

GigaScience

Single molecule, full-length transcript sequencing provides insight into the extreme metabolism of ruby-throated hummingbird *Archilochus colubris* --Manuscript Draft--

Manuscript Number:	GIGA-D-17-00088R2	
Full Title:	Single molecule, full-length transcript sequencing provides insight into the extreme metabolism of ruby-throated hummingbird <i>Archilochus colubris</i>	
Article Type:	Data Note	
Funding Information:	Human Frontier Science Program (#RGP0062/2016)	Dr. Kenneth C. Welch
	Natural Sciences and Engineering Research Council of Canada (CA) (#386466)	Dr. Kenneth C. Welch
Abstract:	<p>Background Hummingbirds oxidize ingested nectar sugars directly to fuel foraging but cannot sustain this fuel use during fasting periods, such as during the night or during long-distance migratory flights. Instead, fasting hummingbirds switch to oxidizing stored lipids, derived from ingested sugars. The hummingbird liver plays a key role in moderating energy homeostasis and this remarkable capacity for fuel switching. Additionally, liver is the principle location of de novo lipogenesis, which can occur at exceptionally high rates, such as during premigratory fattening. Yet understanding how this tissue and whole organism moderates energy turnover is hampered by a lack of information regarding how relevant enzymes differ in sequence, expression, and regulation.</p> <p>Findings We generated a de novo transcriptome of the hummingbird liver using PacBio full-length cDNA sequencing (Iso-Seq), yielding a total of 8.6Gb of sequencing data, or 2.6M reads from 4 different size fractions. We analyzed data using the SMRTAnalysis v3.1 Iso-Seq pipeline, then clustered isoforms into gene families to generate de novo gene contigs using Cogent. We performed orthology analysis to identify closely related sequences between our transcriptome and other avian and human gene sets. Finally, we closely examined homology of critical lipid metabolism genes between our transcriptome data and avian and human genomes.</p> <p>Conclusions We confirmed high levels of sequence divergence within hummingbird lipogenic enzymes, suggesting a high probability of adaptive divergent function in the hepatic lipogenic pathways. Our results leverage cutting-edge technology and a novel bioinformatics pipeline to provide a first direct look at the transcriptome of this incredible organism.</p>	
Corresponding Author:	Winston Timp Johns Hopkins University Baltimore, Maryland UNITED STATES	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Johns Hopkins University	
Corresponding Author's Secondary Institution:		
First Author:	Rachael E. Workman	
First Author Secondary Information:		
Order of Authors:	Rachael E. Workman	
	Alexander M. Myrka	
	Elizabeth Tseng	
	G. William Wong	

	Kenneth C. Welch
	Winston Timp
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Below please find our revisions and response to the reviewers concerns on the manuscript entitled "Single molecule, full-length transcript sequencing provides insight into the extreme metabolism of ruby-throated hummingbird <i>Archilochus colubris</i>"</p> <p>We thank the editors and reviewers for their thoughtful comments on the revision. We have addressed all of the reviewers' concerns and are submitting a detailed response along with an improved version of the manuscript.</p> <p>>>></p> <p>Reviewer #1:</p> <p>I am satisfied with the changes made in response to suggestions (from me and another reviewer) previously, and pleased with the way the manuscript has shaped with the additional Illumina sequences.</p> <p>However, the paper still has inadequacies that need to be fixed. While, genes that should not have been included have been fixed, others that are bona-fide genes have been left out.</p> <p>Several genes that map to the <i>Calypte anna</i> and <i>Aquila chrysaetos</i> transcripts are missing in the <i>acolubris</i> HQD or cogent CDS. An example is SRR5237173.100014 (953 nt). This is 99% similar to 'Calypte anna zinc finger CCCH-type containing 6 (ZC3H6), transcript variant X1, mRNA' (XM 008499961.1). Another example is SRR5237173.117510 (904 nt), which matches 88% identity to the golden eagle (XM 011575986.1) "vesicle transport through interaction with t-SNAREs 1A (VT11A) transcript variant X2, mRNA'. Its acceptable to miss a few genes, but my analysis shows 6000 genes (corresponding to about 22,000 PacBio transcripts) missing by just using <i>Calypte anna</i> and <i>Aquila chrysaetos</i> (not including other avian species). Given the quality and depth of the Pacbio transcriptome, the coverage of known genes should be better.</p> <p>===</p> <p>Thank you for your analysis. First, it should be noted that, although this is a high quality and (relatively) high depth dataset, it originated from a single tissue of a single animal, which naturally limits genes detected. As we stated in the paper, "our <i>A. colubris</i> transcriptome only captured around half of this diversity, likely due to our sample being single tissue, collection time point and individual." Furthermore, our filtering steps are somewhat conservative; there is a filter in place in the Quiver/Arrow step of the pipeline, which requires that a putative isoform have at least two independent, full-length sequences, with predicted accuracy scores of >99%. These very stringent parameters were put in place to preclude the inclusion of potentially spurious transcripts into further analysis. This means that low abundance genes in this tissue and sample may be lost; though Supplemental Figure 1 suggests that we have largely saturated sequencing depth on this particular sample.</p> <p><<<</p> <p>>>></p> <p>Another aspect that has not been touched upon is the quantification of the Illumina transcriptome, which provides an unbiased view of highly expressed genes, as compared to known genes in the hepatic lipogenic pathway. The availability of the sequences from two different technologies provides a platform for comparing and contrasting them.</p> <p>===</p> <p>See "Looking at Pacbio reads..." response below for rationale against comparative quantification.</p> <p><<<</p> <p>>>></p> <p>Using Salmon (https://www.nature.com/nmeth/journal/v14/n4/abs/nmeth.4197.html) taking the <i>acolubris</i> HQD.cds as target, I have quantified the Illumina transcriptome (SRR6148275). The top ten genes are provided below - and the corresponding annotation based on the the BLAST 'nt' database.</p> <p>I1.C1447.F29P0.1522.M.5511 XM 008499421.1 <i>Calypte anna</i> acidic mammalian chitinase-like (LOC103533552)</p> <p>I1.C3140.F2P0.1756.M.20042 XM 010008736.1 <i>Chaetura pelagica</i> apolipoprotein A-I (LOC104398601)</p> <p>I0.C76367.F48P0.832.M.3544 XM 008491568.1 <i>Calypte anna</i> lipocalin-like 1 (LCNL1)</p>

I0.C68578.F2P0.585.M.3451 XM 008498154.1 Calypte anna avidin-like (LOC103532403)
I2.C288336.F2P0.1193.M.25884 KP875235.1 Pygoscelis antarcticus 18S ribosomal RNA gene
I0.C104457.F3P0.427.M.4404 XM 014961892.1 Calidris pugnax apolipoprotein A-II-like (LOC106899953)
I0.C102059.F3P0.428.M.4043 XM 014961892.1 Calidris pugnax apolipoprotein A-II-like (LOC106899953)
I1.C554372.F3P0.1606.M.14434 HM033221.1 Archilochus colubris voucher BIOUG:BIBS RTHU cytochrome
I0.C16075.F5P0.700.M.4006 XM 008504716.1 Calypte anna serum amyloid A protein-like (LOC103538359)
I3.C3772.F3P0.2060.M.41994 XM 008499863.1 Calypte anna apolipoprotein A-I (LOC103533949)

The mammalian chitinase-like is an interesting gene, which seems to have little or no reference in the current literature for hummingbirds. This might be a serum chitinase (<https://www.ncbi.nlm.nih.gov/pubmed/11591385>) associated with the I0 C16075.F5P0.700.M.4006 transcript (a serum amyloid A protein), another highly transcribed gene. Another possibility is it is a gut-chitinase (<https://www.ncbi.nlm.nih.gov/pubmed/12133911>), although I am not sure whether a liver chitinase can make its way to the gut.

===

From the Suzuki [PMID: 11591385] paper, "Avian species have dual sources for chitinase production in the stomach and liver, and mammals may have selected the major source, either the alimentary canal or liver, through evolution. The chitinase produced in hepatocytes circulates in the bloodstream and specifically defends against chitin containing microorganisms via its chitin-binding and chitin-fragmenting abilities." While humans only express chitinase in the gut (and also through certain white blood cells), chickens express the enzyme in both the gut and liver, and other mammals (cows) express the enzyme only in the liver. Suzuki et al. hypothesize that the ancestral state is expression of chitinase in both tissues. While the gut chitinase is used for digestion, expression in liver is believed to contribute to serum chitinase levels and to act as a defense against chitin-containing pathogens. The chitinase-like isoform in our dataset is highly homologous to the chicken liver chitinase. As we have not found any literature describing chitinase expression in the pancreas, nor of liver chitinase being transported to the gut for digestion, it seems likely that the hummingbird liver expresses this protein for circulation and defense against pathogens.

>>>

<<<

The high expression levels of the chitinase gene is not apparent from the Pacbio sequences (SRR5237173), probably since they are not raw reads(?). Do the Pacbio raw reads corroborate with the Illumina results?

===

Looking at the Pacbio raw reads, it appears that this chitinase-like gene is in low abundance, around 100 putative hits as opposed to 80,000 for the most abundant transcripts (acyl-CoA desaturase, steroyl CoA desaturase). This large discrepancy in expression relative to Illumina is likely due to two things; 1) differences in library preparation methodologies, and 2; inter-individual differences. cDNA production methodology is completely different (full length reverse transcription vs random priming, different amounts of PCR cycles) which could be skewing abundance drastically. Additionally, different individuals were used for the two types of sequencing; though this is reasonable for sequence correction using Illumina, it's likely not ideal for quantitation validation. In considering future experiments, we will be sure to perform such an assay to measure the relative quantitation between Illumina and Pacbio (long-read) quantitation; there seems to be too many confounding factors to confidently report this analysis. We are also planning incorporation of spike-in standards (Lexogen SIRVs) for future studies to have an additional control.

>>>

<<<

Minor comment: It would be also useful to annotate the known cds based on (high) homology with other known avian genes (I had to BLAST to 'nt' in order to obtain the

	<p>annotation. === This has been done and uploaded to Zenodo (10.5281/zenodo.1054086), thank you for the suggestion. >>></p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in publicly available repositories (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.</p> <p>Have you have met the above requirement as detailed in our Minimum</p>	Yes

[Standards Reporting Checklist?](#)

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Single molecule, full-length transcript sequencing provides insight into the extreme metabolism of ruby-throated hummingbird *Archilochus colubris*

Rachael E. Workman^{1*}, Alexander M. Myrka^{2*}, Elizabeth Tseng⁴, G. William Wong³, Kenneth C. Welch Jr.²⁺, and Winston Timp¹⁺

¹ Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD, USA

² Department of Biological Sciences, University of Toronto Scarborough, Toronto, Ontario, Canada and Department of Cell & Systems Biology, University of Toronto, Toronto, Ontario, Canada

³ Department of Physiology and Center for Metabolism and Obesity Research, Johns Hopkins University School of Medicine, Baltimore, MD, USA

⁴ Pacific Biosciences, Menlo Park, California, USA

* Co-first author

+ Co-Corresponding author

1
2
3
4
5
6 **Abstract**

7
8 **Background**

9
10
11 Hummingbirds oxidize ingested nectar sugars directly to fuel foraging but cannot sustain this
12 fuel use during fasting periods, such as during the night or during long-distance migratory
13 flights. Instead, fasting hummingbirds switch to oxidizing stored lipids, derived from ingested
14 sugars. The hummingbird liver plays a key role in moderating energy homeostasis and this
15 remarkable capacity for fuel switching. Additionally, liver is the principle location of *de novo*
16 lipogenesis, which can occur at exceptionally high rates, such as during premigratory fattening.
17 Yet understanding how this tissue and whole organism moderates energy turnover is hampered
18 by a lack of information regarding how relevant enzymes differ in sequence, expression, and
19 regulation.
20
21
22

23
24 **Findings**

25
26 We generated a *de novo* transcriptome of the hummingbird liver using PacBio full-length cDNA
27 sequencing (Iso-Seq), yielding a total of 8.6Gb of sequencing data, or 2.6M reads from 4
28 different size fractions. We analyzed data using the SMRTAnalysis v3.1 Iso-Seq pipeline, then
29 clustered isoforms into gene families to generate *de novo* gene contigs using Cogent. We
30 performed orthology analysis to identify closely related sequences between our transcriptome
31 and other avian and human gene sets. Finally, we closely examined homology of critical lipid
32 metabolism genes between our transcriptome data and avian and human genomes.
33
34
35

36 **Conclusions**

37
38 We confirmed high levels of sequence divergence within hummingbird lipogenic enzymes,
39 suggesting a high probability of adaptive divergent function in the hepatic lipogenic pathways.
40 Our results leverage cutting-edge technology and a novel bioinformatics pipeline to provide a
41 first direct look at the transcriptome of this incredible organism.
42
43
44

45 **Keywords**

46
47 Pacbio; single molecule sequencing; Iso-seq; transcriptome; liver; metabolism; hummingbirds
48

49 **Data Description**

50
51 *Background*

52
53 Hummingbirds are the only avian group to engage in sustained hovering flight as a means for
54 accessing floral nectar, their primary caloric energy source. While hovering, small
55 hummingbirds, such as the ruby-throated hummingbird (*Archilochus colubris*), achieve some of
56 the highest mass-specific metabolic rates observed among vertebrates ^{1,2}. Given their
57 specialized, sugar-rich diet, it is not that surprising that hummingbirds are able to fuel this
58
59
60
61
62
63
64
65

1
2
3
4 intense form of exercise exclusively by oxidizing carbohydrates ^{3,4}. This energetic feat is also
5 remarkable in that the source of sugar oxidized by flight muscles during hovering is the same
6 sugar ingested in nectar meals only minutes prior ^{4,5}. In addition, hummingbirds seem equally
7 adept at relying on either glucose or fructose (the two monosaccharides comprising their nectar)
8 ⁶ as a metabolic fuel for flight ⁴. In doing so, they achieve rates of sugar flux through their bodies
9 that are up to 55× greater than non-flying mammals ⁷.
10
11

12
13 Hummingbird flight is not always a solely carbohydrate-fueled endeavor. Lipids are a more
14 energy dense form of fuel storage, and fasted hummingbirds are as capable of fueling hovering
15 flight via the oxidation of onboard lipid stores as they are dietary sugars ⁵. Lipids are likely the
16 sole or predominant fuel used during overnight periods ⁸. Just as flux of sugar through the
17 hummingbird is extremely rapid, the building of lipid stores from dietary sugar is also rapid when
18 needed. For example, ruby-throated hummingbirds can routinely increase their mass by 15% or
19 more between midday and dusk on a given day ⁹. The ruby-throated hummingbird (*A. colubris*)
20 completes an arduous annual migratory journey from breeding grounds as far north as Quebec
21 in Canada to wintering grounds in Central America ¹⁰. Hummingbirds are constrained to fueling
22 long distance migratory flights using onboard lipids. In preparing for such flights, some
23 individuals rapidly build fat stores prior to departure or at migratory stopover points, increasing
24 their mass by 25-40% in as few as four days ^{9,11,12}.
25
26
27
28
29

30 The ability to switch so completely and quickly between fuel types means these animals
31 possess exquisite control over rates of substrate metabolism and biosynthesis in the liver, the
32 principal site of lipogenesis in birds ¹³. While hummingbird liver does indeed exhibit remarkably
33 high activities of lipogenic and other metabolic enzymes ¹⁴, the mechanisms underlying high
34 catalytic rates (high catalytic efficiency and/or high levels of enzyme expression) and control
35 over flux (the role of hierarchical versus metabolic control), remain unclear.
36
37
38

39 Despite long-standing recognition of, and interest in, their extreme metabolism, the lack of
40 knowledge about gene and protein sequences in hummingbirds has limited more detailed and
41 mechanistic analyses. Amplification of hummingbird genetic sequences for sequencing and/or
42 cloning is hampered by the lack of sequence information from closely related groups, making
43 well-targeted primer design difficult. Only two genes have thus far been cloned from any
44 hummingbird: an uncoupling protein (UCP) homolog and insulin ^{15,16}. These two studies offer
45 limited insight into what adaptations in hepatic molecular physiology underlie extreme energy
46 turnover or its regulation. The UCP homolog was cloned from pectoralis (flight muscle) and its
47 functional significance *in vivo* is unclear. The amino acid sequence of hummingbird insulin was
48 found to be largely identical to that from chicken; however, birds are insulin insensitive and lack
49 the insulin-regulated glucose transporter (GLUT) protein GLUT4, making the role of this
50 hormone in the regulation of energy homeostasis in hummingbirds unknown ¹⁷⁻¹⁹.
51
52
53
54
55

56 Recently completed sequencing of the Anna's hummingbird (*Calypte anna*) genome provides a
57 powerful new tool in the arsenal of biologists seeking to understand variation in metabolic
58 physiology in hummingbirds and other groups ²⁰. Despite their extreme catabolic and anabolic
59 capabilities, hummingbirds have the smallest genomes among birds ²¹ and, in general, have
60
61
62
63
64
65

1
2
3
4 among the smallest vertebrate genomes ²². Thus, it seems likely that understanding of
5 transcriptional variation, overlaid on top of genetic variation, is crucial to understanding what
6 makes these organisms such elite metabolic performers.
7
8

9 To this end, we produced the liver transcriptome of the ruby-throated hummingbird, *Archilochus*
10 *colubris*. Because many of the proteins involved in cellular metabolism are quite large, we
11 collaborated with Pacific Biosciences to generate long-read sequences as these would enhance
12 our ability to identify full coding sequences and multiple encoded isoforms. The primary
13 advantage to the PacBio Iso-seq methodology is the capability for full-length transcript
14 sequencing, rendering complete mRNA sequences without the need for assembly. This has
15 been demonstrated in previous studies to dramatically increase detection of alternative splicing
16 events ²³. Additionally, full-length sequences greatly enhance the likelihood of detecting novel or
17 rare splice variants, which is crucial for fully characterizing the transcriptomes of lesser studied,
18 non-model organisms such as the hummingbird.
19
20
21
22

23 **Methods**

24 *Sacrifice and sample preparation*

25
26
27 A wild adult male ruby-throated hummingbird (*Archilochus colubris*) was captured at the
28 University of Toronto Scarborough using modified box traps on July 23rd 2013 at 8:15AM. At
29 the time of its capture, the bird was aged as an “after hatch year” bird, meaning it was at least 1
30 year old. Standard aging techniques make more precise aging of hummingbirds more than 1
31 year old difficult ²⁴. The bird was housed in the University of Toronto Scarborough vivarium and
32 fed NEKTON-Nectar-Plus (Nekton, Tarpon Springs, FL, USA) *ad libitum*, and sacrificed after *ad*
33 *libitum* feeding at 1:22PM on July 16th 2014 (being 2+ years old). On arrival it weighed 2.68g
34 and at the time of sacrifice it weighed 3.11g. Tissues were sampled immediately after
35 euthanization using RNase-free tools. Liver tissue was dissected out and homogenized at 4°C
36 in 1 mL cold Tri Reagent using an RNase free glass tissue homogenizer and RNase free
37 syringes of increasing needle gauge. We used 100 mg of tissue per 1 mL of Tri Reagent
38 (Sigma-Aldrich, St. Louis, MO, USA), and chloroform extraction was performed twice to ensure
39 quality. RNA was precipitated with isopropanol, centrifuged at 12000xg for 10 minutes, washed
40 with ethanol 2x, vacuum dried at room temperature and eluted in RNase free water ²⁵. DNase I
41 (Life Technologies) digestion and spin column cleanup were performed (Ambion Purelink RNA
42 mini kit, Life Tech). RNA concentration and RIN were determined with RNA Bioanalyzer
43 (Agilent). The sample used for Illumina sequencing was harvested using the same methods, but
44 from a different animal. The bird was captured as described above on August 22nd, 2011 at
45 10:50AM. At the time of capture, the bird was aged as “hatch year” and it weighed 2.93g. It was
46 housed and sacrificed as described above on January 25th, 2016 at 10:50AM (being over 4
47 years old).
48
49
50
51
52
53
54
55

56 *Sequencing library preparation*

57
58 Pacific Biosciences’s Iso-Seq sequencing protocol was followed to generate sequencing
59 libraries ²⁶. Briefly, Clontech SMARTER cDNA synthesis kit with Oligo-dT primers was used to
60
61
62
63
64
65

1
2
3
4 generate first and second-strand cDNA from polyA mRNA. After a round of PCR amplification,
5 the amplified cDNA was size selected into 4 size fractions (1-2kb, 2-3kb, 3-6kb, and 5-10kb) to
6 prevent preferential small template sequencing, using the BluePippin (0.75% agarose external
7 marker, Sage Sciences). Additional PCR cycles were used post size-selection to generate
8 adequate starting material, and then SMRTbell hairpin adapters were ligated onto size-selected
9 templates. Each of the 4 size fractions was sequenced on 10 SMRT Cells, for a total of 40
10 SMRT Cells. Sequencing was performed by the JHU HiT Center using P6-C4 chemistry on the
11 RSII sequencer. Illumina sequencing libraries were generated using Lexogen mRNA sense v2
12 Illumina library preparation kit, and sequenced on a single rapid-run lane of Hiseq 4000
13 2x100bp paired end, yielded 153M reads.
14
15
16
17

18 **Analysis Methods**

19 *Data processing, isoform clustering sorting and quality control of liver transcriptome*

20
21
22 We performed initial data processing using SMRTanalysis 3.1 Iso-Seq pipeline employed using
23 a DNANexus interface. From 40 SMRTcells, we produced 440.75 Gb of raw data, which was
24 classified into 3.4 Gb of non-chimeric circular consensus (CCS) reads. CCS reads comprised
25 1.23M full length, 1.27M non-full-length reads; reads were considered full-length if both 5' and 3'
26 cDNA primers as well as the polyA tail signal were detected. Of the four size-selected bins, our
27 average CCS length was 1533, 2464, 3650, and 5444 bp, respectively (Figure 1B). The Iso-Seq
28 pipeline then performed isoform-level clustering (ICE) followed by final polishing using Arrow ²⁷
29 to output high-quality (predicted accuracy $\geq 99\%$), full-length, isoform consensus sequences.
30 The Iso-Seq pipeline produced 238Mb of high quality consensus isoforms (HQD, 94,724 reads),
31 and 2Gb (712,210 reads) of low quality consensus isoforms (summary statistics Figure 1A).
32 BLAST searches were then performed to remove putative contaminants, and coding sequence
33 and protein translation was performed, resulting in 93K HQ and 680K LQ protein sequences. A
34 summary of the analyses performed are displayed in Figure 2A-B, and further details and
35 settings can be found in Supplemental Methods.
36
37
38
39
40
41

42 *Assessing transcriptome completion*

43
44 To estimate the completeness of our liver transcriptome sequencing, we used both subsampling
45 and gene diversity estimation, as well as BUSCO ^{28,29}. BUSCO checks for essential single copy
46 orthologs which should be present in a whole transcriptome dataset for any member of the
47 given lineage. We used both Metazoan and Aves lineages (ortholog sets) to examine
48 transcriptome completion (Figure 2C and Supplemental Table 1), and to ensure that
49 completeness tracked across multiple data processing steps, we analyzed ASD (all sequence
50 data), HQD (high quality data) and CCD (Cogent collapsed data). As expected, *Gallus gallus*
51 and *Calypte anna* genome predicted transcriptomes were nearly complete for both Aves and
52 Metazoan BUSCO sets, and our *A. colubris* transcriptome only captured around half of this
53 diversity, likely due to our sample being a single tissue, collection time point and individual.
54
55
56
57
58

59 Our subsampling approach to estimating transcriptome completeness involved pulling subsets
60 of the circular consensus reads dataset and BLASTing against the predicted *Calypte anna* gene
61
62
63
64
65

1
2
3
4 set. We found that the number of unique genes detected began to saturate when reaching a
5 90% subset of our data, suggesting that additional sequencing would not substantially
6 contribute to transcriptome completion (Supplemental Figure 1). Lower expressed genes may
7 not be detected, but that vast majority of annotated liver expressed genes are likely represented
8 in our data.
9

10 *Agreement with established Anna's hummingbird genomes reveals general clade conservation*

11
12 We aligned transcripts to the *Calypte anna* (Anna's hummingbird) genome using GMAP³⁰. In
13 order to validate transcript coverage and alignment throughout the multiple processing steps,
14 we aligned using not only high quality isoforms (HQD), but also the full consensus isoform
15 dataset (ASD) and gene families predicted by Cogent (CCD, methods in Supplemental methods
16 and below).
17
18

19
20 *Calypte anna* and *Archilochus colubris* are close relatives within the North American Bee
21 (Mellisugini) clade of hummingbirds³¹; *A. colubris* is a member of the Caribbean Sheartails
22 subclade and *C. anna* is of the *Calypte* subclade, which diverged from the from ancestral
23 Mellisugini around early to mid Pliocene³². Given this fairly recent divergence, we expected
24 alignment to perform well. We found an average alignment identity of 94.8%, with 87%
25 transcripts uniquely mapping to the reference. Of the uniquely mapped, 73% covered >90% of
26 the query sequence (alignment length and statistics, Supplemental Figure 2A, 2B),
27 demonstrating high fidelity of aligned reads to reference. When ASD reads were parsed by
28 number of reads of insert supporting each consensus cluster, it was found that generally,
29 alignment identity was high regardless of number of supporting reads. A clear increase in mean
30 alignment identity was found when two or more supporting reads were collapsed (Supplemental
31 Figure 3).
32
33

34
35 When GMAP was performed using only high quality isoforms (filtered for 2+ full-length
36 supporting reads), alignment percentage was 95.7%, with 93.4% of transcripts mapping
37 uniquely to the reference. The average mapped read length was 2411bp (HQD, 2617bp ASD),
38 while the average predicted CDS length for *Calypte anna* was 1386bp. This being said, reads
39 mapped with GMAP contain UTRs. When we predict just the CDS sequences for *A. colubris*
40 using ANGEL³³, the mean length was 981bp. When we BLASTed the unaligned reads to whole
41 NCBI database, they largely mapped back to *Calypte anna* (53%). This result suggests that our
42 mapping parameters were too stringent to map these reads, error rate prevented alignment,
43 unaligned regions are divergent enough between both hummingbirds to preclude alignment, or
44 a combination of the above.
45
46
47
48
49
50
51

52 *Putative gene family prediction and reduction of transcript redundancy reduces data load while* 53 *maintaining transcript diversity*

54
55 To assign transcripts to putative gene families, as well as cluster and eliminate redundant
56 transcripts to produce a unique set of gene isoforms, we utilized the newly developed Cogent³⁴
57 pipeline. Cogent is specifically designed for transcriptome assembly in the absence of a
58 reference genome, allowing for isoforms of the same gene to be distinctly identified from
59
60
61
62
63
64
65

1
2
3
4 different gene families, which are defined as having more than two (possibly redundant)
5 transcript copies. Of the 94,724 HQ consensus isoforms, 91,733 were grouped into 6,725
6 multi-transcript gene families (Figure 3A). The remaining 2,991 sequences were classified as
7 putative single-isoform genes, or “orphans”. Reconstructed contigs were then applied in place of
8 a reference (or *de novo* clustering) to reduce redundant transcripts in the original HQD dataset.
9 From this approach, we were able to reduce our HQ dataset to 14,628 distinct transcript
10 isoforms and 2990 orphan isoforms, for a total of 17,618 isoform sequences (18% of the
11 original). Due to the use of HQD only transcripts (2 full length reads and estimated accuracy
12 >99%), and constraints of transcript collapse, a number of additional isoforms were likely lost in
13 filtering and collapse, reducing transcript diversity. However, without sufficient supporting data
14 the trade-off between between gene diversity and reliability led us to choose reliability. Future
15 studies should examine whether transcript “rescue” from low quality datasets is possible with
16 Illumina validation or additional consensus generation strategies.
17
18
19
20
21

22 Cogent collapsed data is further summarized and most abundant transcripts are detailed in
23 Supplemental Table 2. An average of 1.53 isoforms was found per gene family (Figure 3B), with
24 2624, or 27.4% of the gene families having more than one isoform, including “orphans”. While
25 other studies have found more isoforms per locus, for example 6.56 in *Zea mays*³⁵, that study
26 multiplexed six plant tissues, whereas a lower complexity is to be expected with single tissue
27 analysis. This dataset (Cogent collapsed data, or CCD) was also mapped onto the *Calypte anna*
28 genome assembly³⁶, to demonstrate the effectiveness of this method in reducing transcript
29 redundancy and classifying isoforms (Figure 3C). Cogent gene families were polished using
30 Illumina short read RNAseq data and the error correction algorithm Pilon³⁷ (Supplemental
31 Methods) to obtain higher accuracy reads.
32
33
34
35

36 *Orthologous gene pair predictions and GO annotation show putative unique hummingbird* 37 *orthologs* 38 39

40 To examine protein sequence similarity and divergence between *Archilochus colubris* and other
41 avian species, we used OrthoMCL, which generates reciprocal best hits from comparison
42 species using BLAST all-vs-all, then clustering to group orthologous sequences for each pair of
43 organisms³⁸. OrthoMCL protein sequences were predicted using ANGEL³³, and 119,292 high
44 quality sequences were put into this analysis. We compared our ruby-throated hummingbird,
45 *Archilochus colubris*, to five other birds: *Calypte anna* (Anna’s hummingbird) fellow member of
46 the bee clade of hummingbirds, *Chaetura pelagica* (chimney swift) the closest available
47 outgroup species to the hummingbird clade, and other bird species for which relatively
48 well-annotated genomes and/or transcriptomes are available, *Gallus gallus* (chicken),
49 *Taeniopygia guttata* (zebra finch), and *Melopsittacus undulatus* (budgerigar), as well as *Homo*
50 *sapiens* (human), and *Alligator mississippiensis* (American alligator). Algorithm parameters and
51 data accession numbers are presented in Supplemental methods.
52
53
54
55
56

57 A matrix of ortholog pairings, with duplicate ortholog hits removed, shows the number of
58 orthologous sequences for each species pair (Supplemental Table 3). Orthologs shared
59 between ruby-throated hummingbird and a subset of the other species analyzed are illustrated
60
61
62
63
64
65

1
2
3
4 in Figure 4A. Unsurprisingly, the largest amount of orthologs which pair closely to only one
5 species, i.e., 1:1 orthologs, were found between Anna's and Ruby-throated hummingbird.
6 Surprisingly, the second-largest set was between chicken and ruby-throated hummingbird, as
7 opposed to its closest outgroup species, *Chaetura pelagica*. This is likely due to the
8 completeness of chicken transcriptome annotation, as chicken is the most well-studied avian
9 species. Of the 596 unpaired *A. colubris* protein sequences, 190 paired most closely with
10 *Calypte anna* when compared using BlastP and the majority of matches output (559/594) were
11 less than 50 AA, only a fraction of the average sequence length.
12
13
14

15
16 In order to more closely examine the identity of orthologs in related hummingbird species, gene
17 ontology (GO) annotation was performed on the set of orthologs which were shared between
18 *Calypte anna* and *Archilochus colubris*, but not by the other birds included in the OrthoMCL
19 analysis. This set of 2,376 protein sequences was examined using BlastP and GO analysis
20 performed by Panther^{39,40}. Additional datasets used for GO comparison included 1:1 orthologs
21 for *Gallus gallus* and *A. colubris* (518), and *A. colubris* and *Chaetura pelagica* (430), as well as
22 whole transcriptome data from *C. pelagica* and the Cogent-collapsed dataset from our
23 transcriptome (Supplemental Table 4, Figure 4B).
24
25
26

27
28 As the initial impetus for our investigation centered on the exceptional metabolism and
29 energetics of hummingbirds, we focused our investigation on orthologs tagged as part of the
30 "metabolic process (GO:0008152)" grouping. Of the 1444 orthologs identified in *Archilochus*
31 *colubris* as part of this process grouping, 236 (16.3%) were unique to hummingbirds. Within this
32 top-level grouping, the largest number of genes group under "primary metabolic processes"
33 (GO:0044238)". Of the 1240 orthologs identified within this grouping, 204 (16.3%) are identified
34 as uniquely shared by our hummingbird species. Six GO biological processes are defined under
35 the "primary metabolic processes". Of these processes, the process with the highest proportion
36 of identified *A. colubris* orthologs hitting as unique to the two hummingbird species is "lipid
37 metabolic processes" (GO:0006629; 33 of 114 orthologs, 28.9%), which is significantly enriched
38 relative to the comparative orthology databases of both chicken and human (Statistical
39 overrepresentation test, Panther,³⁹ p-values given in Supplemental table 4). Because we
40 considered it likely that an enrichment in lipid metabolic genes could be a result of our dataset
41 being from liver tissue, we compared enrichment with that of the entire Cogent predicted gene
42 set from the ruby-throated hummingbird transcriptome, and found no significant enrichment
43 using the same tests (Supplemental table 4). Because 1:1 hummingbird orthologs are relatively
44 more abundant in lipid metabolic genes than the sequences which were found to be highly
45 homologous to one or more of the other species compared using OrthoMCL, we predict that
46 lipid metabolic genes are more divergent from the other examined species than other classes of
47 enzymes. Though this alone is not direct evidence of greater selection on proteins within that
48 pathway, it is suggestive. If neutral sequence divergence is assumed to be randomly accrued
49 throughout a species' genome, then greater divergence in enzymes making up "lipid metabolic
50 processes" suggests that closer examination of these proteins for evidence of functional, or
51 even adaptive, divergence is warranted. A phylogenetically-informed analysis of ortholog
52 divergence among taxa is necessary to establish a selection signature, which will become
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 possible in the future with the advance of the B10K project ⁴¹ and larger numbers of avian
5 species in GO databases.
6

7
8 Given the apparent sequence divergence among enzymes involved in “lipid metabolic
9 processes” hinted at by orthology and ontology analyses, we elected to more closely examine
10 enzymes comprising the lipogenic pathway. In liver, fatty acids can be synthesized via the *de*
11 *novo* lipogenesis pathway using acetyl CoA as substrate. These newly synthesized fatty acids
12 can then be esterified onto the glycerophosphate backbone to generate triglycerides via the
13 glycerol-3-phosphate pathway of lipid synthesis. We predicted that key enzymes involved in
14 these two pathways (Figure 5A) would be divergent in hummingbirds given their extraordinary
15 metabolic demands. Eight enzymes involved in this pathway were examined for *Archilochus*
16 *colubris*, *Calypte anna*, *Gallus Gallus*, *Chaetura pelagica*, *Alligator mississippiensis* and *Homo*
17 *sapiens* (accession numbers and details given in Supplemental Table 5). Pairwise protein
18 alignment scores are given in Supplemental Table 6 as well as illustrated in a heatmap shown in
19 Figure 5B, and alignments in Supplemental Data 1. Interestingly, enzymes involved in *de novo*
20 fatty acid synthesis share a higher degree of identity between examined organisms, whereas
21 enzymes involved in triglyceride synthesis tend to be slightly less conserved (Figure 5A). Figure
22 5B also shows normalized abundances of the enzymes of interest in our liver transcriptome
23 dataset, revealing a high expression level of the rate-setting enzyme involved in *de novo*
24 lipogenesis (*ACACA*; acetyl CoA carboxylase). In contrast to the cytosolic *ACACA* enzyme that
25 uses acetyl-CoA as substrates for fatty acid synthesis, *MCAT* encodes a mitochondrial enzyme
26 that uses malonyl-CoA as substrates for fatty acid synthesis. Much less is known about the
27 *MCAT*-dependent pathway of fatty acid synthesis in mitochondria. Interestingly, *MCAT* has the
28 lowest relative abundance in ruby throated hummingbird liver. The relative hepatic expression
29 levels of triglyceride synthesis genes (e.g., *LPIN1* and *DGAT2*) are also much lower compared
30 to genes involved in *de novo* lipogenesis (*ACACA* and *FASN*). It is important to note that most
31 metabolic enzymes are tightly regulated. The relative levels of hepatic lipogenesis enzymes
32 may vary greatly depending on the time of day and the physiological states (fast vs. fed) of the
33 animals.
34
35
36
37
38
39
40
41
42

43 In order to further investigate degree of conservation between key hepatic lipogenesis enzymes
44 in hummingbirds and comparative organisms, we performed conservation analysis and
45 determined ratio of nonsynonymous to synonymous codon changes (dN/dS) as a metric of
46 positive selection, using pairwise alignments followed by the CodeML module in PAML⁴².
47 These ratios are given in Supplemental Table 6 and plotted in a heatmap in Figure 5B. A dN/dS
48 score > 1 denotes genomic regions putatively undergoing positive selection. We found, in
49 general, good conservation of these enzymes among species, with the exception of the 3' and
50 5' ends of alignments. These often had an extended or retracted coding sequence in the case of
51 hummingbirds and *C. pelagica*, which could be related to post-translational modification or
52 selection on pathway regulation ⁴³. Surprisingly, terminal sequence length was variable even
53 between *C. anna* and *A. colubris*, which both belong to the closely-related Bee hummingbird
54 taxon ³¹. Variation in 5' and 3' length may also be an effect of the different methodologies used
55 to produce these sequences, RNA sequencing for *A. colubris*, *G. gallus*, and *H. sapiens*, and
56
57
58
59
60
61
62
63
64
65

1
2
3
4 ORF prediction from genomic data for the other organisms examined. For example, we note in
5 our analysis that *MCAT* appears more conserved between *A. colubris* and *H. sapiens*, than
6 between *A. colubris* and *C. anna*, which could be due not to *A. colubris* actually being more
7 similar to *H. sapiens*, but rather to ORF prediction oversights.
8
9

10 The averaged dN/dS values, while useful for comparison, can be misleading when considered
11 over the entire gene, as 3' and 5' variation can overshadow conserved motifs, and pairwise
12 comparisons (Supplemental Data 1 and 2) are limited in scope. This type of analysis is ideal for
13 very divergent sequences, but less informative for pairs of sequences that are highly similar⁴⁴.
14 Despite this, conservation analysis is still valuable and provides insights connecting nucleotide
15 to amino acid information that alignments alone can miss. For example, lysophosphatidic acid
16 acyltransferase (*ABHD5*), which functions primarily in phosphatidic acid biosynthesis, has
17 reasonable protein alignment scores to all comparative organisms but also shows positive
18 selection acting upon this gene relative to *Calypte anna*, swift, human and alligator, but not
19 chicken (Figure 5B). This led us to more closely examine the coding sequence alignment, where
20 we found that the bulk of differences in coding sequence were attributable to exon 1, with
21 alignment largely becoming synchronous (with the exception of *H. sapiens*, which is widely
22 divergent) by exon 2 and continuing through to the end of the transcript. Although the primary
23 AB hydrolase-1 domain is very well conserved between species, these differences in exon 1
24 could be functionally significant, and honing down to regions of differentiation between
25 comparative species gives us interesting starting points for future investigations, including the
26 cloning and enzyme kinetics studies of *ABHD5*. Additionally, pairwise comparisons provide
27 interesting observations, such as coding strand elongation in the 5' region in *A. colubris* *GPAM*
28 (Supplemental Data 2). This information can be leveraged for future studies examining enzyme
29 structure, function and evolution.
30
31
32
33
34
35
36
37

38 *Transcriptome resource mining could provide functional genomic insights*

39

40 Access to the transcriptome informs the investigation of biological processes and enables the
41 formation of new hypotheses. This is exemplified by the serendipitous observation that
42 hummingbird glucose transporter 2 (*GLUT2*) lacks a N-glycosylation site due to an asparagine
43 to aspartic acid amino acid substitution. This missing glycosylation site was also seen in the
44 available Anna's hummingbird genome. All class 1 glucose transporters studied in model
45 vertebrates contain one N-glycosylation site located on the large extracellular loop of the protein
46⁴⁵. In *GLUT2* the associated glycan interacts with the glycan-galectin lattice of the cell,
47 stabilizing cell surface expression⁴⁶. Removal of the N-glycan of *GLUT2* in rat pancreatic β cells
48 results in the sequestering of cell-surface *GLUT2* in lipid rafts and this sequestered *GLUT2*
49 exhibits a reduction in glucose transport activity by approximately 25%⁴⁶. This reduction in
50 transport is thought to occur through interaction of the *GLUT* with lipid raft-bound stomatin^{46,47}.
51 In mammals, *GLUT2* serves a glucose-sensing role in the pancreatic β cells and is required for
52 the regulation of blood glucose through insulin and glucagon⁴⁸. The lack of N-glycosylation of
53 *GLUT2* may contribute to observed high blood glucose concentration in hummingbirds⁴⁹.
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7 Another serendipitous observation was the highly abundant chitinase-like transcript noted from
8 Illumina sequencing results. While humans express chitinase in the gut, but not the liver,
9 chickens express the enzyme in both gut and liver, and other mammals (cows) express the
10 enzyme only in the liver⁵⁰. Suzuki et al. hypothesize that the ancestral state is expression of
11 chitinase in both tissues. While the gut chitinase is used for digestion, expression in liver is
12 believed to contribute to serum chitinase levels and to act as a defense against chitin-containing
13 pathogens⁵⁰; perhaps the second animal had an infection? The chitinase-like isoform in our
14 dataset is highly homologous to the chicken liver chitinase-like transcript.
15
16

17 *Re-use potential*

18
19
20 In conclusion, our results have leveraged cutting-edge technology to provide a compelling first
21 direct look at the transcriptome of this incredible organism. By using PacBio sequencing, we
22 have been able to generate full length cDNA transcripts from the hummingbird liver.
23 Transcriptome data generated using the Iso-seq methodology, when coupled to recently
24 developed sophisticated gene synthesis techniques⁵¹, will allow simple generation of relevant
25 isoforms for biochemical experiments. Some of the key metabolic enzymes identified from our
26 work as being unique to either *A. colubris* or at most common to *C. anna* and *A. colubris* can
27 now be quickly cloned and expressed. Follow up studies will allow for biochemical studies of
28 proteins generated directly from our transcriptome data, measuring their enzymatic properties,
29 e.g. k_{cat} or V_{max} , as compared to other avian or mammalian analogues^{14,52,53}. Expressed proteins
30 may also be used for structural biology studies, applying either x-ray crystallography or cryoEM
31 to generate structural maps of the proteins, then examine how the structure compares to other
32 analogues.
33
34
35
36
37

38 *Availability of supporting data*

39
40
41 Filtered fastq of clustered CCS reads deposited in SRA accession number SRP099041.
42 Predicted Cogent gene families, coding sequence and annotations, peptide and untranslated
43 region data are available at 10.5281/zenodo.1054086. All other data available upon request.
44
45

46 *Availability and requirements*

47
48 Project name: Ruby_iseq

49
50 Project home page: <https://github.com/reworkman/hummingbird>

51
52 Operating system: Unix

53
54 Programming language: Bash, Python, R

55
56 Other requirements: BUSCO, GMAP, Blast+, ANGEL, CLUSTAL, Cogent, and their
57 dependencies
58
59
60
61
62
63
64
65

1
2
3
4 License: None
5

6 **Acknowledgments**

7
8 Pacific Biosciences for reagents and SMRTcells as well as technical support. M. Schatz, E.
9 Jarvis, J. Korlach, Y. Guo for discussion. HFSP grant #RGP0062/2016. Natural Sciences and
10 Engineering Research Council of Canada Discovery Grant (#386466) to KCW.
11
12

13 **Disclosure Declaration**

14
15 W.T. and R.W. have received travel funds to speak at symposia organized by Pacific
16 Biosciences. Bulk of reagents for IsoSeq were provided by Pacific Biosciences.
17
18

19 **Figure legends**

20
21 **Figure 1.** Transcriptome dataset quality control reveals good throughput, read length, and
22 transcriptome completion. Average read lengths and isoform counts for 4 sequenced size
23 fractions given in **A**, and read length distribution for all sequence data (ASD, is all sequence
24 data, high quality (HQ) and low quality (LQ) isoforms) on x vs read counts on y plotted in **B**, with
25 black line representing Mb data greater than read length. For example, at 2000bp, 4000Mb of
26 sequence data was larger than 2000bp. **C.** BUSCO transcriptome assessment results displayed
27 for *Archilochus colubris* (ruby-throated hummingbird, all sequence data ASD, high quality
28 sequence data HQD), Cogent-collapsed data (CCD), *Calypste anna* (Anna's hummingbird),
29 *Gallus gallus* Thomas (chicken single-tissue transcriptome²⁶) illustrate transcriptome completion
30 relative to predicted single copy ortholog datasets for both the Class Aves and Kingdom
31 Metazoa.
32
33
34
35
36

37 **Figure 2.** Analysis pipeline details, as well as amount of data present at each step (in green
38 text). **A** Raw sequence reads from a Pacbio RSII sequencer (bax.h5, bas.h5) were sorted into
39 full and non-full length reads of insert (ROI) using a classification algorithm that identified full
40 length reads with forward and reverse primers, as well as a poly-A tail. Iterative clustering for
41 isoforms (ICE) was performed on full length reads, and non-full length reads were recruited to
42 perform ARROW polished on the consensus isoforms. Polishing sorted reads into high and
43 low-quality bins, and either high quality data (HQD), all sequence data (ASD) or both sets of
44 data, were carried on to further applications (**B**). Applications include ORF and protein
45 sequence generation from high quality (HQD) and low quality (LQD) consensus isoforms,
46 alignment to *C. anna* reference with GMAP of both high quality data (HQD) and
47 Cogent-collapsed data (CCD), detection of orthologous sequences (orth groups) using
48 OrthoMCL, and prediction of gene families (gene fam) using Cogent. Numbers of available
49 reads at each analysis step is displayed in green in each bubble.
50
51
52
53
54

55 **Figure 3.** Reducing transcript redundancy and predicting gene families using Cogent software.
56 **A.** Gene families predicted and classified by relationship to *Calypste anna* genome assembly
57 shown, along with statistics for alignment using GMAP software which show excellent alignment
58 to closely related hummingbird reference species *Calypste anna*. Cogent comparison cases
59
60
61
62
63
64
65

1
2
3
4 highlight the relationships between predicted gene families and *C. anna* reference (column
5 captioned “In ref”), and demonstrate the additional information given to an assembly by
6 transcriptome information. Number of isoforms predicted per gene family (unigene) given in **B**
7 shows relatively low isoform diversity in this tissue. **C** Alignment of the MATR3 gene
8 demonstrates the redundancy reducing capabilities of the Cogent software, which was reduced
9 from 11 semi-redundant reads to 3 unique isoforms using this pipeline.
10
11

12
13 **Figure 4.** Orthology analysis. The proteomes of five birds (Anna’s hummingbird:*Calypste anna*,
14 Zebrafinch: *Tinamus guttatus*, Chicken: *Gallus gallus*, Swift: *Chaetura pelagica*, and
15 Budgeridger:*Melopsitticus undulatus*), one mammal (*Homo sapiens*) and one reptile (*Alligator*
16 *mississippiensis*) were compared against *Archilochus colubris* using OrthoMCL to detect
17 homologous sequences. A Venn diagram illustrating sequences with best reciprocal blast hits
18 between the given species and *A. colubris* is shown in **A**. Bar chart illustrates
19 observed/expected ratios of metabolism enzymes (GO group: metabolic process 0008152) for
20 comparison groups (statistical overrepresentation test) for selected OrthoMCL groups using
21 Panther. Datasets input either include the entire proteome of target species (swift all, anna’s all)
22 or distinct set of homologs shared two groups (Ex. *A. colubris* + *C. anna* are homologs shared
23 between these two species but not any of the other comparison groups). Asterisks denote
24 significant overrepresentation of metabolic process proteins relative to expected baseline
25 ($p < 0.05$)(**B**).
26
27
28
29
30

31 **Figure 5.** Pathway analysis of key enzymes in hepatic lipogenesis. **A** An overview of the
32 relationship between the investigated genes and their roles in triacylglycerol, phospholipid and
33 fatty acid synthesis, **B** heat maps illustrating percent amino acid identity of these proteins
34 relative to *Archilochus colubris* predicted sequences, abundances ($\log_2(\text{reads per } 10000)$)
35 transformed) of their transcripts, and dN/dS (ratio of synonymous to nonsynonymous gene
36 mutations). Taken together, these illustrate the complex relationships between target proteins
37 and identity, conservation and abundance.
38
39
40

41 **References:**

- 42
43
44 1. Suarez, R. K. Hummingbird flight: sustaining the highest mass-specific metabolic rates
45 among vertebrates. *Experientia* **48**, 565–570 (1992).
46
47
48 2. Chai, P. & Dudley, R. Limits to flight energetics of hummingbirds hovering in hypodense
49 and hypoxic gas mixtures. *J. Exp. Biol.* **199**, 2285–2295 (1996).
50
51
52
53 3. Suarez, R. K. *et al.* Fuel selection in rufous hummingbirds: ecological implications of
54 metabolic biochemistry. *Proc. Natl. Acad. Sci. U. S. A.* **87**, 9207–9210 (1990).
55
56
57
58 4. Chen, C. C. W. & Welch, K. C. Hummingbirds can fuel expensive hovering flight completely
59
60
61
62
63
64
65

- 1
2
3
4 with either exogenous glucose or fructose. *Funct. Ecol.* **28**, 589–600 (2014).
5
6
7 5. Welch, K. C., Jr, Altshuler, D. L. & Suarez, R. K. Oxygen consumption rates in hovering
8 hummingbirds reflect substrate-dependent differences in P/O ratios: carbohydrate as a
9 'premium fuel'. *J. Exp. Biol.* **210**, 2146–2153 (2007).
10
11
12
13 6. Baker, H. G. Sugar Concentrations in Nectars from Hummingbird Flowers. *Biotropica* **7**,
14 37–41 (1975).
15
16
17
18 7. Welch, K. C., Jr & Chen, C. C. W. Sugar flux through the flight muscles of hovering
19 vertebrate nectarivores: a review. *J. Comp. Physiol. B* **184**, 945–959 (2014).
20
21
22
23 8. Powers, D. R., Brown, A. R. & Van Hook, J. A. Influence of normal daytime fat deposition
24 on laboratory measurements of torpor use in territorial versus nonterritorial hummingbirds.
25 *Physiol. Biochem. Zool.* **76**, 389–397 (2003).
26
27
28
29 9. Hou, L., Verdirame, M. & Welch, K. C. Automated tracking of wild hummingbird mass and
30 energetics over multiple time scales using radio frequency identification (RFID) technology.
31 *J. Avian Biol.* **46**, 1–8 (2015).
32
33
34
35
36
37 10. Weidensaul, S., Robinson, T. R., Sargent, R. R., Sargent, M. B. & Poole, A. Ruby-throated
38 hummingbird (*Archilochus colubris*). *Birds of North America Online* (A. Poole, Editor).
39 *Cornell Lab of Ornithology, Ithaca, NY, USA.* <http://bna.birds.cornell.edu/bna/species/204>
40 (2013).
41
42
43
44
45
46 11. Carpenter, F. L., Hixon, M. A., Beuchat, C. A., Russell, R. W. & Paton, D. C. Biphase Mass
47 Gain in Migrant Hummingbirds: Body Composition Changes, Torpor, and Ecological
48 Significance. *Ecology* **74**, 1173–1182 (1993).
49
50
51
52
53 12. Hou, L. & Welch, K. C., Jr. Premigratory ruby-throated hummingbirds, *Archilochus colubris*,
54 exhibit multiple strategies for fuelling migration. *Anim. Behav.* **121**, 87–99 (2016).
55
56
57
58 13. Hermier, D. Lipoprotein metabolism and fattening in poultry. *J. Nutr.* **127**, 805S–808S
59
60
61
62
63
64
65

- 1
2
3
4 (1997).
5
6
7 14. Suarez, R. K., Brownsey, R. W., Vogl, W., Brown, G. S. & Hochachka, P. W. Biosynthetic
8 capacity of hummingbird liver. *Am. J. Physiol.* **255**, R699–702 (1988).
9
10
11 15. Vianna, C. R. *et al.* Cloning and functional characterization of an uncoupling protein
12 homolog in hummingbirds. *Physiol. Genomics* **5**, 137–145 (2001).
13
14
15 16. Fan, L., Gardner, P., Chan, S. J. & Steiner, D. F. Cloning and analysis of the gene encoding
16 hummingbird proinsulin. *Gen. Comp. Endocrinol.* **91**, 25–30 (1993).
17
18
19 17. Welch, K. C., Jr, Allalou, A., Sehgal, P., Cheng, J. & Ashok, A. Glucose transporter
20 expression in an avian nectarivore: the ruby-throated hummingbird (*Archilochus colubris*).
21 *PLoS One* **8**, e77003 (2013).
22
23
24 18. Braun, E. J. & Sweazea, K. L. Glucose regulation in birds. *Comp. Biochem. Physiol. B*
25 *Biochem. Mol. Biol.* **151**, 1–9 (2008).
26
27
28 19. Polakof, S., Mommsen, T. P. & Soengas, J. L. Glucosensing and glucose homeostasis:
29 from fish to mammals. *Comp. Biochem. Physiol. B Biochem. Mol. Biol.* **160**, 123–149
30 (2011).
31
32
33 20. Jarvis, E. D. *et al.* Phylogenomic analyses data of the avian phylogenomics project.
34 *Gigascience* **4**, 4 (2015).
35
36
37 21. Gregory, T. R., Andrews, C. B., McGuire, J. A. & Witt, C. C. The smallest avian genomes
38 are found in hummingbirds. *Proc. Biol. Sci.* **276**, 3753–3757 (2009).
39
40
41 22. Hughes, A. L. & Hughes, M. K. Small genomes for better flyers. *Nature* **377**, 391 (1995).
42
43
44 23. Abdel-Ghany, S. E. *et al.* A survey of the sorghum transcriptome using single-molecule long
45 reads. *Nat. Commun.* **7**, 11706 (2016).
46
47
48 24. Pyle, P. Identification Guide to North American Birds, Part I: Columbidae to Ploceidae:
49 Peter Pyle: 9780961894023: Books - Amazon.ca. Available at:
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 <https://www.amazon.ca/Identification-Guide-North-American-Birds/dp/0961894024>.

5
6
7 (Accessed: 1st October 2017)

- 8
9 25. Chomczynski, P. & Mackey, K. Short technical reports. Modification of the TRI reagent
10 procedure for isolation of RNA from polysaccharide- and proteoglycan-rich sources.
11
12 *Biotechniques* **19**, 942–945 (1995).
13
14
15 26. Thomas, S., Underwood, J. G., Tseng, E., Holloway, A. K. & Bench To Basinet CvDC
16 Informatics Subcommittee. Long-read sequencing of chicken transcripts and identification
17 of new transcript isoforms. *PLoS One* **9**, e94650 (2014).
18
19
20 27. Genomic Consensus. *Github* Available at:
21
22 <https://github.com/PacificBiosciences/GenomicConsensus>.
23
24
25 28. Assessing genome assembly and annotation completeness with Benchmarking Universal
26 Single-Copy Orthologs (BUSCO). Available at: <https://gitlab.com/ezlab/busco>.
27
28
29 29. Simão, F. A., Waterhouse, R. M., Ioannidis, P., Kriventseva, E. V. & Zdobnov, E. M.
30 BUSCO: assessing genome assembly and annotation completeness with single-copy
31 orthologs. *Bioinformatics* **31**, 3210–3212 (2015).
32
33
34 30. Wu, T. D. & Watanabe, C. K. GMAP: a genomic mapping and alignment program for mRNA
35 and EST sequences. *Bioinformatics* **21**, 1859–1875 (2005).
36
37
38 31. McGuire, J. A., Witt, C. C., Altshuler, D. L. & Remsen, J. V., Jr. Phylogenetic systematics
39 and biogeography of hummingbirds: Bayesian and maximum likelihood analyses of
40 partitioned data and selection of an appropriate partitioning strategy. *Syst. Biol.* **56**,
41 837–856 (2007).
42
43
44 32. Licona-Vera, Y. & Ornelas, J. F. The conquering of North America: dated phylogenetic and
45 biogeographic inference of migratory behavior in bee hummingbirds. *BMC Evol. Biol.* **17**,
46 126 (2017).
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
- 2
- 3
- 4 33. Tseng, E. Robust Open Reading Frame prediction (ANGLE re-implementation). Available
- 5 at: <https://github.com/PacificBiosciences/ANGEL>.
- 6
- 7
- 8
- 9 34. Coding Genome Reconstruction using Iso-Seq data. *Elizabeth Tseng* Available at:
- 10 <https://github.com/Magdoll/Cogent>.
- 11
- 12
- 13 35. Wang, B. *et al.* Unveiling the complexity of the maize transcriptome by single-molecule
- 14 long-read sequencing. *Nat. Commun.* **7**, 11708 (2016).
- 15
- 16
- 17 36. Zhang, G. *et al.* Genomic data of the Anna's Hummingbird (*Calypte anna*). (2014).
- 18 doi:10.5524/101004
- 19
- 20
- 21
- 22 37. Walker, B. J. *et al.* Pilon: an integrated tool for comprehensive microbial variant detection
- 23 and genome assembly improvement. *PLoS One* **9**, e112963 (2014).
- 24
- 25
- 26 38. Li, L., Stoeckert, C. J., Jr & Roos, D. S. OrthoMCL: identification of ortholog groups for
- 27 eukaryotic genomes. *Genome Res.* **13**, 2178–2189 (2003).
- 28
- 29
- 30 39. Mi, H., Muruganujan, A., Casagrande, J. T. & Thomas, P. D. Large-scale gene function
- 31 analysis with the PANTHER classification system. *Nat. Protoc.* **8**, 1551–1566 (2013).
- 32
- 33
- 34 40. Mi, H. *et al.* PANTHER version 11: expanded annotation data from Gene Ontology and
- 35 Reactome pathways, and data analysis tool enhancements. *Nucleic Acids Res.* **45**,
- 36 D183–D189 (2017).
- 37
- 38 41. Zhang, G. *et al.* Genomics: Bird sequencing project takes off. *Nature* **522**, 34 (2015).
- 39
- 40 42. Yang, Z. PAML 4: phylogenetic analysis by maximum likelihood. *Mol. Biol. Evol.* **24**,
- 41 1586–1591 (2007).
- 42
- 43
- 44 43. Jacob, E. & Unger, R. A tale of two tails: why are terminal residues of proteins exposed?
- 45 *Bioinformatics* **23**, e225–30 (2007).
- 46
- 47 44. Kryazhimskiy, S. & Plotkin, J. B. The population genetics of dN/dS. *PLoS Genet.* **4**,
- 48 e1000304 (2008).
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60
- 61
- 62
- 63
- 64
- 65

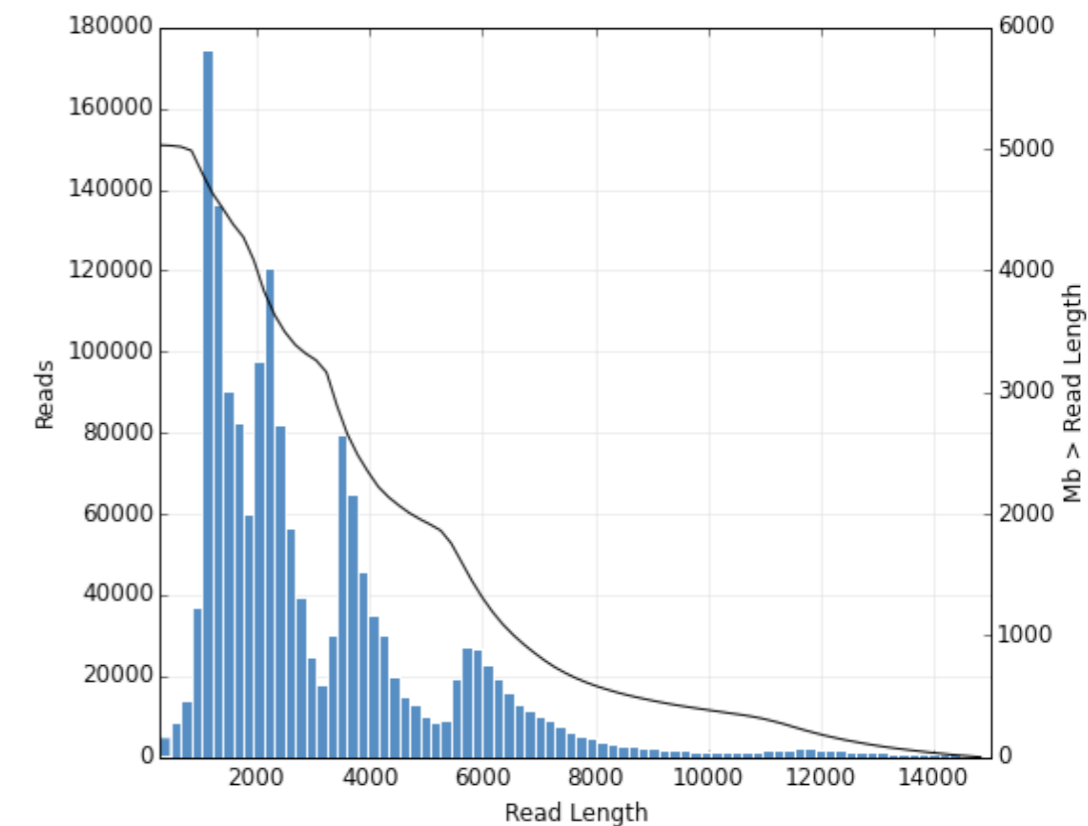
- 1
2
3
4 45. Joost, H. G. & Thorens, B. The extended GLUT-family of sugar/polyol transport facilitators:
5 nomenclature, sequence characteristics, and potential function of its novel members
6
7
8 (review). *Mol. Membr. Biol.* **18**, 247–256 (2001).
9
- 10
11 46. Ohtsubo, K. *et al.* N-Glycosylation modulates the membrane sub-domain distribution and
12 activity of glucose transporter 2 in pancreatic beta cells. *Biochem. Biophys. Res. Commun.*
13
14 **434**, 346–351 (2013).
15
16
17
- 18 47. Zhang, J. Z., Abbud, W., Prohaska, R. & Ismail-Beigi, F. Overexpression of stomatin
19 depresses GLUT-1 glucose transporter activity. *Am. J. Physiol. Cell Physiol.* **280**,
20
21 C1277–83 (2001).
22
23
24
- 25 48. Thorens, B. & Mueckler, M. Glucose transporters in the 21st Century. *Am. J. Physiol.*
26
27 *Endocrinol. Metab.* **298**, E141–5 (2010).
28
29
- 30 49. Beuchat, C. A. & Chong, C. R. Hyperglycemia in hummingbirds and its consequences for
31 hemoglobin glycation. *Comp. Biochem. Physiol. A Mol. Integr. Physiol.* **120**, 409–416
32
33 (1998).
34
35
- 36 50. Suzuki, M. *et al.* Cellular expression of gut chitinase mRNA in the gastrointestinal tract of
37 mice and chickens. *J. Histochem. Cytochem.* **50**, 1081–1089 (2002).
38
39
40
- 41 51. Kosuri, S. & Church, G. M. Large-scale de novo DNA synthesis: technologies and
42 applications. *Nat. Methods* **11**, 499–507 (2014).
43
44
45
- 46 52. Suarez, R. K., Welch, K. C., Jr, Hanna, S. K. & Herrera M, L. G. Flight muscle enzymes and
47 metabolic flux rates during hovering flight of the nectar bat, *Glossophaga soricina*: further
48 evidence of convergence with hummingbirds. *Comp. Biochem. Physiol. A Mol. Integr.*
49
50 *Physiol.* **153**, 136–140 (2009).
51
52
53
- 54 53. FernándezM.J., BozinovicF. & SuarezR.K. Enzymatic flux capacities in hummingbird flight
55 muscles: a ‘one size fits all’ hypothesis. *Can. J. Zool.* **89**, 985–991 (2011).
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

A]

Size Fraction	1-2kb	2-3kb	3-6kb	5-10kb	Total
# of SMRTcells	10	10	10	10	40
Reads of Insert (ROI)	688,069	591,050	735,670	625,194	2,639,983
Avg length ROI (bp)	1533	2464	3650	5444	
ROI Yield (Mbp)	1055	1457	2685	3404	8601
Filtered full length reads	430,381	306,841	272,781	193,906	1,203,909
# Consensus Isoforms	359,981	163,618	209,969	121,109	807,114
HQ consensus isoforms	41,763	25,776	24,735	7,436	94,724
% HQ	11.60%	15.75%	11.78%	6.14%	11.74%
Avg HQ length	1315	2329	3629	5491	
LQ consensus isoforms	321,101	135,415	186,523	113,162	712,210
% LQ	89.20%	82.76%	88.83%	93.44%	88.56%
Avg LQ length	1503	2621	4170	6718	

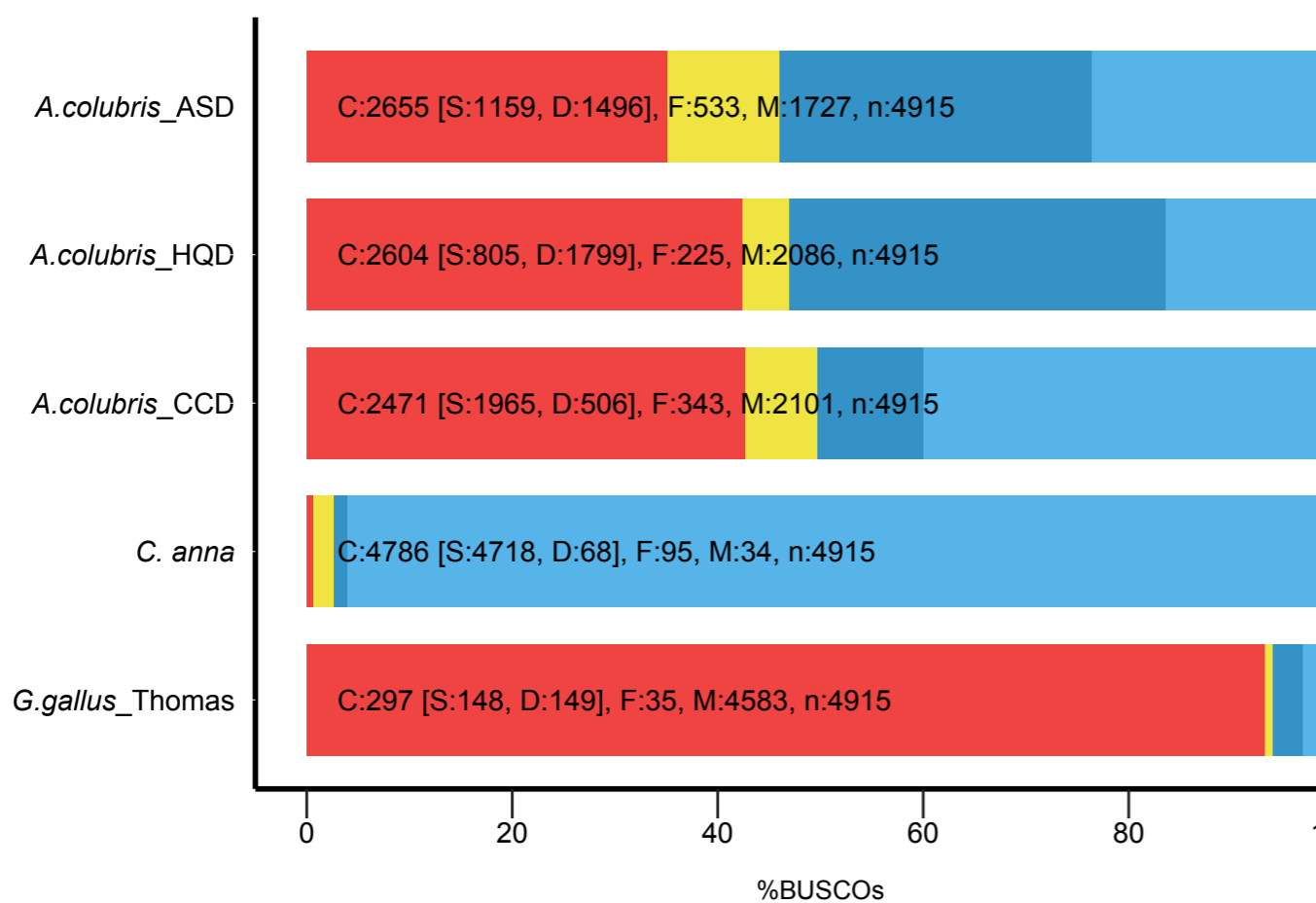
B]



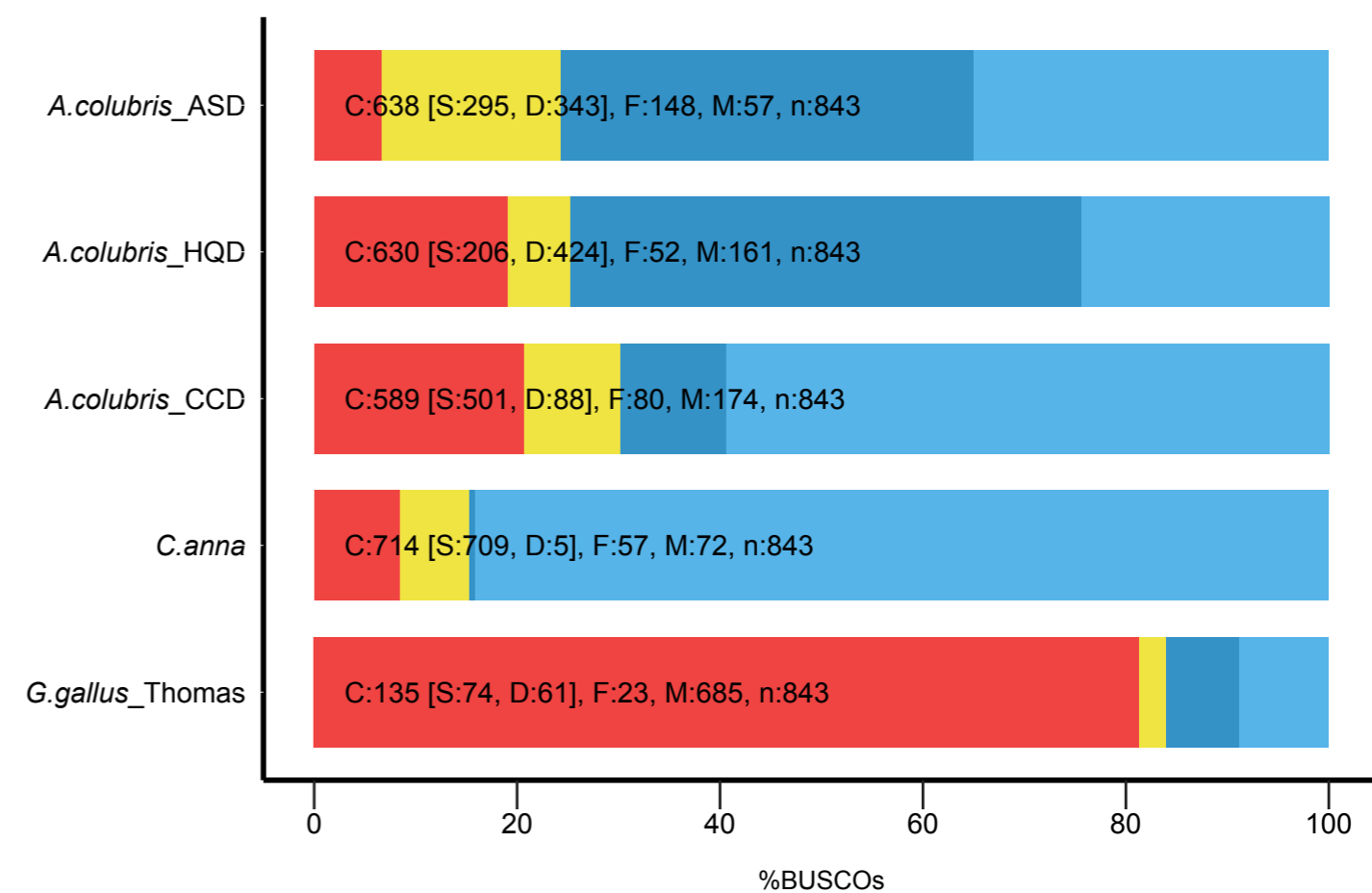
C] BUSCO ASSESSMENT RESULTS



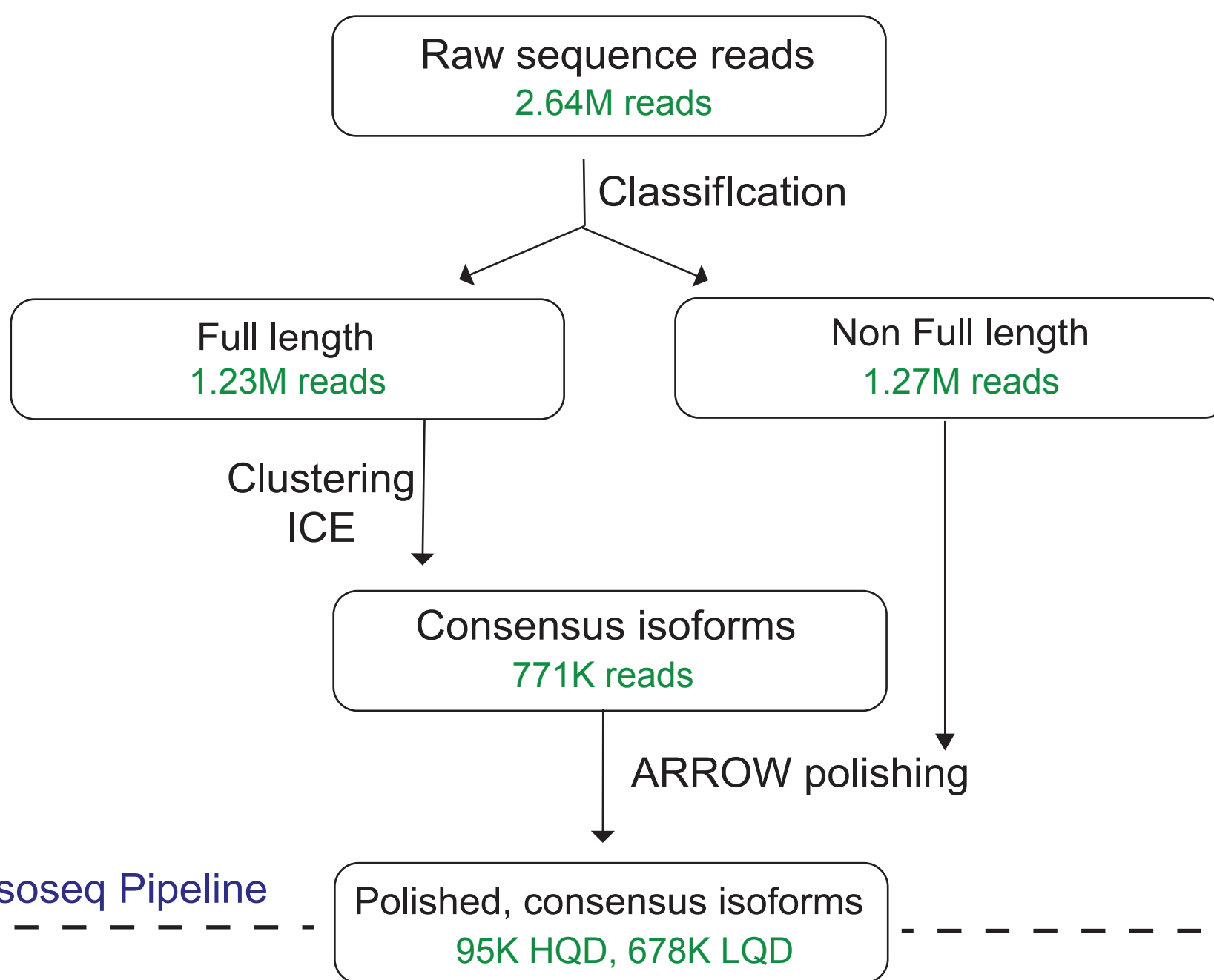
Aves



Metazoan



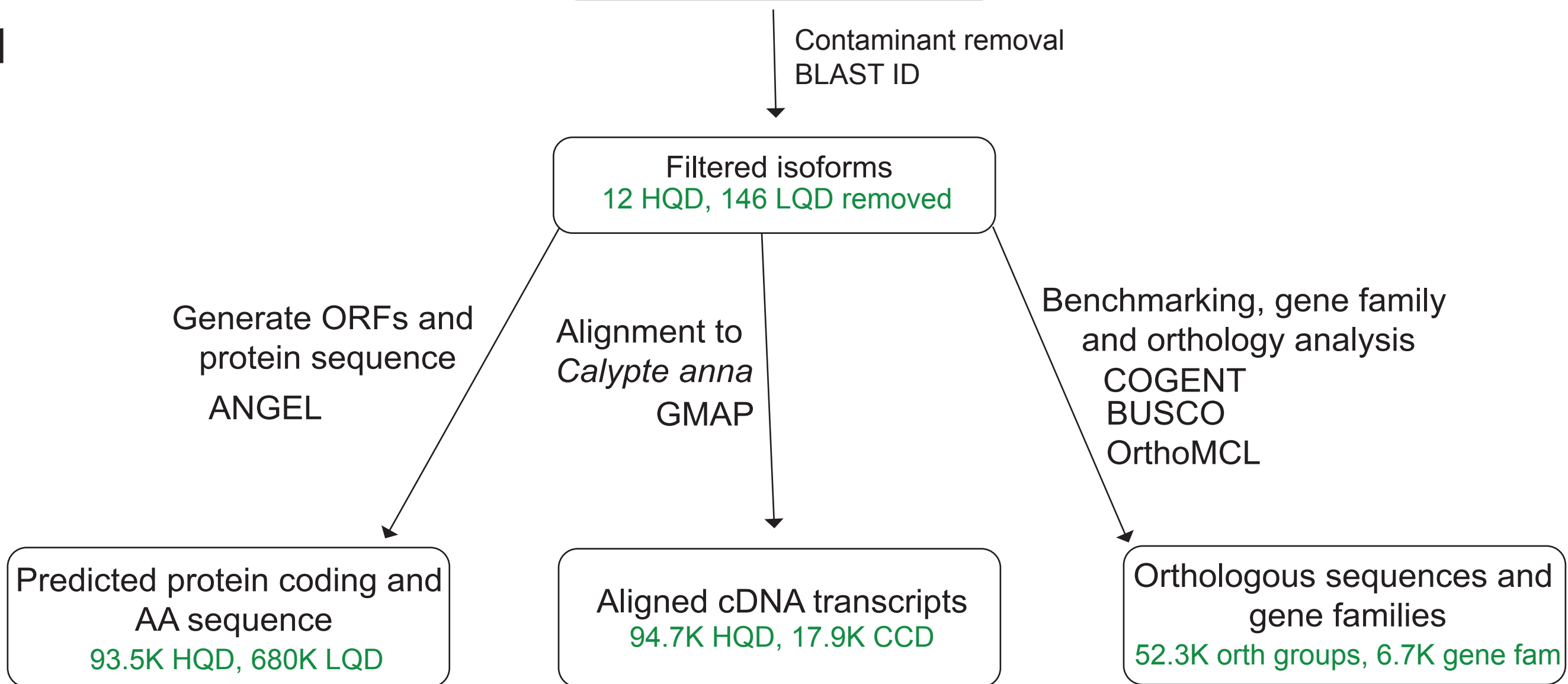
A]



Pacbio SMRT Analysis Isoseq Pipeline

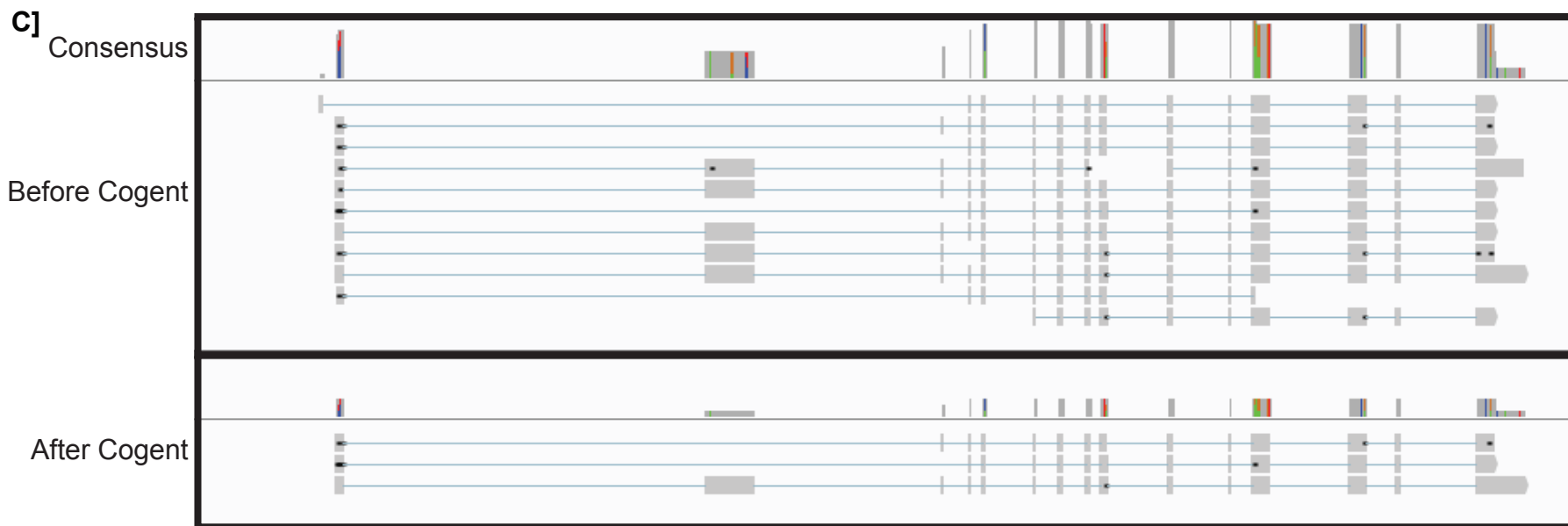
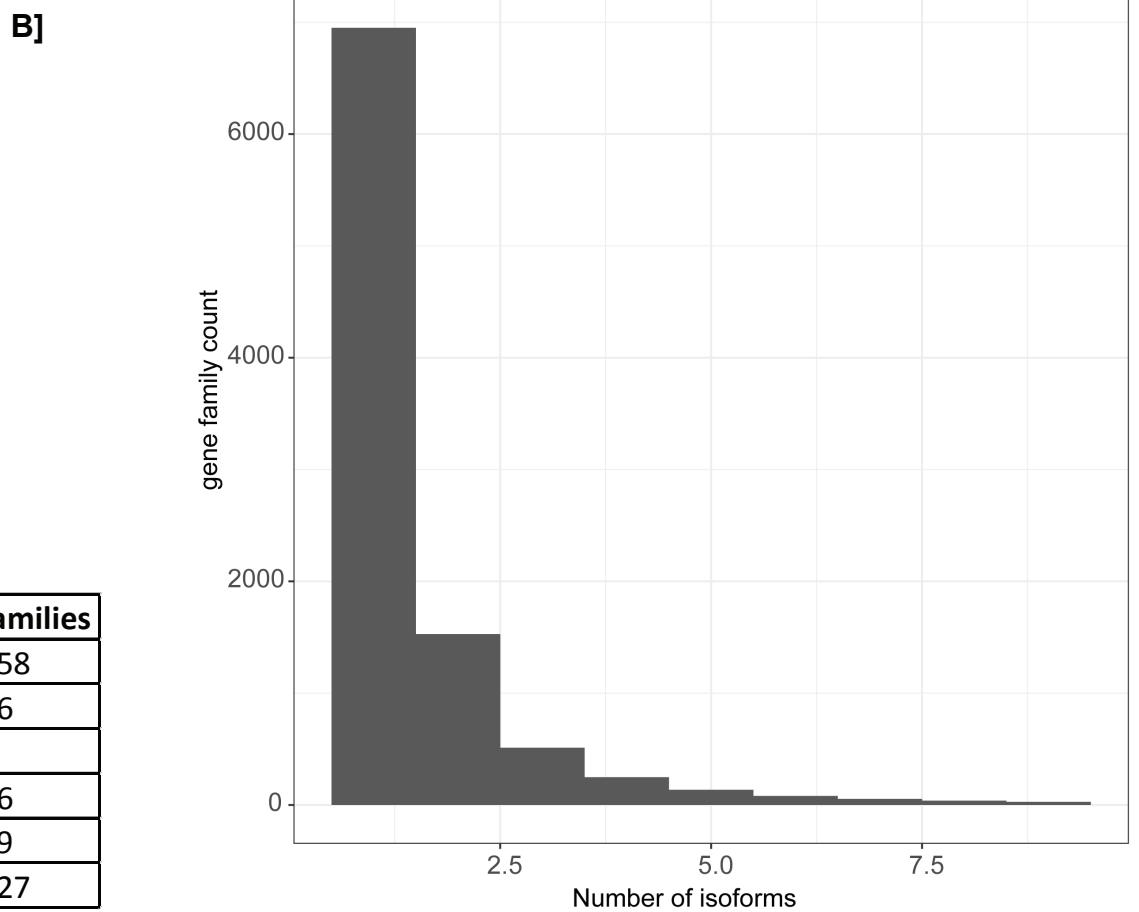
Applications

B]



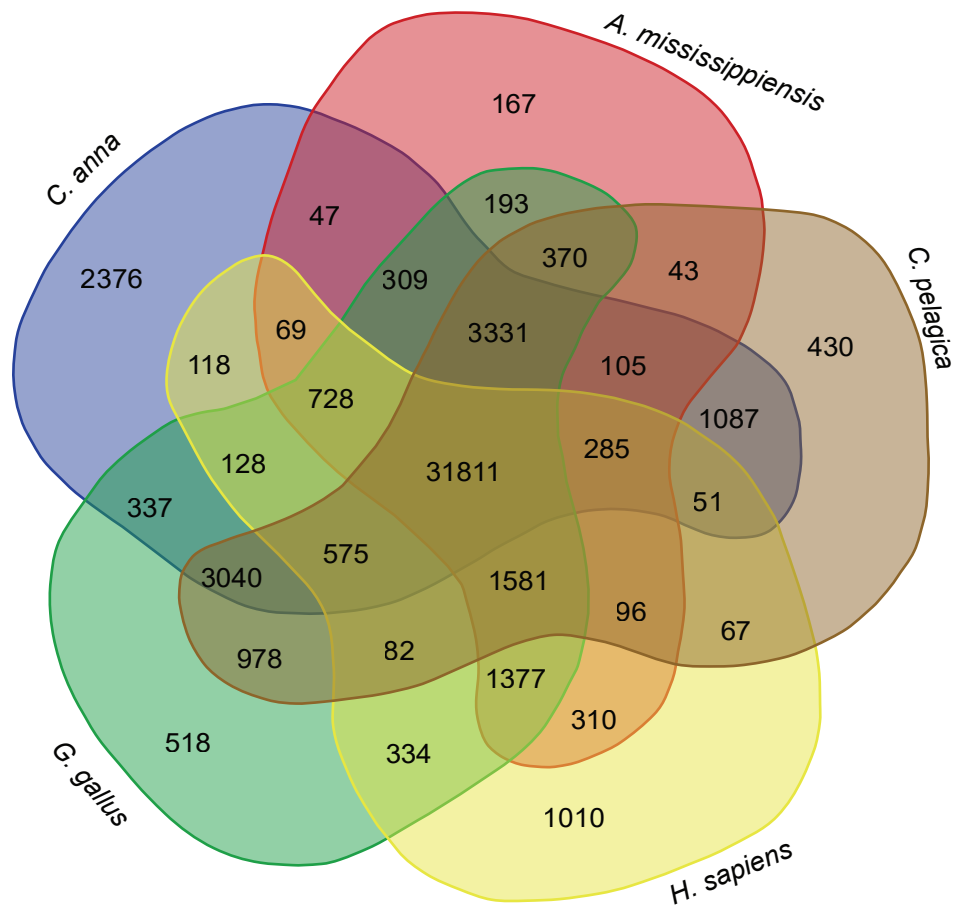
A)

Clustering results	Counts		
Total high quality (HQ) Isoforms	94724		
Total grouped by Cogent	91733		
Orphan seqs (likely single-isoform)	2991		
Gene families predicted	6727		
Gene family alignment: GMAP	Counts	Percent	
Unaligned	1068	5.97%	
Multi-mapped	2614	14.62%	
Uniquely Mapped	15262	85.37%	
qCoverage = 100%	10076	56.36%	
qCoverage >= 99%:	14018	78.41%	
qCoverage >= 90%	14559	81.44%	
Total number transcripts	17877	100.00%	
Cogent comparison cases	In Cogent	In Ref	#families
Single gene locus	1	1	5258
Missing gene, possible broken	1	>1	176
Missing gene	1	0	38
Unresolvable to 1 contig	>1	1	836
Possible multi-loci gene	>1	>1	419
Total gene families			6727

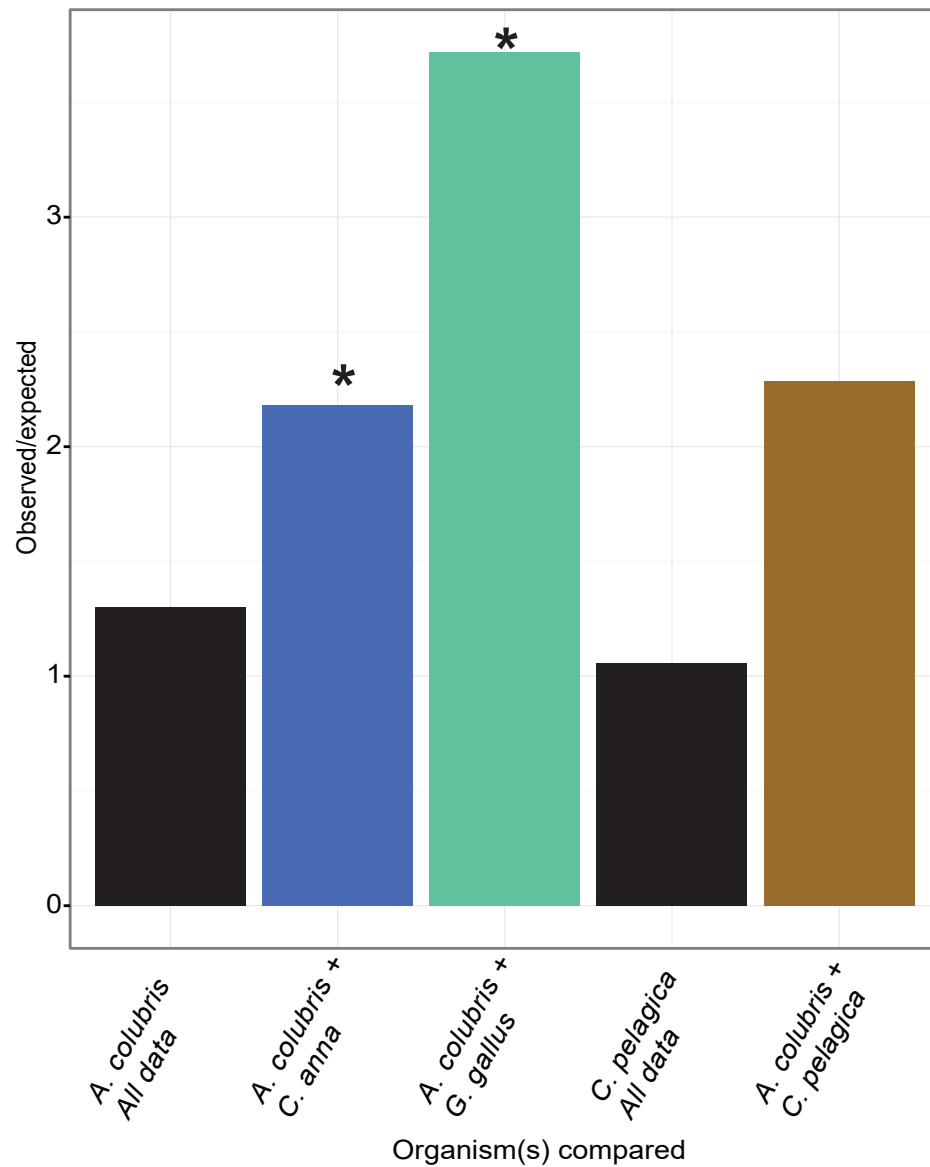


Cogent family 14912
 MATR3 gene
 FALCON assembly
 scaffold 000168F
 860,227-921,621

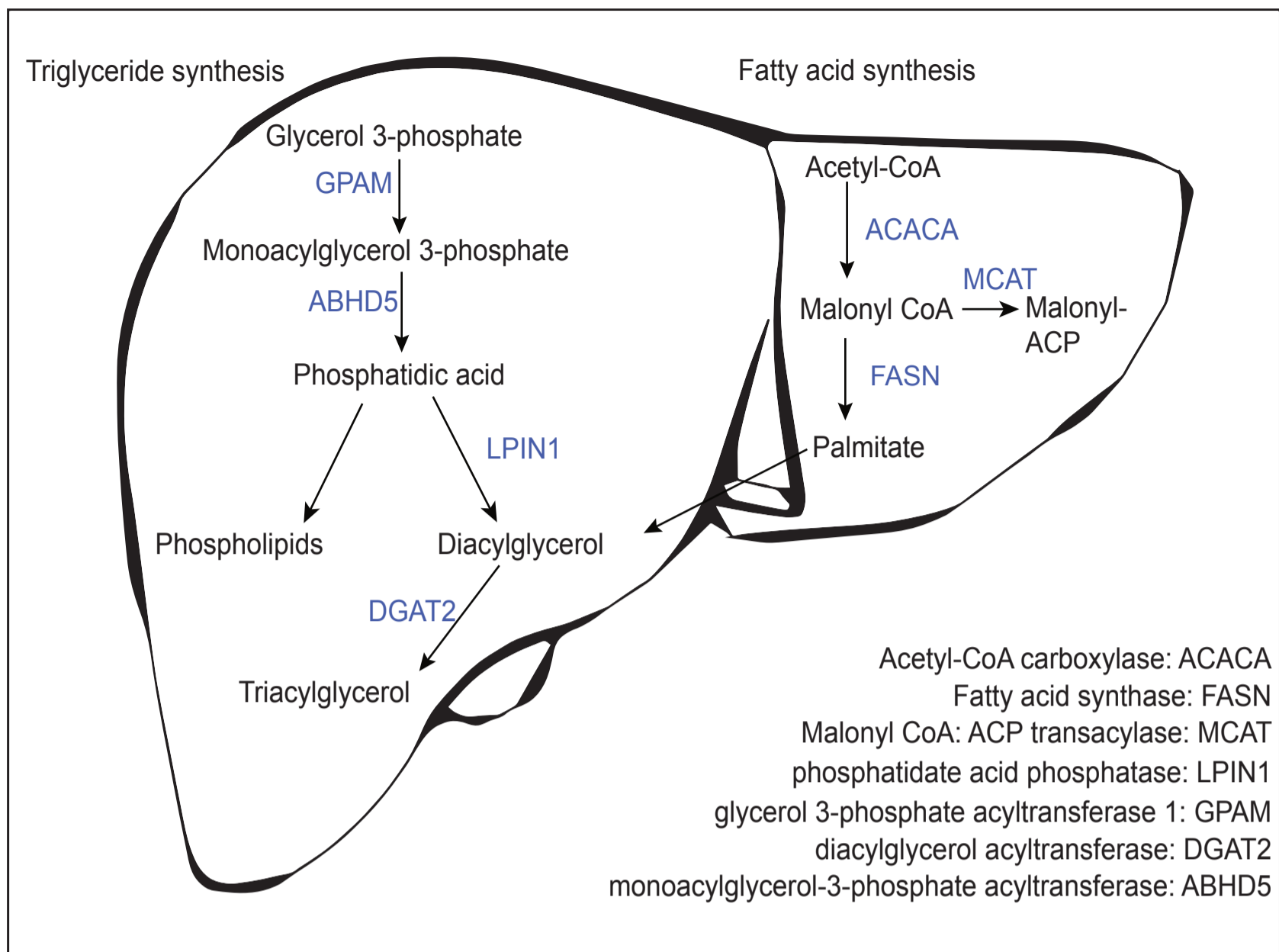
A] OrthoMCL predicted orthologs to *A. colubris*



B] Panther overrepresentation test: Metabolic process proteins abundance



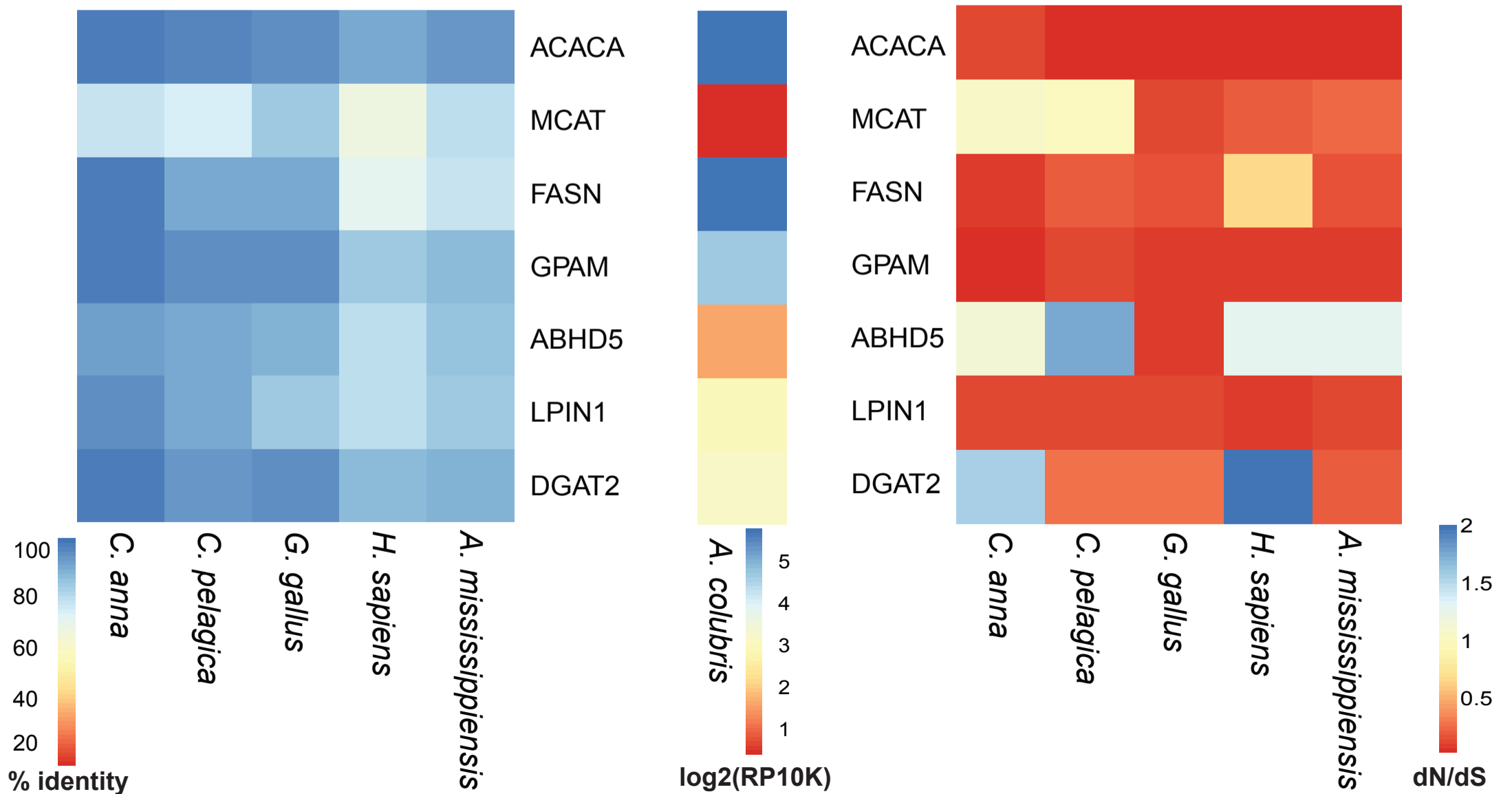
A]



B] Protein alignment

Abundance

Conservation analysis





[Click here to access/download](#)

Supplementary Material

[Supplemental_methods_workman_180102.pdf](#)





Click here to access/download
Supplementary Material
171009_suppdata.pdf





Department of Biomedical Engineering

Clark Hall 118A
3400 N. Charles St.
Baltimore MD 21218
wtimp@jhu.edu
www.timplab.com

Winston Timp
Assistant Professor

1/2/2018

Dear Dr. Zauner:

Below please find our revisions to the manuscript entitled "Single molecule, full-length transcript sequencing provides insight into the extreme metabolism of ruby-throated hummingbird *Archilochus colubris*".

We thank the editors and reviewers for their thoughtful comments on the revision. We have addressed all of the reviewers' concerns and are submitting a detailed response along with an improved version of the manuscript. Specifically:

We believe the missing genes are due to two factors – the fact that this transcriptome is generated from a single tissue, and hence several genes are likely not present and that our filtering to ensure transcript quality is relatively conservative, potentially losing bona fide genes.

We have added a brief section discussing chitinase-like gene expression in the liver of our hummingbird sample.

We also examined quantitation – but feel it's not appropriate on these samples, as we unfortunately were unable to generate Illumina data from the same sample as PacBio data – the Illumina data was from a different individual which may result in altered expression – for example, we believe this could explain the chitinase expression level.

Please feel free to contact us with any further concerns or issues. Thank you for your kind consideration of the manuscript.

Sincerely,

A handwritten signature in black ink that reads "Winston Timp". The signature is written in a cursive, flowing style.

Winston Timp
Assistant Professor
Department of Biomedical Engineering
Johns Hopkins University