

2 DR. SONAL SINGHAL (Orcid ID : 0000-0001-5407-5567)

4 Article type : Resource Article

6
8 **Squamate Conserved Loci (SqCL): a unified set of conserved loci for phylogenomics and population genetics of squamate reptiles**

10 Sonal Singhal^{1*}, Maggie Grundler¹, Guarino Colli², Daniel L. Rabosky¹

12 ¹Museum of Zoology and Department of Ecology and Evolutionary Biology, University of Michigan, Ann Arbor, MI 48109

14 ²Departamento de Zoologia, Universidade de Brasília, 70910-900 Brasília, Brazil

16 * corresponding author, sonal.singhal1@gmail.com

18 **Abstract**

The identification of conserved loci across genomes, along with advances in target capture methods and high-throughput sequencing, has helped spur a phylogenomics revolution by enabling researchers to gather large numbers of homologous loci across clades of interest with minimal upfront investment in locus design. Target capture for vertebrate animals is currently dominated by two approaches – anchored hybrid enrichment (AHE) and ultraconserved elements (UCE) – and both approaches have proven useful for addressing questions in phylogenomics, phylogeography, and population genomics. However, these two sets of loci have minimal overlap with each other; moreover, they do not include many traditional loci that have been used for phylogenetics. Here, we combine across UCE, AHE, and traditional phylogenetic gene locus sets to generate the Squamate Conserved Loci (SqCL) set, a single integrated probe set that can generate high-quality and highly complete data across all three loci types. We use these probes to generate data for 44 phylogenetically-disparate taxa that collectively span approximately 33%

This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the [Version of Record](#). Please cite this article as [doi: 10.1111/1755-0998.12681](https://doi.org/10.1111/1755-0998.12681)

This article is protected by copyright. All rights reserved

30 of terrestrial vertebrate diversity. Our results generated an average of 4.29 Mb across 4709 loci per
individual, of which an average of 2.99 Mb was sequenced to high enough coverage ($\geq 10\times$) to use for
32 population genetic analyses. We validate the utility of these loci for both phylogenomic and population
genomic questions, provide a comparison among these locus sets of their relative usefulness, and suggest
34 areas for future improvement.

36 **Running Title**

SqCL: a unified locus set

38

Keywords

40 Target capture, phylogenomics, comparative population genomics, squamate reptiles

42 **Introduction**

For researchers working on biodiversity genomics, a primary challenge in project design is deciding
44 which portion of the genome to sequence for the organisms of interest. Given that whole-genome
sequencing remains prohibitively expensive for most organisms and most projects (but see (Therkildsen
& Palumbi 2016)), sequencing part of the genome allows researchers to affordably sample both more
46 individuals and species. There are many approaches to subsetting the genome for sequencing, including
48 transcriptome sequencing, restriction-aided digest methods (e.g., RAD sequencing), and target sequence
capture, each of which poses benefits and challenges (Jones & Good 2016). In the phylogenetics
50 community, targeting and sequencing conserved elements – *i.e.*, anchored hybrid enrichment (AHE;
(AHE, Lemmon *et al.* 2012) and ultraconserved elements (UCE; (Faircloth *et al.* 2012b) – has been applied
52 to infer phylogenies across broad phylogenetic scales (Crawford *et al.* 2012; Prum *et al.* 2015), resolve
rapid radiations (Giarla & Esselstyn 2015; Meiklejohn *et al.* 2016), and characterize phylogeographic
54 patterns (Brandley *et al.* 2015; Smith *et al.* 2014). Because these loci are fairly conserved across broad
phylogenetic scales (*i.e.*, all of arthropods (Faircloth *et al.* 2015) or all of angiosperms (Budenhagen *et al.*
56 2016)), researchers can use a common set of publicly available probe sequences for all their species of
interest, thus saving energy, time, and money.

58

The approaches targeting AHEs and UCEs are conceptually similar, though they are implemented
60 differently. In both approaches, the basic premise is to identify regions of the genome that are conserved
across deep phylogenetic scales and to design probes specific to these regions for use in target capture.

62 AHE loci are long (>1 kb), the probes targeting these loci cover most of the locus sequence, and the probe
sequences are about 15% divergent among organisms diverged across 200 million years (i.e., snakes and
64 geckos, Fig. S1A; (Zheng & Wiens 2016)). UCE loci tend to be shorter (500 - 800 bp), the probes only cover
the highly-conserved central 100 - 200 bp of these loci, and the probe sequences are very conserved (<5%
66 across snakes and geckos; Fig. S1A).

68 Research groups targeting conserved loci have focused either on AHEs or UCEs in generating data for
their clades of interest, either of which offers more than enough data to resolve most phylogenetic
70 questions. Unfortunately, AHEs and UCEs only have minor overlap in target loci. This creates a divide in
the field. Historically, researchers targeted a common set of mitochondrial and nuclear loci across diverse
72 species, enabling researchers to combine across data sets to create deeper, more fully-sampled trees (c.f.,
(Jetz *et al.* 2012; Pyron *et al.* 2013)). However, fully utilizing existing datasets is challenging if different
74 research groups have targeted distinct and largely independent locus sets. In this study, we create a
single inclusive locus set with applications to comparative population genomics, phylogeography, and
76 phylogenomics of squamate reptiles (lizards and snakes, ~200 million years of evolutionary history,
(Zheng & Wiens 2016)). This locus set – the squamate conserved locus set (SqCL) – combines across three
78 major sets of loci: AHEs, UCEs, and traditional genes used in squamate phylogenetics. We then test this
locus set on a phylogenetically diverse set of 56 individuals representing 44 squamate species, confirming
80 its efficacy and its usefulness for both population-scale and phylogenetic studies. We further highlight
areas of improvement in how these data are collected and analyzed.

82

Methods

84 *Samples*

To test the efficacy of the SqCL set, we targeted 16 of the most species-rich families in squamates that
86 span the entire phylogenetic breadth of the clade, resulting in 56 individuals from 44 species that reflect
the full diversity of 10,000+ squamate species. Importantly, this sampling consisted of multiple closely-
88 related congeneric species (Fig. 1; Table S1), allowing us to test how these markers resolved both
shallower and deeper phylogenetic relationships. These individuals were all collected as part of
90 ecological and macroecological studies in the Brazilian Cerrado over a ten year span from September 2005
to October 2015 (Colli *et al.* 2002). Full details on the samples used can be found in Table S1.

92

Probe design

94 To design the probes for the SqCL set, we started with publicly available sequences for each locus set. For
the AHEs, we used the sequence data for the AHE v2, as published in Ruane *et al.* (2015). This marker set
96 consists of 394 loci as identified from multiple vertebrate genomes, of which five loci had no match in
Anolis carolinensis. For the UCEs, we used the probe set Tetrapods-UCE-5Kv1 (accessed from
98 www.ultraconserved.org on 10 February 2016). Because some UCE probes overlap, we assembled them
using cap3 (Huang & Madan 1999), to result in 5061 unique targets. For the standard genes used in
100 squamate phylogenetics, we downloaded data matrices from two recent phylogenetic studies; these
datasets included 44 genes from approximately 160 tips (Wiens *et al.* 2012) and 12 genes across 4161 tips
102 (Pyron *et al.* 2013). These two gene sets had four overlapping genes, resulting in 52 genes of which five
were mitochondrial. Because mitochondrial DNA has much higher copy number than nuclear DNA,
104 capturing both genomic types simultaneously can lead to an excess of sequence reads mapping to the
mitochondrial genome (Bi *et al.* 2012). As such, we dropped these five mitochondrial genes, giving us a
106 total of 47 "traditional" nuclear loci that have traditionally been obtained using Sanger sequencing.

108 We then used blast (Camacho *et al.* 2009) to identify loci across the three sets that significantly overlapped
with each other; we identified and dropped 28 duplicate loci. For the remaining 5469 targets, we used
110 blast to search for homologous regions of this genome across 11 publicly available squamate reptile
genomes (Table S2), extracted the matching regions, and aligned across these regions using mafft v7.294
112 (Katoh & Standley 2013). We used these alignments to characterize how divergent the targeted sequences
were across genomes. We found that although the target sequences exhibited less divergence among
114 snakes, they tended to show equal divergence among "lizards" and between any given "lizard" and any
given snake. Given this, for every target, we included sequence representatives from two divergent
116 clades within the phylogeny, to better capture some of this variation in target sequence identity across
clades. For AHEs, we used both sequence from *Anolis carolinensis* and from either *Calamaria pavementata* or
118 *Python molurus*, as originally published in Ruane *et al.* (2015). For UCEs, we extended the central probes
until we accumulated more than 15% sequence divergence across a rolling mean of 10 bp. Previous
120 studies (Hugall *et al.* 2015) have shown that, beyond 15% sequence divergence, capture efficiency begins
to decline. We then extracted sequence data with these expanded coordinates from *A. carolinensis* and
122 *Gekko japonicus*; for the few targets for which we could not identify a homolog from *G. japonicus*, we
instead used *Ophisaurus gracilis*. For traditional phylogenetic genes, we used sequence data from *G.*
124 *japonicus* and *Boa constrictor*. We then screened all targets against the RepeatMasker database, identifying
7 targets that matched significantly to repeats. The final set consisted of 5462 targets, each represented by

126 two squamate sequences. Probes were designed across these targets at $\sim 2\times$ tiling density by MYcroarray
Inc (Ann Arbor, Michigan) resulting in 38,431 probes. These probes were then further filtered to remove
128 probes that matched to repeats or to multiple places in the *Anolis carolinensis* genome. The final probe set
consisted of 37,517 probes targeting 2.25 Mb of unique sequence. The total assembled sequence should be
130 greater as UCEs are designed to capture flanking regions.

132 *Data Collection*

From each individual, we extracted high molecular-weight DNA using a high-salt DNA extraction
134 method (Aljanabi & Martinez 1997) and then measured DNA quantity using a QuBit dsDNA BR Assay
Kit (ThermoFisher, cat. no. Q32850) and DNA purity with a NanoDrop (ThermoScientific). MYcroarray
136 then produced dual-barcoded libraries for each sample. Roughly 1.0 to 1.6 ng genomic DNA were
sheared using a Qsonica Q800RS sonicator, then size-selected to approximately 450 bp modal lengths
138 with SPRI beads. Sheared DNA were then end-repaired and adapter-ligated with the NEBNext Ultra
DNA Library Prep Kit for Illumina (NEB, cat. no. E7370), and index-amplified with custom primers for 6
140 cycles using HotStart HiFi Readymix (Kapa Biosystems, cat. no. KR0370). Following amplification,
libraries were purified and quantified with Quant-iT PicoGreen dsDNA assay (ThermoFisher, cat. no.
142 P7589). Roughly 100 ng across each of 8 individuals was pooled. Because capture efficiency typically has
phylogenetic signal (Bi *et al.* 2012; Cosart *et al.* 2011), we reduced bias by pooling individuals by
144 taxonomic family. These pools of 800ng were then dried to 7uL via vacuum centrifugation and used as
template for standard capture reactions following the MYbaits protocol v3. We modified the protocol
146 slightly to include xGEN Universal Blockers (Integrated DNA Technologies, cat. no. 1046636 and
1046639), which have been shown to improve target capture efficiency by up to $4\times$ (Portik *et al.* 2016).
148 Following a 12-cycle post-capture PCR, all 56 individuals were combined with an additional 8 frog
samples from another study (Larson, unpublished) and sequenced by Hudson Alpha on one 100 paired-
150 end run of a HiSeq 2500 v4.

152 *Data Analysis*

Our pipeline for SqCL facilitates both population genetic and phylogenomic analysis, reflecting the
154 potential use of these loci for questions at both shallow and deep scales of divergence. This pipeline is
influenced by the publicly available PHYLUCE pipeline (Faircloth 2015) but includes two primary
156 modifications. First, we implement species tree methods for phylogenetic inference, because these
methods generally outperform concatenated based approaches (Kubatko & Degnan 2007; Warnow 2011)

158 but see (Springer & Gatesy 2016). Second, we incorporate industry-standard SNP calling, filtering, and
phasing to enable population genetic analyses. This pipeline, along with documentation explaining its
160 implementation, is available at <https://github.com/singhal/SqCL>.

162 Following de-multiplexing, we removed adapters and low-quality regions from the reads using
Trimmomatic v0.36 (Bolger *et al.* 2014) and then merged overlapping reads using PEAR v0.9.10 (Zhang *et al.*
164 *al.* 2014). We then assembled the reads using default settings on the program Trinity v2.2.0 (Grabherr *et al.*
2011); for the few samples requiring memory in excess of 64 Gb, we used *in silico* read normalization to
166 thin the original data set. We matched contigs in each individual assembly to the original targets using
blat v36 (Kent 2002). We identified two types of matches: in the first, the contig and the target have a one-
168 to-one unique match; in the second, the target matches to multiple contigs with match scores within 10
orders of magnitude of the best match. We then used these match designations to create a pseudo-
170 reference genome (PRG) for each species. Here, we identified all contigs across all individuals in a given
species that match to a given target and then retained either the longest contig or the best matching contig
172 if it was a significantly better match than the next best matching contig (>3 orders of magnitude). We then
implemented phylogenomic and population genomic analyses as detailed below.

174
Assessing informativeness for phylogenomics

176 To facilitate phylogenomic analyses, we first extracted homologs for our target loci from *Gallus gallus*
(Hillier *et al.* 2004) for use as an outgroup. We then used mafft to generate alignments for each locus
178 sampled for ≥ 4 species (Kato & Standley 2013) and trimmed alignments to remove regions of low
quality using GBLOCKS (Castresana 2000). We inferred gene trees for each alignment using RAxML
180 v8.2.4 (Stamatakis 2006).

182 For each locus, we measured (1) phylogenetic informativeness, (2) certainty of gene trees inferred with
that locus, and (3) how clock-like a locus is. Empirical results show that maximizing these metrics can
184 improve the accuracy of topology and branch length inference. First, empirical results suggest that the
ideal loci are phylogenetically informative across evolutionary time – i.e., they should contain variable
186 markers at recent time scales while not exhibiting homoplasy at deeper time scales (Dornburg *et al.* 2014;
Gilbert *et al.* 2015). To characterize phylogenetic informativeness (PI) for each locus, we used the method
188 introduced by Townsend (2007) and implemented in TAPIR v1.1 (Faircloth *et al.* 2012a; Guindon *et al.*
2010). Measuring phylogenetic informativeness requires an ultrametric tree. We used the tree inferred

190 from our combined ASTRAL + BEAST analysis (see below); we rescaled the tree using the R package
`geiger` to have a root age that reflects estimates from the literature (Harmon *et al.* 2008; Zheng & Wiens
192 2016). Second, empirical results suggest gene trees should exhibit high tree certainty (Blom *et al.* 2016;
Salichos & Rokas 2013). We used a tree certainty measure that calculates the relative frequency of each
194 bipartition in a set of trees with respect to the frequency for the most common conflicting bipartition
(Salichos & Rokas 2013). Higher scores reflect a topology that shows greater stability across replicates. We
196 used RAxML to both infer 100 bootstraps for each locus and to calculate tree certainty across these
bootstraps (Salichos *et al.* 2014). Finally, empirical results suggest trees inferred with more clock-like
198 genes are more accurate (Doyle *et al.* 2015). We measured clocklikeness following the approach outlined
in Doyle *et al.* 2015, in which we compared tree likelihoods estimated by PAUP v4 for a gene tree forced
200 to be ultrametric to one that was not (Swofford 2003).

202 We then employed three phylogenetic approaches. First, we generated concatenated alignments by
marker type, defined partitions using PartitionFinder2 with the 'rcluster' algorithm (Lanfear *et al.* 2016),
204 and then inferred phylogenies with RAxML v8.2.4 (Stamatakis 2006). Second, we implemented a species
tree approach. RAxML generates fully bifurcating gene trees even if some nodes have no support. Using
206 the di2multi function in the R package `ape` (Paradis *et al.* 2004), we first collapsed all such nodes in the
gene trees; these nodes have branch lengths $<1e-5$. We then used these gene trees to infer species tree
208 using ASTRAL v4.10.7 (Mirarab & Warnow 2015). ASTRAL only infers tree topology, so to infer branch
lengths, we used BEAST v2.4.5 (Bouckaert *et al.* 2014). To ensure reasonable run times, we randomly
210 subsampled the data sets to 100 loci each and ran 5 independent samples. We did not set fossil or
mutation rate priors as we were interested primarily in comparing relative branch lengths. Because we
212 were only interested in inferring branch lengths, we fixed the topology to the ASTRAL tree by turning off
the subtree-slide, Wilson-Balding and narrow and wide exchange operators. We used an uncorrelated
214 relaxed clock across branches and ran each locus set for 100e6 steps with a 20% burn-in. Trees were
visualized and compared using the R packages `ggtree` and `treescape` (Jombart *et al.* 2015; Yu *et al.*
216 2016). Third, because the SqCL set does not target any mitochondrial loci, we used the program MITObim
v1.8 (Hahn *et al.* 2013) to reconstruct partial to whole mitochondrial genomes from by-catch reads. For
218 each individual, we identified their closest phylogenetic relative from the 271 squamates that have
publically available mitochondrial genomes and used this genome as the seed genome. We then
220 generated a concatenated alignment of the mitochondrial gene sequences and used RAxML to infer the
mitochondrial gene tree.

222

Assessing informativeness for population genomics

224 To facilitate population genetic analyses, we aligned trimmed reads from each individual to its PRG
using bwa v0.7.12 (Li 2013), fixed mate-pair information using samtools v1.3.1 (Li *et al.* 2009), marked
226 duplicate read pairs using picard v2.4.1 (accessed from <https://broadinstitute.github.io/picard/>), and
identified and realigned indels using GATK v3.6 (McKenna *et al.* 2010). We then called a raw set of
228 variants across all individuals in a species using GATK in UnifiedGenotyper mode, filtered the variants
to retain only high-quality variants occurring at sites $\geq 10\times$, and used this filtered variant set to perform
230 base quality score recalibration of the read alignment files. We used GATK's UnifiedGenotyper to call
both non-variant and variants from these recalibrated alignment files and filtered the variants to remove
232 low-quality sites and to set genotypes to missing where coverage was $< 10\times$. Finally, we used GATK's
ReadBackedPhasing to phase variants. The resulting variants were used to infer nucleotide diversity (π ;
234 (Tajima 1983)) and F_{ST} (Reich *et al.* 2009).

236 **Results and Discussion**

Data Quality

238 The data collected were of high-quality and confirmed the efficacy of the SqCL probe set. Of the 5462
targets, only 150 targets failed (7 AHE loci, 140 UCE loci, and 2 genes); we define failed loci as those that
240 were recovered at $< 10\times$ coverage for all individuals (Fig. 1). In total, we were able to generate an average
of 4.29 Mb across 4709 loci of sequence data per individual, of which an average of 69.8% was sequenced
242 to high coverage ($> 10\times$). We were able to assemble most targets in most individuals, leading to a fairly
complete data set particularly for AHE loci (Fig. S2). Because missing data can often complicate
244 phylogenomic inference (Hosner *et al.* 2016; Wiens 2003), researchers using this locus set should be able to
restrict analyses to just well-sampled loci and should still have sufficient data to power most
246 phylogenetic analyses.

248 The 5312 captured targets are distributed across the nuclear genome and across all chromosomes in
Anolis carolinensis (Fig. S3). This dispersed genomic distribution makes it likely that these loci are
250 independently evolving, as assumed in many population genomic and phylogenomic analyses (Brito &
Edwards 2009). For one individual (here, we chose *A. brasiliensis* because it is closely related to the
252 squamate species for which we have the best annotated genome, *A. carolinensis*), we determined the
percentage of coding loci in the capture dataset. Of the 5.90 Mb of assembled sequence for *A. brasiliensis*,

254 5.31 Mb could be aligned to the *A. carolinensis* genome, of which 825 Kb (15.5%) spanned exons and 2.32
Mb (43.8%) fell within gene coordinates. Because only a fraction of the assembled sequence is coding,
256 these loci are not appropriate for researchers interested in some molecular evolutionary questions (*i.e.*,
looking at substitution rates for non-synonymous vs. synonymous sites) although other questions (*i.e.*,
258 levels of heterozygosity in natural populations) can still be addressed. In all subsequent analyses, we
analyze both population genetic and phylogenetic inference across all sequence.

260
We calculated several other quality metrics, including capture efficiency (or, the proportion of sequenced
262 reads that map onto targeted loci), the number of total loci recovered, mean locus length, mean coverage
across loci, and percentage of duplicate reads (Fig. 2, Fig. S4, Fig. S5, Table S1). In general, we see good
264 results for all metrics across all of the diversity sampled. Most notably, on average 93% of AHE loci were
captured at an average length of 1556 bp, 82% of genes at 1040 bp, and 86% of UCEs at 841 bp. Our
266 experiment had a relatively high average capture efficiency rate (60.0%) – capture efficiencies reported in
the literature for AHEs and UCEs can range from 10% to 80% (Faircloth *et al.* 2012b; McCormack *et al.*
268 2016; Ruane *et al.* 2015). On average, our locus assemblies were 30% and 70% longer than the total target
length for AHEs and genes, respectively, and these assemblies were of high quality – 80% of our paired
270 reads mapped properly. Snakes generally performed less well than other squamates, particularly for UCE
loci (Fig. 2). This reduced data quality is partially because one of our eight pools performed poorly
272 during the target capture step of the lab experiment. A linear model found that the pool identity best
explained variation in the number of loci assembled across individuals (adjusted $r^2=0.47$; Fig. S6). Pool
274 number and taxonomy are conflated because we pooled individuals by families. However, both pools 1
and 2 consisted solely of species from the family Dipsadidae, and yet, they had markedly different
276 success rates.

278 Probe design also explains some of this variation in capture efficiency across individuals. The probe
design included a gecko and an anole for UCEs and a snake and an anole for AHEs, which we believe led
280 to geckos' hybridization with UCEs outcompeting their hybridization with AHEs and vice versa for the
snakes. The data confirm this hypothesis; we see geckos have lower AHE recovery compared to
282 squamates as a whole but see no performance reduction for UCEs, and snakes have lower UCE recovery
compared to squamates as a whole but have no performance reduction for AHEs (Fig. 2). As such, we
284 suggest future users use a modified version of this initial probe set (SqCL v2; available at
github.com/singhal/SqCL), in which we include, for 96% of loci, a representative sequence from the lizard

286 *Anolis carolinensis* and the snake *Python molurus*. The remaining 4% have poor matches to either the *A.*
288 *carolinensis* or the *P. molurus* genomes, so we instead use sequence from *Gallus gallus*, one of seven snake
species, or the lizard *Ophisaurus gracilis*.

290 Perhaps the biggest area for improvement is to increase library complexity. Library complexity measures
292 how many of the reads in a library share identical start sites; lower complexity libraries lead to more
sequenced reads being exact duplicates of existing reads. Anywhere from 22% to 54% of our reads were
294 marked as duplicates via computational methods, and duplication rates were correlated to library
complexity ($r = -0.349$, $p = 0.009$). Our libraries should be low complexity because we targeted a subset of
296 the genome, but we see variance around that expectation – libraries in this experiment achieve saturation
at different sequencing depths (Fig. S7). Improving library complexity, both by using higher quality
298 DNA, increasing conversion rates during library generation, and increasing capture efficiency – allowing
us to reduce the number of PCR cycles used to amplify libraries – should make these experiments more
efficacious, ensuring that more reads sequenced are unique and can be retained for downstream analyses.

300 Standard quality metrics for target capture experiments have not yet been reported for either AHE or
302 UCE loci, such as the correlation in coverage across loci across individuals and the variance in coverage
across loci within an individual. Another standard measure, sensitivity, or the percentage of bases of the
304 original target that are at least covered by one read, is less relevant to report here given that UCEs are
designed to capture loci much longer than the original probe sequence. These metrics are particularly
306 useful when target capture loci are used for population genomics, because variant calling quality is
sensitive to sequencing depth (Nielsen *et al.* 2011). If coverage is uneven across loci and across probes, it
308 can lead to sparse data matrices. Thus, in an ideal capture set, variance in coverage across loci within an
individual would be minimal. Although we expect some probes will work better because of their GC-
310 content, melting temperature, and divergence across loci, minimizing variance helps ensure a more
complete data matrix across loci and individuals. Concordantly, in an ideal target set, average coverage
312 across loci should be correlated across individuals, reflecting the differential efficacy of targets. Low
correlations indicate that variance across samples is because of experimental error. We report an average
314 coefficient of variation of 1.20 across loci within individuals, with lower values for AHEs (0.86) and genes
(0.98) than UCEs (1.23) (Fig. S8). Coverage across loci and across individuals was correlated at an average
316 $r = 0.373$ (Fig. S9), and we recovered higher correlations for AHEs ($r = 0.48$) and genes ($r = 0.65$) than UCEs
($r = 0.36$).

318

We see both higher coefficients of variation across loci coverage and lower correlation among individuals than has been reported in exome capture experiments (Bi *et al.* 2012; Bragg *et al.* 2015; Hugall *et al.* 2015; Portik *et al.* 2016). Most exome capture experiments are conducted at a much narrower taxonomic scale (i.e., across species diverged tens of millions of years) than the taxonomic scale used here (i.e., hundreds of millions of years). This increased variance could simply reflect the increased divergence between the probes and the target genomic sequences. To test this hypothesis, we fit a linear model for which factors best predict how well a given locus worked across all individuals, including factors such as the average divergence of the probe sequence across the species considered, the number of probes used for that species, the GC and repeat content of the probes and the loci themselves, and the type of locus (i.e., AHE, UCE, or gene). Our best-fitting single-variable model showed that more divergent probes lead to lower rates of locus recovery (Fig. S10). To ameliorate these effects, future work could reconstruct the ancestral sequence for a given probe across the species of interest and include this sequence in probe sets. A similar approach allowed researchers to target exome data successfully across 250 million years of evolution in the invertebrate class Ophiuroidea (Hugall *et al.* 2015). Making these improvements would indubitably help, but our linear model explains a relatively small portion of the variance ($r^2=0.09$, $p<0.001$; Fig. S10). Future work should consider how we can reduce variance in assembly success and coverage across loci to improve the completeness of our data sets.

336

Data Informativeness for Phylogenetics

Because previous work has clearly shown the utility of AHE and UCE markers for phylogenetic inference (Faircloth *et al.* 2012b; Lemmon *et al.* 2012), we focus our discussion on how marker type influences phylogenetic inference. First, although the probes targeting UCE loci are much more conserved than the probes targeting AHE loci, the sequence divergence of the loci themselves are comparable across locus types (Fig. S1). Further, these loci show broad distributions in how variable they are across sampled individuals (Fig. S1). Because the evolutionary rates of these loci vary, these loci should be able to resolve both broad and shallow radiations. In fact, as others have found, both locus types contain many variable sites across both broad and more shallow radiations (Fig. 3) – the average AHE, gene and UCE locus contains 0.44, 0.39, and 0.43 variable sites/bp across the broad array of squamates sampled. However, where these variable sites occur across loci varies by locus type. AHEs exhibit a fairly uniform density of variable sites across the length of the locus and, as reported previously (Faircloth *et al.* 2012b), UCEs show a U-shaped pattern, with the density of variable sites increasing away from the locus center. Further,

348

350 UCEs show a decline in variable site density at loci ends. Because assembled locus length varies across
individuals, the percent of missing characters at any given column of an alignment increases towards the
352 alignment ends (Fig. S11). This pattern underscores the importance of trimming alignments to remove
regions with high density of missing data (Lemmon & Lemmon 2013).

354 We then explored how these alignments – and their resulting gene trees – differ across several metrics
356 that empirical data suggest can influence phylogenetic inference. First, we inferred phylogenetic
informativeness (PI) (Townsend 2007). PI profiles for the three marker types are comparable (Fig. S12),
358 and all markers are able to resolve deep relationships. None of the marker types shows appreciable
declines after they reach their maximum informativeness, unlike what is typically seen in more quickly
360 evolving loci, like mitochondrial genes (Dornburg *et al.* 2014). As such, all three marker types should be
useful for phylogenetic inference. Second, we measured tree certainty as measured by Salichos and Rokas
362 (2013). Our results showed that AHE and gene markers have greater tree certainty than UCE markers
(Fig. S13). This difference in part reflects a trade-off between locus length and tree certainty. Longer loci
364 (like AHE loci) tend to result in better resolved gene trees (Arcila *et al.* 2017; Blom *et al.* 2016), though
they are also more likely to contain recombination events that violate most gene tree inference methods.
366 Finally, we characterized how well these loci fit to a clock-like model for molecular evolution, finding
UCE markers appear to be more clock-like than AHEs (Fig. S14). No one marker type emerges as superior
368 to the others across these metrics. Rather, these loci exhibit significant variation across these metrics,
suggesting that sampling more loci will allow users to carefully filter loci as required by their desired
370 analysis.

372 We then inferred phylogenies for these loci using concatenated and species tree approaches. We do not
discuss our concatenated results (Fig. S15) because concatenation (particularly with phylogenomic data)
374 can often converge on the wrong tree with high support (Kubatko & Degnan 2007). Our species tree
analyses recover largely similar topologies across the three marker types (Fig. 4), particularly between the
376 topologies inferred with AHE and UCE markers. Across all comparisons, the nodes that disagree also
tend to have low support. Additionally, our divergence dating analyses across marker types showed that
378 branch length estimates were highly correlated across inferred trees (Fig. S16), a result that is
unsurprising given that raw pairwise genetic divergences between tips are also highly correlated (Fig.
380 S17). In sum, phylogenetic inference – both in topology and branch length estimation – is robust to
marker type. Further, while the increased data content of both AHE and UCE marker sets allow us to

382 resolve some tricky nodes in the phylogeny, some nodes remain poorly resolved. Future work will
explore (1) filtering loci to see if filtered data sets lead to more resolved trees (Blom *et al.* 2016; Doyle *et al.*
384 2015; Salichos & Rokas 2013), (2) using methods that resolve tricky nodes by constraining the topology
space explored (Arcila *et al.* 2017), or (3) accounting for phylogenetic uncertainty in any tree-based
386 analyses. Further explorations into the causes for this topological discordance are beyond the scope of
this study.

388 *Recovering the mtDNA genome*

390 Including mitochondrial targets in the probe set is not recommended. Because of the difference in copy
number between mitochondrial and nuclear genomes in vertebrates, mitochondrial DNA generally
392 outcompetes nuclear DNA for binding, leading to far greater coverage of the mtDNA genome than the
nDNA genome (Bi *et al.* 2012). However, mtDNA is the traditional workhorse for phylogenetics, and
394 genealogical discordance between mtDNA and nuclear data is often used as a marker for introgression
between taxa (Toews & Brelsford 2012). As such, we evaluated our ability to recover mtDNA from these
396 taxa. We were able to assemble portions of the mtDNA genome for 55 of our 56 individuals, although two
of these individuals had no sequence data for any of the 13 mtDNA polypeptide genes. Of the remaining
398 53 individuals, we recovered 89.2% of the total length of the mitochondrial genome. We used these data
to infer a mtDNA gene tree (Fig. S21), which differs from the SqCL-based tree at deeper nodes though it
400 recovers many of the same species relationships within families. The quality of our mtDNA assembly
(here, measured by the portion of sequence that was recovered) was negatively correlated with capture
402 efficiency ($r=0.453$, $p < 0.001$). In individuals with more reads mapping on target, there are fewer reads
randomly sequenced from the mtDNA genome and, thus, recovering a complete mtDNA genome is less
404 likely.

406 *Data Informativeness for Population Genomics*

The utility of AHE and UCE loci for population genomics and phylogeography has already been reported
408 in a number of papers (Brandley *et al.* 2015; Harvey *et al.* 2016; Zarza *et al.* 2016). Here, we further
compare and contrast across patterns of variation across the locus types. Although we expect AHEs,
410 UCEs, and conserved genes to be less variable than other locus types used in population genomics – *i.e.*,
exons or RAD data (Bragg *et al.* 2015; Harvey *et al.* 2016) – we recover sufficient variation across all three
412 locus types to power population genomic analyses (Fig. 5). In particular, summarizing across all data
types, we were able to generate robust estimates of isolation-by-distance slopes for both *C. modesta* and *B.*

414 *moojeni* (Fig. 6), illustrating the utility of these markers to study population-level processes. The average
AHE, gene and UCE locus contains 0.0035 (10-90% distribution: 0.001-0.006), 0.0025 (0.0 – 0.004), and
416 0.0042 (0.001-0.008) segregating sites per bp across *Colobosaura modesta*, the lizard for which we sampled 6
individuals, and 0.0019 (0.0-0.004), 0.0022 (0.0 – 0.007), and 0.0021 (0.0-0.004) segregating sites per bp
418 across *Bothrops moojeni*, the snake for which we sampled 5 individuals. The pattern of SNP density
mimics the pattern of variable site density (Fig. 3, 5). As seen with our phylogenetic results, locus design
420 influences both coverage and patterns of variation across the length of loci (Fig. S18). Despite this,
patterns of both genetic diversity and differentiation were highly correlated across marker types (Fig. S19,
422 Fig. S20). The slope of these relationships generally deviated from unity, which reflects these loci's
different evolutionary histories. Selection, recombination, and their interaction likely influence effective
424 population sizes across these markers differentially (Charlesworth 2009).

426 *Practicality of Approach*

We pooled fewer individuals to a lane than most other target capture experiments, which regularly
428 multiplex 100 individuals to a single lane of sequencing. Thus, we sequenced our libraries to a much
greater depth than is typical. To test how reduced sequencing would affect the quality of the data
430 recovered, we conducted a series of subsampling experiments in which we took the 11 individuals in
Colobosaura modesta and *Bothrops moojeni* and randomly sampled 5e5, 1e6, 1.5e6, and 2e6 paired-reads (for
432 a total of 1e6, 2e6, 3e6, and 4e6 reads). With current sequencing yields on the Illumina HiSeq 2500 v4
sequencing platform of approximately 250 million paired-reads, this represents pooling of approximately
434 500, 250, 166 and 125 individuals to a single sequencing lane. Even with significantly reduced sequencing,
we still assembled a large number of loci for a given individual, with only modest improvements for
436 additional sequencing beyond 2e6 reads (Fig. 7). However, sequencing more reads led to a linear increase
in the number of sites with sufficiently high coverage to call variants (Fig. 7). Researchers interested in
438 population genomic analyses might want to use lower levels of multiplexing than those interested solely
in phylogenomics. This analysis is contingent on both capture efficiency and library complexity, and
440 improving the number of reads mapping on target and / or reducing the library duplication rate will
allow researchers to multiplex even further.

442
Using the SqCL probe set presents additional costs. More probes must be synthesized than if either locus
444 set was used in isolation. In our study, the cost for probes per sample increased from \$25 for solely
capturing UCEs to \$31.25 for the entire SqCL set. Further, sequencing both loci requires further

446 investment in sequencing than sequencing either set alone. However, our subsetting experiment (Fig. 7)
suggests that researchers should still be able to multiplex at similar levels as used in other projects using
448 AHE and UCE loci (Meiklejohn *et al.* 2016; Prum *et al.* 2015), despite the increase in overall target length.
Thus, we anticipate that using SqCL loci will result in only modest increases in cost for a given project,
450 while generating a much more inclusive dataset.

452 *Conclusions*

The AHE and UCE locus sets made an important contribution to the field of biodiversity genomics by
454 allowing researchers to efficiently query homologous loci across a diversity of organisms. However, the
presence of two largely non-overlapping locus sets has created an unfortunate divide, in that many
456 research groups have invested in either AHEs or UCEs for their clade of interest. This lack of overlap will
hinder future attempts at synthesis in both population genomics and phylogenetics, limiting the utility of
458 existing datasets. We have provided a simple resolution to this problem by presenting a probe set that
includes AHEs, UCEs, and ~50 additional loci that have served as "workhorse genes" for squamate
460 phylogenetics. Because target capture also often allows us to recover the mitochondrial genome (Fig. S21),
the SqCL probe set thus provides maximal integration with most existing phylogenetic data. These data
462 also allow population and conservation genetics researchers to generate datasets that can ultimately be
integrated into broad-scale comparative analyses. We advocate the use of this integrated probe set for
464 questions currently being addressed with UCEs or AHEs, thus ensuring that future data sets for
squamate reptiles are compatible with much of the existing phylogenetic data generated over the last
466 thirty years. Importantly, both population genomic and phylogenetic inferences are robust across marker
types.

468 Although we refined the AHE, UCE and traditional gene sets for their application to squamate
470 phylogenetics only, our approach can easily be applied to other tetrapod systems and could be used to
create a probe set of general use across the tetrapod phylogeny, thus further supporting the development
472 of community-wide, inclusive locus set for use in phylogenomics and comparative population genomics.
This study took effort to customize these probe sets for squamates; however, published probe sequences
474 could simply be synthesized and applied to tetrapod systems of interest (Faircloth *et al.* 2012b; Lemmon *et al.*
et al. 2012). AHE probes tend to diverge more quickly across phylogenetic distance than UCE probes (Fig.
476 S1A). To ensure efficient capture, researchers should ideally synthesize AHE probes specific to their
broad clade of interest (*i.e.*, amphibians, reptiles or mammals).

478

Author Contributions

480 SS was involved with project design, lab work, data analysis, and paper writing, MG helped with lab
work, GC contributed samples, and DLR helped design the project. All authors read and approved the
482 manuscript.

484 Data Accessibility

- Raw reads: associated with BioProject PRJNA382381
- 486 • Probe sequences for SqCL v1 and v2 available at <https://github.com/singhal/SqCL>
- Assemblies for all 44 species available at 10.5061/dryad.r0q02
- 488 • VCF files for the two species for which we called variants available at 10.5061/dryad.r0q02
- Scripts used in probe design and data analysis, along with README, available at
490 https://github.com/singhal/SqCL_analysis

492 Acknowledgements

For logistical support, we thank Alison Devault & Jake Enk at MYcroarray, Marcella Baiz, Robbin
494 Murrell, Gabriel Costa, Izabella Paim da Silva, and Fabricius Maia Domingos. Brant Faircloth and Mozes
Blom provided useful advice. This work is funded by a grant from the David and Lucile Packard
496 Foundation to DLR and a NSF Postdoctoral Fellowship in Research Biology to SS. GRC thanks
Coordenação de Apoio à Formação de Pessoal de Nível Superior – CAPES, Conselho Nacional do
498 Desenvolvimento Científico e Tecnológico – CNPq, and Fundação de Apoio à Pesquisa do Distrito
Federal – FAPDF for financial support.

500

Figures

502 Figure 1: A phylogeny for the 44 species used to test the SqCL set along with a matrix, in which each
column represents one of the 5,462 loci targeted. Green columns indicate anchored hybrid enrichment
504 (AHE) loci, purple ultraconserved element (UCE) loci, orange traditional phylogenetic genes, and gray
indicates missing data. Loci are arrayed in order of most to least complete. The phylogeny topology was
506 inferred using ASTRAL-II and BEAST2 for 2,815 loci that were 95% complete across all taxa and rooted
with *Gallus gallus* (not shown). Gray dots mark nodes with >0.95 local posterior probability.

508 - `heat_map.pdf`

510 Figure 2: Several metrics of data quality summarized across the three types of loci in the SqCL set
(anchored hybrid enrichment loci: AHE; traditional phylogenetic genes: gene; ultraconserved elements:
512 UCE) across the 16 squamate families sampled. Results show the SqCL set works well across taxa that last
shared a common ancestor more than 200 million years ago. Data quality metrics are: the percent of loci
514 targeted that were recovered, the mean locus length, and mean coverage across the locus. Shown are
median values and the 95% percentile range across individuals sampled for that family. Not all points are
516 shown with confidence intervals because we only sampled one species in some families. A version of this
figure showing patterns across additional metrics is shown in Fig. S3.

518 - `data_quality_by_loci.pdf`

520 Figure 3: Density of variable sites in multi-species alignments for the three types of loci in the SqCL set
(anchored hybrid enrichment: AHE; traditional phylogenetic genes: gene, ultraconserved element: UCE).
522 Black dots indicate variable site density for the 44 squamate taxa sequenced in this study; gray dots
variable site density for the 30 snake taxa sequenced. The squamates span 200 million years of
524 divergence, and the snakes span 120 million years of divergence (Zheng & Wiens 2016). The frequency of
variable sites differs between the two comparisons, reflecting the difference in phylogenetic depth.
526 Different loci types exhibit different variable density patterns across the length of the loci, which is both a
function of locus design and variation in levels of missingness across the locus alignment.

528 - `PIC_density_types.pdf`

530 Figure 4: Species trees inferred using anchored hybrid enrichment (AHE), traditional phylogenetic (gene)
loci, and ultraconserved element (UCE) loci that were 95% complete across the 44 species with ASTRAL-
532 II. Trees were rooted with *Gallus gallus* (not shown). Gray boxes mark clades that exhibit unstable
topologies across marker sets, and the matrix shows normalized Robinson-Foulds distances between
534 trees. Nodes with <0.95 local posterior probability are shown in red. Topologies are largely concordant
across marker sets, and conflicting nodes generally have low support.

536 - `speciestrees_across_markers.pdf`

538 Figure 5: Density of single nucleotide polymorphisms (SNPs) for the three types of loci in the SqCL set
(anchored hybrid enrichment: AHE; traditional phylogenetic genes: gene, ultraconserved element: UCE).
540 Black dots indicate single nucleotide polymorphism (SNP) density across 6 individuals of the lizard
Colobosaura modesta; gray dots SNP density across 5 individuals of the snake *Bothrops moojeni*. Different

542 loci types exhibit different SNP densities across the length of the loci, which is a both a function of locus
design and average sequencing coverage across the loci length.

544 - `variant_density_types.pdf`

546 Figure 6: Isolation-by-distance estimates for *Colobosaura modesta* and *Bothrops moojeni*. Each point rep-
resents a pairwise comparison between two individuals. F_{ST} estimates are based on an average of 37K
548 variant sites. The two species have very different isolation-by-distance relationships, illustrating the
power of SqCL markers to address questions about population-level variation.

550 - `IBD.pdf`

552 Figure 7: Results of an *in silico* experiment testing the effect of reducing sequencing depth on the number
of loci assembled and the number of sites with $\geq 10\times$ coverage, or those sites at we call single nucleotide
554 polymorphisms (SNPs). For the 5 individuals in the species *Bothrops moojeni* (shown in gray) and the 6
individuals in *Colobosaura modesta* (in white), we used SeqTK (<https://github.com/lh3/seqtk>) to randomly
556 subsample 5e5, 1e6, 1.5e6, and 2e6 paired-reads (for a total of 1e6, 2e6, 3e6, and 4e6 reads) and analyzed
the data using our bioinformatics pipeline. In this study, we sequenced an average of 3.5e6 paired reads
558 for these 11 individuals. We could reduce sequencing depth by 70% and still recover 86.7% of the loci.
Decreasing sequencing depth, however, does decrease the number of sites recovered at high coverage.

560 - `subset_experiment.pdf`

562 References

Aljanabi SM, Martinez I (1997) Universal and rapid salt-extraction of high quality genomic DNA for PCR-
564 based techniques. *Nucleic acids research* **25**, 4692-4693.

Arcila D, Orti G, Vari R, *et al.* (2017) Genome-wide interrogation advances resolution of recalcitrant
566 groups in the tree of life. *Nature Ecology & Evolution* **1**, 0020.

Bi K, Vanderpool D, Singhal S, *et al.* (2012) Transcriptome-based exon capture enables highly cost-
568 effective comparative genomic data collection at moderate evolutionary scales. *BMC genomics* **13**,
1.

570 Blom MP, Bragg JG, Potter S, Moritz C (2016) Accounting for uncertainty in gene tree estimation:
summary-coalescent species tree inference in a challenging radiation of Australian lizards.

572 *Systematic biology*, syw089.

- 574 Bolger AM, Lohse M, Usadel B (2014) Trimmomatic: a flexible trimmer for Illumina sequence data.
Bioinformatics **30**, 2114-2120.
- 576 Bouckaert R, Heled J, Kühnert D, *et al.* (2014) BEAST 2: a software platform for Bayesian evolutionary
analysis. *PLoS Comput Biol* **10**, e1003537.
- 578 Bragg JG, Potter S, Bi K, Moritz C (2015) Exon capture phylogenomics: efficacy across scales of
divergence. *Molecular ecology resources*.
- 580 Brandley MC, Bragg JG, Singhal S, *et al.* (2015) Evaluating the performance of anchored hybrid
enrichment at the tips of the tree of life: a phylogenetic analysis of Australian Eugongylus group
scincid lizards. *BMC Evolutionary Biology* **15**, 1.
- 582 Brito PH, Edwards SV (2009) Multilocus phylogeography and phylogenetics using sequence-based
markers. *Genetica* **135**, 439-455.
- 584 Budenhagen C, Lemmon AR, Lemmon EM, *et al.* (2016) Anchored Phylogenomics of Angiosperms I:
Assessing the Robustness of Phylogenetic Estimates. *bioRxiv*, 086298.
- 586 Camacho C, Coulouris G, Avagyan V, *et al.* (2009) BLAST+: architecture and applications. *BMC
bioinformatics* **10**, 1.
- 588 Castresana J (2000) Selection of conserved blocks from multiple alignments for their use in phylogenetic
analysis. *Molecular Biology and Evolution* **17**, 540-552.
- 590 Charlesworth B (2009) Effective population size and patterns of molecular evolution and variation. *Nature
Reviews Genetics* **10**, 195-205.
- 592 Colli GR, Bastos RP, Araujo AF, Oliveira P, Marquis R (2002) The character and dynamics of the Cerrado
herpetofauna. *The Cerrados of Brazil: ecology and natural history of a Neotropical savanna*, 223-241.
- 594 Cosart T, Beja-Pereira A, Chen S, *et al.* (2011) Exome-wide DNA capture and next generation sequencing
in domestic and wild species. *BMC genomics* **12**, 347.
- 596 Crawford NG, Faircloth BC, McCormack JE, *et al.* (2012) More than 1000 ultraconserved elements provide
evidence that turtles are the sister group of archosaurs. *Biology Letters* **8**, 783-786.
- 598 Dornburg A, Townsend JP, Friedman M, Near TJ (2014) Phylogenetic informativeness reconciles ray-
finned fish molecular divergence times. *BMC evolutionary biology* **14**, 169.
- 600 Doyle VP, Young RE, Naylor GJ, Brown JM (2015) Can we identify genes with increased phylogenetic
reliability? *Systematic biology* **64**, 824-837.
- 602 Faircloth BC (2015) PHYLUCE is a software package for the analysis of conserved genomic loci.
Bioinformatics, btv646.

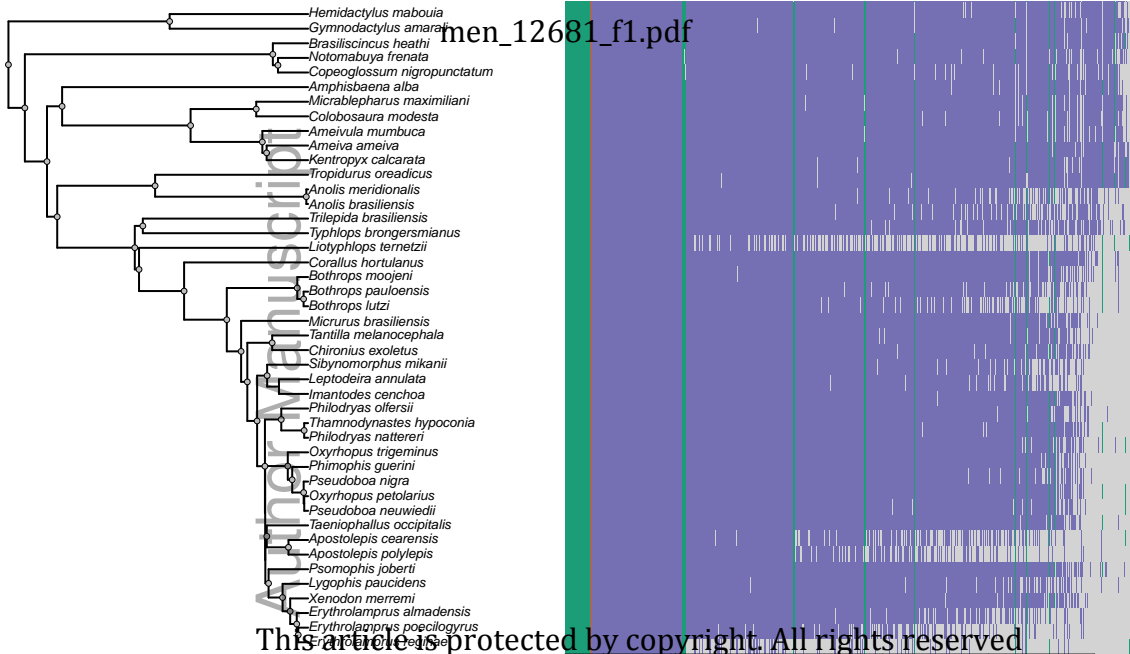
- 604 Faircloth BC, Branstetter MG, White ND, Brady SG (2015) Target enrichment of ultraconserved elements
from arthropods provides a genomic perspective on relationships among Hymenoptera.
606 *Molecular ecology resources* **15**, 489-501.
- Faircloth BC, Chang J, Alfaro ME (2012a) TAPIR enables high-throughput estimation and comparison of
608 phylogenetic informativeness using locus-specific substitution models. *arXiv preprint*
arXiv:1202.1215.
- 610 Faircloth BC, McCormack JE, Crawford NG, *et al.* (2012b) Ultraconserved elements anchor thousands of
genetic markers spanning multiple evolutionary timescales. *Systematic biology*, sys004.
- 612 Giarla TC, Esselstyn JA (2015) The challenges of resolving a rapid, recent radiation: empirical and
simulated phylogenomics of Philippine shrews. *Systematic biology* **64**, 727-740.
- 614 Gilbert PS, Chang J, Pan C, *et al.* (2015) Genome-wide ultraconserved elements exhibit higher
phylogenetic informativeness than traditional gene markers in percomorph fishes. *Molecular*
616 *phylogenetics and evolution* **92**, 140-146.
- Grabherr MG, Haas BJ, Yassour M, *et al.* (2011) Full-length transcriptome assembly from RNA-Seq data
618 without a reference genome. *Nature biotechnology* **29**, 644-652.
- Guindon S, Dufayard J-F, Lefort V, *et al.* (2010) New algorithms and methods to estimate maximum-
620 likelihood phylogenies: assessing the performance of PhyML 3.0. *Systematic biology* **59**, 307-321.
- Hahn C, Bachmann L, Chevreur B (2013) Reconstructing mitochondrial genomes directly from genomic
622 next-generation sequencing reads—a baiting and iterative mapping approach. *Nucleic acids*
research, gkt371.
- 624 Harmon LJ, Weir JT, Brock CD, Glor RE, Challenger W (2008) GEIGER: investigating evolutionary
radiations. *Bioinformatics* **24**, 129-131.
- 626 Harvey MG, Smith BT, Glenn TC, Faircloth BC, Brumfield RT (2016) Sequence Capture Versus Restriction
Site Associated DNA Sequencing for Shallow Systematics. *Systematic biology*, syw036.
- 628 Hillier LW, Miller W, Birney E, *et al.* (2004) Sequence and comparative analysis of the chicken genome
provide unique perspectives on vertebrate evolution. *Nature* **432**, 695-716.
- 630 Hosner PA, Faircloth BC, Glenn TC, Braun EL, Kimball RT (2016) Avoiding missing data biases in
phylogenomic inference: An empirical study in the landfowl (Aves: Galliformes). *Molecular*
632 *biology and evolution* **33**, 1110-1125.
- Huang X, Madan A (1999) CAP3: A DNA sequence assembly program. *Genome research* **9**, 868-877.
- 634 Hugall AF, O'Hara TD, Hunjan S, Nilsen R, Moussalli A (2015) An exon-capture system for the entire
class Ophiuroidea. *Molecular biology and evolution*, msv216.

- 636 Jetz W, Thomas G, Joy J, Hartmann K, Mooers A (2012) The global diversity of birds in space and time.
Nature **491**, 444-448.
- 638 Jombart T, Kendall M, Almagro-Garcia J, Colijn C (2015) treescape: Statistical Exploration of Landscapes
of Phylogenetic Trees. *R package version 1*, 15.
- 640 Jones MR, Good JM (2016) Targeted capture in evolutionary and ecological genomics. *Molecular ecology*
25, 185-202.
- 642 Katoh K, Standley DM (2013) MAFFT multiple sequence alignment software version 7: improvements in
performance and usability. *Molecular biology and evolution* **30**, 772-780.
- 644 Kent WJ (2002) BLAT—the BLAST-like alignment tool. *Genome research* **12**, 656-664.
- Kubatko LS, Degnan JH (2007) Inconsistency of phylogenetic estimates from concatenated data under
646 coalescence. *Systematic biology* **56**, 17-24.
- Lanfear R, Frandsen PB, Wright AM, Senfeld T, Calcott B (2016) PartitionFinder 2: new methods for
648 selecting partitioned models of evolution for molecular and morphological phylogenetic
analyses. *Molecular Biology and Evolution*, msw260.
- 650 Lemmon AR, Emme SA, Lemmon EM (2012) Anchored hybrid enrichment for massively high-
throughput phylogenomics. *Systematic biology*, sys049.
- 652 Lemmon EM, Lemmon AR (2013) High-throughput genomic data in systematics and phylogenetics.
Annual Review of Ecology, Evolution, and Systematics **44**, 99-121.
- 654 Li H (2013) Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM. *arXiv*,
1303.3997.
- 656 Li H, Handsaker B, Wysoker A, et al. (2009) The sequence alignment/map format and SAMtools.
Bioinformatics **25**, 2078-2079.
- 658 McCormack JE, Tsai WL, Faircloth BC (2016) Sequence capture of ultraconserved elements from bird
museum specimens. *Molecular ecology resources*.
- 660 McKenna A, Hanna M, Banks E, et al. (2010) The Genome Analysis Toolkit: a MapReduce framework for
analyzing next-generation DNA sequencing data. *Genome research* **20**, 1297-1303.
- 662 Meiklejohn KA, Faircloth BC, Glenn TC, Kimball RT, Braun EL (2016) Analysis of a Rapid Evolutionary
Radiation Using Ultraconserved Elements: Evidence for a Bias in Some Multispecies Coalescent
664 Methods. *Systematic biology*, syw014.
- Mirarab S, Warnow T (2015) ASTRAL-II: coalescent-based species tree estimation with many hundreds of
666 taxa and thousands of genes. *Bioinformatics* **31**, i44-i52.

- 668 Nielsen R, Paul JS, Albrechtsen A, Song YS (2011) Genotype and SNP calling from next-generation
sequencing data. *Nature Reviews Genetics* **12**, 443-451.
- 670 Paradis E, Claude J, Strimmer K (2004) APE: analyses of phylogenetics and evolution in R language.
Bioinformatics **20**, 289-290.
- 672 Portik DM, Smith LL, Bi K (2016) An evaluation of transcriptome - based exon capture for frog
phylogenomics across multiple scales of divergence (Class: Amphibia, Order: Anura). *Molecular
ecology resources*.
- 674 Prum RO, Berv JS, Dornburg A, et al. (2015) A comprehensive phylogeny of birds (Aves) using targeted
next-generation DNA sequencing. *Nature*.
- 676 Pyron RA, Burbrink FT, Wiens JJ (2013) A phylogeny and revised classification of Squamata, including
4161 species of lizards and snakes. *BMC Evolutionary Biology* **13**, 1.
- 678 Reich D, Thangaraj K, Patterson N, Price AL, Singh L (2009) Reconstructing Indian population history.
Nature **461**, 489-494.
- 680 Ruane S, Raxworthy C, Lemmon A, Moriarty Lemmon E, Burbrink F (2015) Comparing species tree
estimation with large anchored phylogenomic and small Sanger-sequenced molecular datasets:
682 an empirical study on Malagasy pseudoxyrhophiine snakes. *BMC Evolutionary Biology* **15**,
doi:10.1186/s12862-12015-10503-12861.
- 684 Salichos L, Rokas A (2013) Inferring ancient divergences requires genes with strong phylogenetic signals.
Nature **497**, 327-331.
- 686 Salichos L, Stamatakis A, Rokas A (2014) Novel information theory-based measures for quantifying
incongruence among phylogenetic trees. *Molecular Biology and Evolution*, msu061.
- 688 Smith BT, Harvey MG, Faircloth BC, Glenn TC, Brumfield RT (2014) Target capture and massively
parallel sequencing of ultraconserved elements (UCEs) for comparative studies at shallow
690 evolutionary time scales. *Systematic biology* **63**, 83-95.
- Springer MS, Gatesy J (2016) The gene tree delusion. *Molecular phylogenetics and evolution* **94**, 1-33.
- 692 Stamatakis A (2006) RAxML-VI-HPC: maximum likelihood-based phylogenetic analyses with thousands
of taxa and mixed models. *Bioinformatics* **22**, 2688-2690.
- 694 Swofford DL (2003) PAUP*. Phylogenetic analysis using parsimony (* and other methods). Version 4.
- Tajima F (1983) Evolutionary relationship of DNA sequences in finite populations. *Genetics* **105**, 437-460.
- 696 Therkildsen NO, Palumbi SR (2016) Practical low - coverage genomewide sequencing of hundreds of
individually barcoded samples for population and evolutionary genomics in nonmodel species.
698 *Molecular ecology resources*.

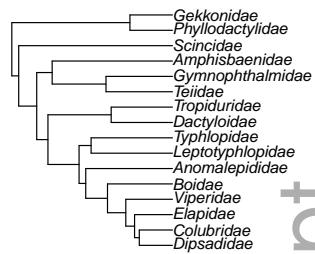
- Toews DP, Brelsford A (2012) The biogeography of mitochondrial and nuclear discordance in animals.
700 *Molecular ecology* **21**, 3907-3930.
- Townsend JP (2007) Profiling phylogenetic informativeness. *Systematic biology* **56**, 222-231.
- 702 Warnow T (2011) Concatenation analyses in the presence of incomplete lineage sorting. *PLoS currents* **7**.
- Wiens JJ (2003) Missing data, incomplete taxa, and phylogenetic accuracy. *Systematic biology* **52**, 528-538.
- 704 Wiens JJ, Hutter CR, Mulcahy DG, *et al.* (2012) Resolving the phylogeny of lizards and snakes (Squamata)
with extensive sampling of genes and species. *Biology Letters* **8**, 1043-1046.
- 706 Yu G, Smith DK, Zhu H, Guan Y, Lam TTY (2016) ggtree: an r package for visualization and annotation of
phylogenetic trees with their covariates and other associated data. *Methods in Ecology and*
708 *Evolution*.
- Zarza E, Faircloth BC, Tsai WL, *et al.* (2016) Hidden histories of gene flow in highland birds revealed with
710 genomic markers. *Molecular Ecology* **25**, 5144-5157.
- Zhang J, Kobert K, Flouri T, Stamatakis A (2014) PEAR: a fast and accurate Illumina Paired-End reAd
712 mergeR. *Bioinformatics* **30**, 614-620.
- Zheng Y, Wiens JJ (2016) Combining phylogenomic and supermatrix approaches, and a time-calibrated
714 phylogeny for squamate reptiles (lizards and snakes) based on 52 genes and 4162 species.
Molecular phylogenetics and evolution **94**, 537-547.
- 716

men_12681_f1.pdf

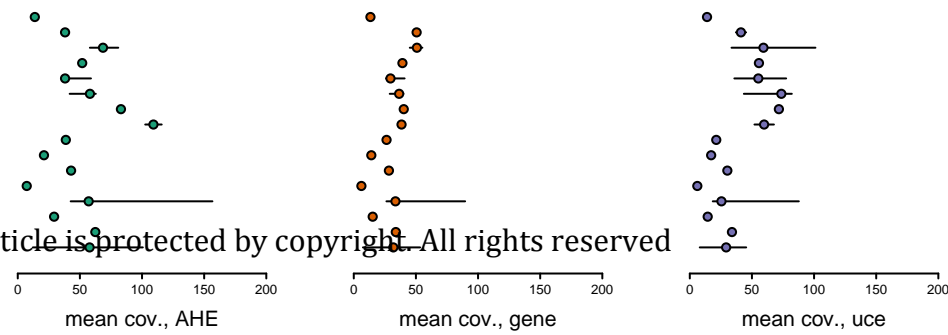
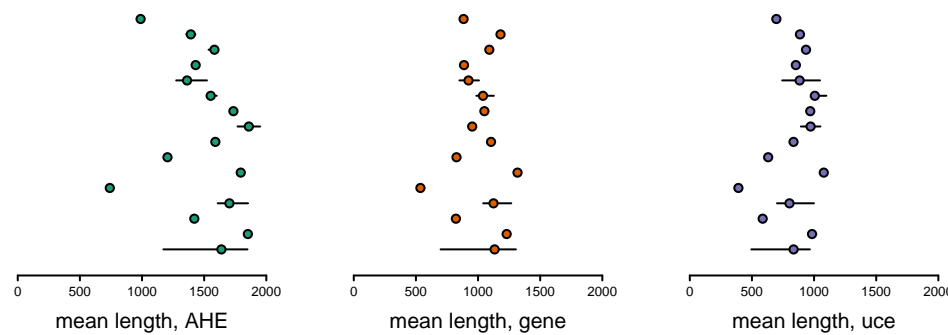
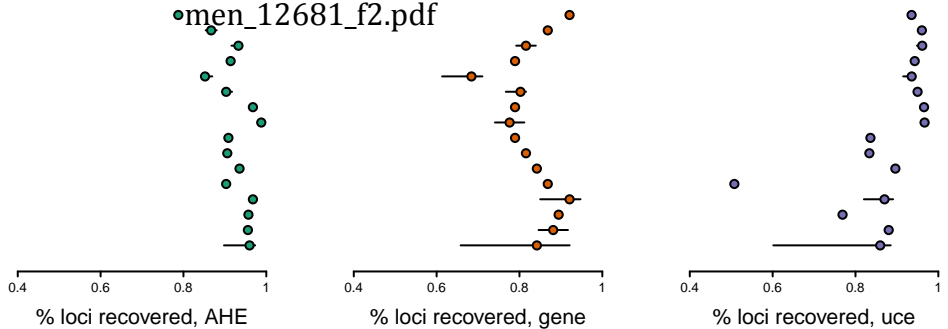


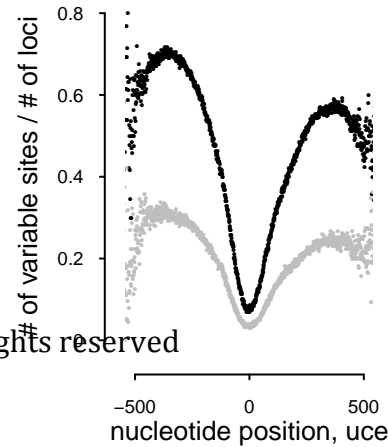
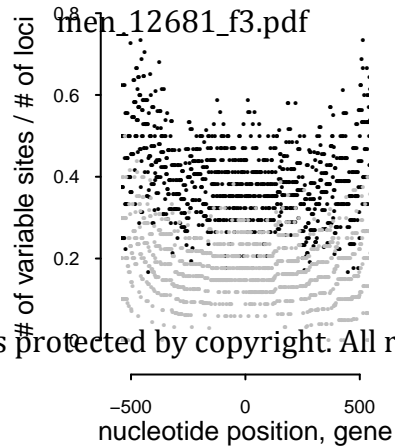
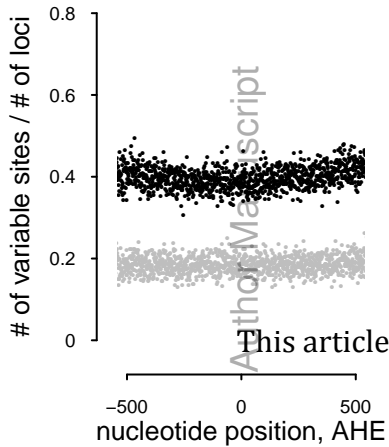
This article is protected by copyright. All rights reserved

1000 2000 3000 4000 5000
number of loci



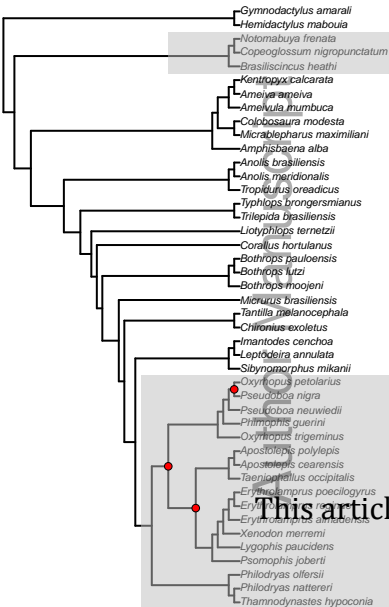
Author Manuscript





This article is protected by copyright. All rights reserved

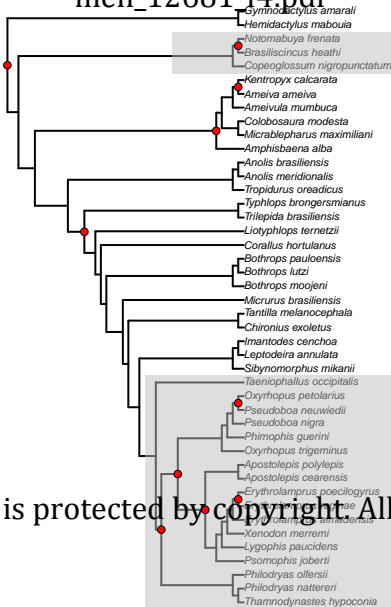
AHE



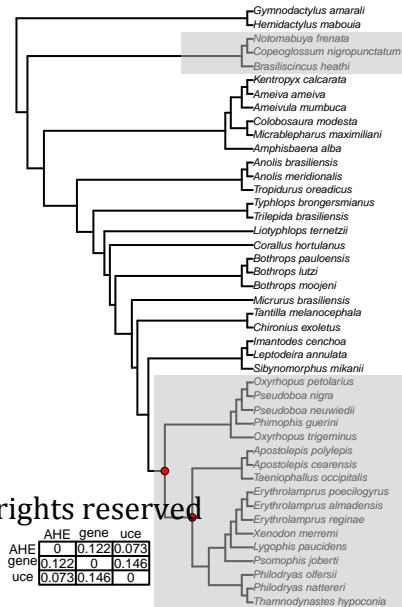
gene

men_12681

f4.pdf

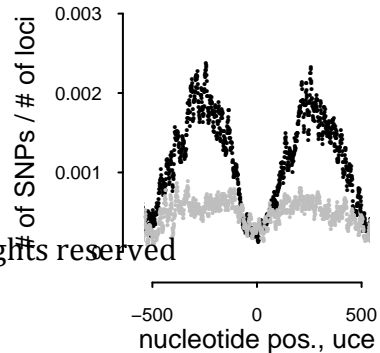
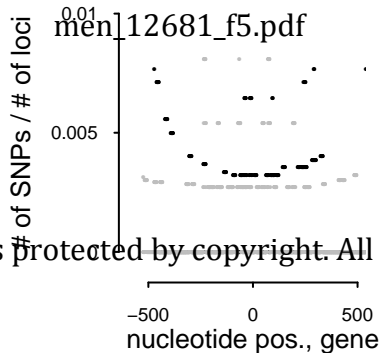
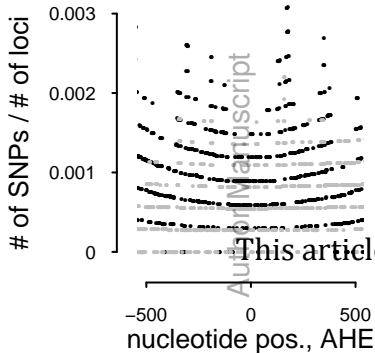


uce



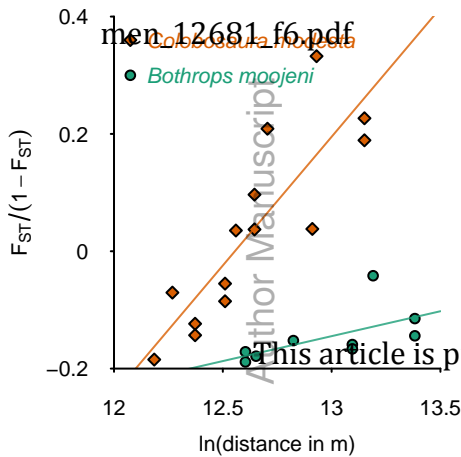
	AHE	gene	uce
AHE	0	0.122	0.073
gene	0.122	0	0.146
uce	0.073	0.146	0

This article is protected by copyright. All rights reserved



This article is protected by copyright. All rights reserved

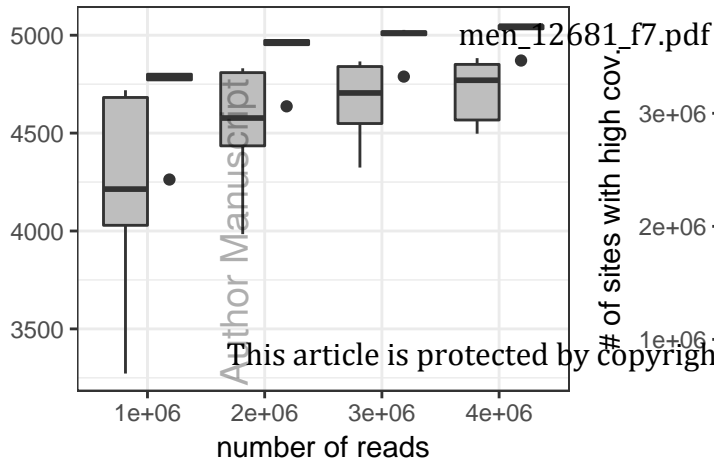
men_12681_f6.pdf



Author Manuscript

This article is p

of loci assembled



of sites with high cov.

