_87 uniques_76209 uniques_58 uniques_53 FELIDAE_Acinonyx_jubatus-FELIDAE_Lynx_rufus* FELIDAE Neofelis nebulosa-FELIDAE_Uncia_uncia FELIDAE_Panthera_tigris unique FELIDAE_Panthera_onca^ FELIDAE_Panthera_pardus FELIDAE_Felis_catus FELIDAE_Panthera_leo~ uniques FELIDAE_Prionailurus_bengalensis/ FELIDAE_Puma_concolor uniques 45630 uniques_57 uniques 8807 uniques_2397 uníques 297* uniques 98210 uniques 42251 uniques 34634 uniques 82037 uniques 1021 HERPESTIDAE_Herpestes_javanicus CANIDAE_Vulpes_zerda-CANIDAE_Vulpes_vulpes CANIDAE_Vulpes_coisac CANIDAE_Vulpes_coisac CANIDAE_Nyctereutes_procyonoides CANIDAE_Cuon_alpinus CANIDAE_Cuon_alpinus CANIDAE_Canis_latrans CANIDAE_Canis_lycaom CANIDAE_Canis_lupus CANIDAE_Chrysocyon_brachyurus/ URSIDAE_Ailuropoda_melanoleuca URSIDAE Tremarctos ornatus URSIDAE Arctodus simus URSIDAE Melursus_ursinus URSIDAE_Helarctos_malayanus URSIDAE_Ursus_thibetanus URSIDAE Ursus maritimus URSIDAE Ursus maritimus URSIDAE Ursus_arctos URSIDAE Ursus_spelaeus URSIDAE_Ursus_americanus

bagpipe_phylo output: Top 14 ...

Shortest

Query	Support	Shared Taxon	Distance to Reference Leaf
uniques_1307	44	Melogale_moschata	0.0000080
uniques_2811	65	Bos_indicus	0.0000080
uniques_5424	22	Capricornis_sumatraensis	0.0000080
uniques_617	94	Martes_flavigula	0.0000080
uniques_961	79	Homo_sapiens	0.0000080
uniques_80200	77	Ursus_thibetanus	0.01261555
uniques_741	48	Canis_latrans	0.013661082
uniques_7756	87	Papio_anubis	0.026794699
uniques_1993	94	Papio_anubis	0.037582712
uniques_4232	69	Manis_pentadactyla	0.043805173
uniques_20547	14	Felidae	0.04743098
uniques_2229	25	Rattus_rattus	0.04824355
uniques_2357	33	Caprinae	0.051507736
uniques_468	30	Macaca	0.07541466

bagpipe_phylo results: shared taxonomies assigned to queries

Query	Support	Shared Taxon	Shortest Distance to Reference Leaf
uniques_20547	14	Felidae	0.04743098
uniques_68378	31	Panthera_tigris	0.159440927



Evolutionary Placement Algorithms

- New data assigned to existing tree
 - EPA in RAXML
 - pplacer



Evolutionary Placement in RAXML



Evolutionary Placement using pplacer



Conclusion

- Building an automated pipeline allows for easy replication / updating
- Customization:
 - can avoid the fitting an unsuited method
 - novelty!

Acknowledgments

- PI: 朱朝东
- Collaborators for the current work:
 - Anna Papadopoulou & Jesus Gomez-Zurita
 - Douglas Yu





Thanks for listening

