

GigaScience

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

Manuscript Number:					
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:	<table border="1"><tr><td>Human Frontier Science Program (#RGP0062/2016)</td><td>Kenneth C. Welch</td></tr><tr><td>Natural Sciences and Engineering Research Council of Canada (CA) (#386466)</td><td>Kenneth C. Welch</td></tr></table>	Human Frontier Science Program (#RGP0062/2016)	Kenneth C. Welch	Natural Sciences and Engineering Research Council of Canada (CA) (#386466)	Kenneth C. Welch
Human Frontier Science Program (#RGP0062/2016)	Kenneth C. Welch				
Natural Sciences and Engineering Research Council of Canada (CA) (#386466)	Kenneth C. Welch				
Abstract:	<p>Hummingbirds can support their high metabolic rates exclusively by oxidizing ingested sugars, which is unsurprising given their sugar-rich nectar diet and use of energetically expensive hovering flight. However, they cannot rely on dietary sugars as a fuel during fasting periods, such as during the night, at first light, or when undertaking long-distance migratory flights, and must instead rely exclusively on onboard lipids. This metabolic flexibility is remarkable both in that the birds can switch between exclusive use of each fuel type within minutes and in that de novo lipogenesis from dietary sugar precursors is the principle way in which fat stores are built, sometimes at exceptionally high rates, such as during the few days prior to a migratory flight. The hummingbird hepatopancreas is the principle location of de novo lipogenesis and likely plays a key role in fuel selection, fuel switching, and glucose homeostasis. Yet understanding how this tissue, and the whole organism, achieves and moderates high rates of energy turnover is hampered by a fundamental lack of information regarding how genes coding for relevant enzymes differ in their sequence, expression, and regulation in these unique animals. To address this knowledge gap, 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, including classification of reads and clustering of isoforms (ICE) followed by error-correction (Arrow). With COGENT, we clustered different isoforms into gene families to generate de novo gene contigs. We performed orthology analysis to identify closely related sequences between our transcriptome and other avian and human gene sets. We also aligned our transcriptome against the <i>Calypte anna</i> genome where possible. Finally, we closely examined homology of critical lipid metabolic genes between our transcriptome data and avian and human genomes. 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 have leveraged cutting-edge technology and a novel bioinformatics pipeline to provide a compelling 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:	
Opposed Reviewers:	
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>	Yes

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

1
2
3
4 Single molecule, full-length transcript sequencing provides insight into the extreme metabolism
5 of ruby-throated hummingbird *Archilochus colubris*
6

7
8 Rachael E. Workman^{1*}, Alexander M. Myrka^{2*}, Elizabeth Tseng⁴, G. William Wong³, Kenneth C.
9 Welch Jr.²⁺, and Winston Timp¹⁺
10

11
12 ¹ Department of Biomedical Engineering, Johns Hopkins University, Baltimore, MD, USA

13
14 ² Department of Biological Sciences, University of Toronto Scarborough, Toronto, Ontario,
15 Canada and Department of Cell & Systems Biology, University of Toronto, Toronto, Ontario,
16 Canada

17
18 ³ Department of Physiology and Center for Metabolism and Obesity Research, Johns Hopkins
19 University School of Medicine, Baltimore, MD, USA

20
21 ⁴ Pacific Biosciences, Menlo Park, California, USA

22
23 * Co-first author

24
25 + Co-Corresponding author
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

Abstract

Hummingbirds can support their high metabolic rates exclusively by oxidizing ingested sugars, which is unsurprising given their sugar-rich nectar diet and use of energetically expensive hovering flight. However, they cannot rely on dietary sugars as a fuel during fasting periods, such as during the night, at first light, or when undertaking long-distance migratory flights, and must instead rely exclusively on onboard lipids. This metabolic flexibility is remarkable both in that the birds can switch between exclusive use of each fuel type within minutes and in that *de novo* lipogenesis from dietary sugar precursors is the principle way in which fat stores are built, sometimes at exceptionally high rates, such as during the few days prior to a migratory flight. The hummingbird hepatopancreas is the principle location of *de novo* lipogenesis and likely plays a key role in fuel selection, fuel switching, and glucose homeostasis. Yet understanding how this tissue, and the whole organism, achieves and moderates high rates of energy turnover is hampered by a fundamental lack of information regarding how genes coding for relevant enzymes differ in their sequence, expression, and regulation in these unique animals. To address this knowledge gap, 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, including classification of reads and clustering of isoforms (ICE) followed by error-correction (Arrow). With COGENT, we clustered different isoforms into gene families to generate *de novo* gene contigs. We performed orthology analysis to identify closely related sequences between our transcriptome and other avian and human gene sets. We also aligned our transcriptome against the *Calypte anna* genome where possible. Finally, we closely examined homology of critical lipid metabolic genes between our transcriptome data and avian and human genomes. 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 have leveraged cutting-edge technology and a novel bioinformatics pipeline to provide a compelling first direct look at the transcriptome of this incredible organism.

Introduction

Hummingbirds are the only avian group to engage in sustained hovering flight as a means for accessing floral nectar, their primary caloric energy source. While hovering, small hummingbirds, such as the ruby-throated hummingbird (*Archilochus colubris*), achieve some of the highest mass-specific metabolic rates observed among vertebrates (Suarez 1992; Chai and Dudley 1996). Given their specialized, sugar-rich diet, it is not that surprising that hummingbirds are able to fuel this intense form of exercise exclusively by oxidizing carbohydrates (Suarez et al. 1990; Chen and Welch 2014). This energetic feat is also remarkable in that the source of sugar oxidized by flight muscles during hovering is the same sugar ingested in nectar meals only minutes prior (Chen and Welch 2014; Welch et al. 2007). In addition, hummingbirds seem equally adept at relying on either glucose or fructose (the two monosaccharides comprising their nectar (Baker 1975) as a metabolic fuel for flight (Chen and Welch 2014). In doing so, they

1
2
3
4 achieve rates of sugar flux through their bodies that are up to 55× greater than non-flying
5 mammals (Welch and Chen 2014).
6

7
8 Hummingbird flight is not always a solely carbohydrate-fueled endeavor. Lipids are a more
9 energy dense form of fuel storage, and fasted hummingbirds are as capable of fueling hovering
10 flight via the oxidation of onboard lipid stores as they are dietary sugars (Welch et al. 2007).
11 Lipids are likely the sole or predominant fuel used during overnight periods (Powers et al. 2003).
12 Just as flux of sugar through the hummingbird is extremely rapid, the building of lipid stores from
13 dietary sugar is also rapid when needed. For example, ruby-throated hummingbirds can
14 routinely increase their mass by 15% or more between midday and dusk on a given day (Hou et
15 al. 2015). The ruby-throated hummingbird (*A. colubris*) completes an arduous annual migratory
16 journey from breeding grounds as far north as Quebec in Canada to wintering grounds in
17 Central America (Weidensaul et al. 2013). Hummingbirds are constrained to fueling long
18 distance migratory flights using onboard lipids. In preparing for such flights, hummingbirds
19 rapidly build fat stores prior to departure or at migratory stopover points, increasing their mass
20 by 25-40% in as few as four days (Carpenter et al. 1993; Hou et al. 2015; Hou and Welch
21 2016).
22
23
24
25
26

27 The ability to switch so completely and quickly between fuel types means these animals
28 possess remarkably exquisite control over rates of substrate metabolism and biosynthesis in the
29 liver, the principal site of lipogenesis in birds (Hermier 1997). While hummingbird liver does
30 indeed exhibit remarkably high activities of lipogenic and other metabolic enzymes (Suarez et
31 al. 1988), the mechanisms underlying high catalytic rates (high catalytic efficiency and/or high
32 levels of enzyme expression) and control over flux (the role of hierarchical versus metabolic
33 control), sensu (ter Kuile and Westerhoff 2001), remain unclear.
34
35
36

37 Despite long-standing recognition of, and interest in, their extreme metabolism, the lack of
38 knowledge about gene and protein sequences in hummingbirds has limited more detailed and
39 mechanistic analyses. Amplification of hummingbird genetic sequences for sequencing and/or
40 cloning is hampered by the lack of sequence information from closely related groups, making
41 well-targeted primer design difficult. Only two genes have thus far been cloned from any
42 hummingbird: an uncoupling protein (UCP) homolog and insulin (Vianna et al. 2001; Fan et al.
43 1993). These two studies offer limited insight into what adaptations in hepatopancreatic
44 molecular physiology underlie extreme energy turnover or its regulation. The UCP homolog was
45 cloned from pectoralis (flight muscle) and its functional significance *in vivo* is unclear. The amino
46 acid sequence of hummingbird insulin was found to be largely identical to that from chicken;
47 however, birds are insulin insensitive and lack the insulin-regulated glucose transporter (GLUT)
48 protein GLUT4, making the role of this hormone in the regulation of energy homeostasis in
49 hummingbirds unknown (Welch et al. 2013; Braun and Sweazea 2008; Polakof et al. 2011).
50
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 (Jarvis et al. 2015). Despite their extreme
59 catabolic and anabolic capabilities, hummingbirds have the smallest genome among birds
60
61
62
63
64
65

1
2
3
4 (Gregory et al. 2009) and, in general, have among the smallest vertebrate genomes (Hughes
5 and Hughes 1995). Thus, it seems likely that understanding of transcriptional variation, overlaid
6 on top of genetic variation, is crucial to understanding what makes these organisms such elite
7 metabolic performers.
8
9

10 To this end, we produced the first high-coverage transcriptome of any single avian tissue, the
11 liver of the ruby-throated hummingbird, *Archilochus colubris*. Because many of the proteins
12 involved in cellular metabolism are quite large, we collaborated with Pacific Biosciences to
13 generate long-read sequences as these would enhance our ability to identify full coding
14 sequences and multiple encoded isoforms. The primary advantage to the Pacbio Iso-seq
15 methodology is the capability for full-length transcript sequencing, rendering complete mRNA
16 sequences without the need for assembly. This has been demonstrated in previous studies to
17 drastically increase detection of alternative splicing events (Abdel-Ghany et al. 2016).
18 Additionally, full-length sequences greatly enhance the likelihood of detecting novel or rare
19 splice variants, which is crucial for fully characterizing the transcriptomes of lesser studied,
20 non-model organisms such as the hummingbird.
21
22
23
24
25

26 **Results**

27 *First high coverage single-tissue avian transcriptome quality control and validation*

28
29
30 The transcriptome described in this manuscript represents the first long read, single tissue avian
31 transcriptome completed at high coverage. Four size fractions (1-2kb, 2-3kb, 3-6kb, 5-10kb) of
32 sample were sequenced on 40 SMRT Cells, producing 440.75 Gb of raw data, 3.4Gb of
33 full-length non-chimeric reads after circular consensus sequence (CCS) generation and filtering
34 for full-length read classification. Of the four size-selected bins, our average CCS length was
35 1533, 2464, 3650, and 5444 bp, respectively (Figure 1B).
36
37
38

39 To analyse these data, we employed the Pacbio sa3_ds_iseq pipeline (SMRTAnalysis 3.1.0,
40 (Minoche et al. 2015; Gordon et al. 2015), code here:
41 <https://github.com/PacificBiosciences/SMRT-Analysis>) using the DNANexus cloud computing
42 platform. A summary of analyses performed using both the SMRTAnalysis pipeline and
43 downstream are displayed in Figure 2A-B. Using this pipeline, we classified 3.4Gb of
44 non-chimeric CCS reads into 1,236,437 full-length (48%) and 1,665,929 non-full length reads
45 where reads were determined to be full-length if both the 5' and 3' cDNA primers as well as the
46 polyA tail signal were detected. The Iso-Seq pipeline then performed isoform-level clustering
47 (ICE) followed by final polishing using the Arrow algorithm to output high-quality (predicted
48 accuracy $\geq 99\%$), full-length, isoform consensus sequences. The Iso-Seq pipeline produced
49 238Mb (807,104 reads) of high quality consensus isoforms (HQD, 94,724 reads), and 2Gb
50 (712,210 reads) of low quality consensus isoforms (summary statistics Figure 1A).
51
52
53
54
55

56 To determine completeness of our transcriptome assembly relative to established
57 transcriptomes, as well as relative retention of transcripts between multiple data processing
58 steps, we used benchmarking software BUSCO (Benchmarking Universal Single Copy
59 Orthologs). BUSCO searches for a list of conserved orthologous genes assumed to be present
60
61
62
63
64
65

1
2
3
4 in all completed transcriptome assemblies for members of the given clade. We utilized the
5 Metazoan and Aves lineage datasets (Simão et al. 2015) to determine assembly completeness
6 for not only our ASD (all-sequence data) dataset, but also the HQD (High quality dataset) and
7 COGENT-collapsed dataset (CCD, described below) in order to ensure that sequence diversity
8 was not lost when filtering. We then compared this result with the chicken transcriptome (*Gallus*
9 *gallus*) (ftp://ftp.ncbi.nih.gov/genomes/Gallus_gallus), the predicted transcriptome from the
10 recently published Anna's hummingbird (*Calypte anna*) genome
11 (<http://gigadb.org/dataset/101004>), and *Gallus gallus* from a single-tissue Pacbio Iso-Seq
12 dataset (<http://journals.plos.org/plosone/article?id=10.1371/journal.pone.0094650>). Busco v2.0
13 BETA was used (downloaded 9/1/2016, <https://gitlab.com/ezlab/busco>), and results are
14 summarized in **Figure 2C and Supplemental Table 1**. Notably, the NCBI *Gallus gallus*
15 transcriptome is nearly complete for both Aves and Metazoan sets, with the predicted
16 transcriptome from *Calypte anna* not far behind. In contrast, our transcriptome contains only
17 about half of the Aves set and slightly more of the metazoan set; many genes are not
18 substantially expressed in a single tissue. Notably our ortholog detection is dramatically
19 improved over the other single-tissue Pacbio data available, (Thomas et al. 2014), which only
20 captured 6% of the Aves set.

21
22
23
24
25
26
27
28 To estimate the completeness of our liver transcriptome sequencing, we took subsets of the
29 circular consensus reads dataset and BLASTed against the predicted *Calypte anna* gene set.
30 We found that the number of unique genes detected began to saturate when reaching a 90%
31 subset of our data, suggesting that we had near complete transcriptome sequencing
32 (**Supplemental Figure 1**). Lower expressed genes may not be detected, but that vast majority of
33 liver expressed genes are likely represented in our data.

34 35 36 37 *Agreement with established Anna's hummingbird genome reveals general clade conservation*

38
39 We aligned transcripts to the *Calypte anna* (Anna's hummingbird) genome using both all
40 consensus isoforms (ASD) and high quality isoforms only (HQD) (Korlach et al. 2017), as well
41 as reads collapsed by COGENT (CCD, methods detailed below). As *Calypte anna* is a close
42 relative within the same Bee clade of hummingbirds (McGuire et al. 2007), we expected
43 alignment to perform well. We found an average alignment identity of 94.8%, with 87%
44 transcripts uniquely mapping to the reference. Of the uniquely mapped, 73% covered >90% of
45 the query sequence (alignment length and statistics, **Supplemental Figure 2A, 2B**),
46 demonstrating high fidelity of aligned reads to reference. When ASD reads were parsed by
47 number of reads of insert supporting each consensus cluster, it was found that generally,
48 alignment identity was high regardless of number of supporting reads. A clear increase in mean
49 alignment identity was found when two or more supporting reads were collapsed (**Supplemental**
50 **Figure 3**).

51
52
53
54
55
56 When GMAP was performed using only high quality isoforms (filtered for 2+ full-length
57 supporting reads), alignment percentage was 95.7%, with 93.4% of transcripts mapping
58 uniquely to the reference. The average mapped read length was 2411bp (HQD, 2617bp ASD),
59 while the average predicted CDS length for *Calypte anna* was 1386bp. This being said, reads
60
61
62
63
64
65

1
2
3
4 mapped with GMAP contain UTRs. When we predict just the CDS sequences for *A. colubris*
5 using ANGEL (<https://github.com/PacificBiosciences/ANGEL>), the mean length was 981bp.
6 When we BLASTed the unaligned reads to whole NCBI database, they largely mapped back to
7 *Calypte anna* (53%). This result suggests that our mapping parameters were too stringent to
8 map these reads, error rate prevented alignment, unaligned regions are divergent enough
9 between both hummingbirds to preclude alignment, or some combination of the above.
10
11

12 *Putative gene family prediction and reduction of transcript redundancy reduces data load while* 13 *maintaining transcript diversity* 14 15

16
17 To assign transcripts to putative gene families, as well as cluster and eliminate redundant
18 transcripts to produce a unique set of gene isoforms, we utilized the newly developed COGENT
19 (<https://github.com/Magdoll/Cogent>, Liz Tseng, pre-print pending) pipeline. COGENT is
20 specifically designed for transcriptome assembly in the absence of a reference genome,
21 allowing for isoforms of the same gene to be distinctly identified from different gene families,
22 which are defined as having more than two (possibly redundant) transcript copies. Of the 94,724
23 HQ consensus isoforms, 91,733 were grouped into 6,725 gene families (Figure 3A). The
24 remaining 2,991 sequences were classified as putative single-isoform genes, or “orphans”.
25 Reconstructed contigs were then applied in place of a reference (or de novo clustering) to
26 reduce redundant transcripts in the original HQD dataset. From this approach, we were able to
27 reduce our HQ dataset to 14,628 distinct transcript isoforms and 2990 orphan isoforms, for a
28 total of 17,618 isoform sequences (18% of the original). Data is further summarized in
29 Supplemental Table 2. An average of 1.53 isoforms was found per gene family (Figure 3B), with
30 2624, or 27.4% of the gene families having more than one isoform, including “orphans”. While
31 other studies have found more isoforms per locus, for example 6.56 in (Wang et al. 2016), that
32 study multiplexed six plant tissues, whereas a lower complexity is to be expected with single
33 tissue analysis. This dataset (COGENT collapsed data, or CCD) was also mapped onto the
34 *Calypte anna* genome assembly (<http://gigadb.org/dataset/101004>), to demonstrate the
35 effectiveness of this method in reducing transcript redundancy and classifying isoforms (Figure
36 3C).
37
38
39
40
41
42
43

44 *Orthologous gene pair predictions and GO annotation show putative unique hummingbird* 45 *orthologs* 46 47

48 To examine protein sequence similarity and divergence between *Archilochus colubris* and other
49 avian species, we used OrthoMCL, which generates reciprocal best hits from comparator
50 species using BLAST all-vs-all, then clustering to group orthologous sequences for each pair of
51 organisms (Li et al. 2003). OrthoMCL protein sequences were predicted using ANGEL, and
52 119,292 high quality sequences were put into this analysis. We compared our ruby-throated
53 hummingbird, *Archilochus colubris*, to five other birds: *Calypte anna* (Anna’s hummingbird)
54 fellow member of the bee clade of hummingbirds, *Chaetura pelagica* (chimney swift) the closest
55 available outgroup species to the hummingbird clade, and other bird species for which relatively
56 well-annotated genomes and/or transcriptomes are available, *Gallus gallus* (chicken),
57
58
59
60
61
62
63
64
65

1
2
3
4 *Taeniopygia guttata* (zebra finch), and *Melopsittacus undulatus* (budgerigar), as well as *Homo*
5 *sapiens* (human), and *Alligator mississippiensis* (American alligator).
6

7
8 A matrix of ortholog pairings, with duplicate ortholog hits removed, shows counts of number of
9 orthologous sequences for each species pair (Supplemental Table 3). Orthologs shared
10 between ruby-throated hummingbird and a subset of the other species analyzed are illustrated
11 in Figure 4A. Unsurprisingly, the largest amount of orthologs which pair closely to only one
12 species, i.e., 1:1 orthologs, were found between Anna's and Ruby-throated hummingbird
13 sequences. Surprisingly, the second-largest set was between chicken and ruby-throated
14 hummingbird, as opposed to its closest outgroup species, *Chaetura pelagica*. This is likely due
15 to the completeness of chicken transcriptome annotation, as chicken is the most well-studied
16 avian species. Of the 596 unpaired *A. colubris* protein sequences, 190 paired most closely with
17 *Calypte anna* when compared using BlastP and the majority of matches output (559/594) were
18 less than 50 AA, only a fraction of the average sequence length.
19
20
21
22

23 In order to more closely examine the identity of 1:1 orthologs in related hummingbird species,
24 gene ontology (GO) annotation was performed on a specific set of orthologs which were shared
25 between *Calypte anna* and *Archilochus colubris*, but not by the other birds included in the
26 OrthoMCL analysis. This set of 2,376 protein sequences was run using BlastP and GO analysis
27 performed by Panther (Mi et al. 2013, 2017). Additional datasets used for GO comparison
28 included 1:1 orthologs for *Gallus gallus* and *A. colubris* (518), and *A. colubris* and *Chaetura*
29 *pelagica* (430), as well as whole transcriptome data from *C. pelagica* and COGENT-collapsed
30 dataset from our transcriptome (Supplemental Table 4, Figure 4B).
31
32
33
34

35 As the initial impetus for our investigation centered on the exceptional metabolism and
36 energetics of hummingbirds, we focused our investigation on orthologs tagged as part of the
37 "metabolic process (GO:0008152)" grouping. Of the 1444 orthologs identified in *Archilochus*
38 *colubris* as part of this process grouping, 236 (16.3%) were unique to hummingbirds. Within this
39 top-level grouping, the largest number of genes group under "primary metabolic processes"
40 (GO:0044238)". Of the 1240 orthologs identified within this grouping, 204 (16.3%) are identified
41 as uniquely shared by our hummingbird species. Six GO biological processes are defined under
42 the "primary metabolic processes". Of these processes, the process with the highest proportion
43 of identified *A. colubris* orthologs hitting as unique to the two hummingbird species is "lipid
44 metabolic processes" (GO:0006629; 33 of 114 orthologs, 28.9%), which is significantly enriched
45 relative to the comparative orthology databases of both chicken and human. Because we
46 considered it likely that an enrichment in lipid metabolic genes could be a result of our dataset
47 being from liver tissue, we compared enrichment with that of the entire COGENT predicted gene
48 set from the ruby-throated hummingbird transcriptome, and found no significant enrichment.
49 This suggests a higher degree of divergence within this class of enzymes than would be
50 predicted statistically.
51
52
53
54
55
56

57 The proportion of identified genes within a biological pathway classified as orthologs unique to
58 hummingbirds should not be taken as direct evidence of greater selection on proteins within that
59 pathway. Yet, if neutral sequence divergence is assumed to be randomly accrued throughout a
60
61
62
63
64
65

1
2
3
4 species' genome, then greater divergence in enzymes making up "lipid metabolic processes"
5 suggests that closer examination of these proteins for evidence of functional, or even adaptive,
6 divergence is warranted. A phylogenetically-informed analysis of ortholog divergence among
7 taxa is necessary to establish a selection signature, which will become possible in the future
8 with the advance of the B10K project (Zhang et al. 2015) and larger numbers of avian species in
9 GO databases.
10
11

12 13 *Hepatic lipogenesis case study demonstrates utility of long-read transcriptomics data in* 14 *evaluating biology* 15

16
17 Given the apparent sequence divergence among enzymes involved in "lipid metabolic
18 processes" hinted at by orthology and ontology analyses, we elected to more closely examine
19 sequence divergence in enzymes comprising the lipogenic pathway. Lipogenesis, the process
20 by which fatty acids are produced from acetyl-CoA and subsequently triglycerides are
21 synthesized, involves several key enzymes which we examined in closer detail. We predicted
22 that this pathway (Figure 5A) would be divergent in hummingbirds given their extraordinary
23 metabolic demands. Eight enzymes involved in this pathway were examined for *Archilochus*
24 *colubris*, *Calypte anna*, *Gallus Gallus*, *Chaetura pelagica*, *Alligator mississippiensis* and *Homo*
25 *sapiens* (accession numbers and details given in Supplemental Table 5). Pairwise protein
26 alignment scores are given in Supplemental Table 6 as well as illustrated in a heatmap shown in
27 Figure 5B, and alignments in Supplemental Data 1. Interestingly, enzymes with higher identity
28 between examined organisms are involved in the fatty acid synthesis arc of metabolism, while
29 triglyceride synthesis enzymes tend to be less conserved (Figure 5A). While poor pairwise
30 protein alignment between *A. colubris* and all examined species, such as with DGAT2, is
31 suggestive of misannotation or splice variation in our transcriptome, cases with variable
32 alignment identities provide interesting targets for further investigation.
33
34
35
36
37
38

39 In order to further investigate degree of conservation between key metabolic enzymes in
40 hummingbirds and comparator organisms, we performed conservation analysis and determined
41 ratio of nonsynonymous to synonymous codon changes (dN/dS) as a metric of positive
42 selection, using pairwise alignments followed by codeml module in PAML4 (Yang 2007). These
43 ratios are given in Supplemental Table 6. We found general conservation of these enzymes
44 between organisms, with the exception of the 3' and 5' ends of alignments. These often had an
45 extended or retracted coding sequence in the case of hummingbirds and *C. pelagica*, which
46 could potentially be post-translational modification or selection on pathway regulation (Jacob
47 and Unger 2007). Surprisingly, terminal sequence length was variable even between *C. anna*
48 and *A. colubris*, which both belong to the closely-related Bee hummingbird taxid (McGuire et al.
49 2014). Variation in 5' and 3' length may also be an effect of the different methodologies used to
50 produce these sequences, RNA sequencing for *A. colubris* and *G. gallus*, and ORF prediction
51 from genomic data for the other organisms examined.
52
53
54
55
56

57 The averaged dN/dS values, while useful for comparison, are misleading when considered over
58 the entire gene, as 3' and 5' variation can overshadow conserved motifs and vice versa, and
59 pairwise comparisons are limited in scope. In addition, this type of analysis is ideal for very
60
61
62
63
64
65

1
2
3
4 divergent sequences, and less informative for pairs of sequences that are highly similar
5 (Kryazhimskiy and Plotkin 2008). Despite this, conservation analysis is still valuable and
6 provides relative conservation metrics between these organisms for specific enzymes in a
7 pathway of interest, as well as enabling identification of specific regions of high variation within
8 coding sequences. Additionally, pairwise comparisons provide interesting observations, such as
9 coding strand elongation in the 5' region in *A. colubris* GPAM and GPAT4 (AGPAT6)
10 (Supplemental Data 2). This information will be carried forth into future studies more closely
11 examining enzyme structure, function and evolution.
12
13
14

15
16 Access to the transcriptome informs the investigation of biological processes and enables the
17 formation of new hypotheses. This is exemplified by the serendipitous observation that
18 hummingbird glucose transporter 2 (GLUT2) lacks a N-glycosylation site due to an asparagine
19 to aspartic acid amino acid substitution. This missing glycosylation site was also seen in the
20 available Anna's hummingbird genome. All class 1 glucose transporters studied in model
21 vertebrates contain one N-glycosylation site located on the large extracellular loop of the protein
22 (Joost and Thorens 2001). In GLUT2 the associated glycan interacts with the glycan-galectin
23 lattice of the cell, stabilizing cell surface expression (Ohtsubo et al. 2013). Removal of the
24 N-glycan of GLUT2 in rat pancreatic β cells results in the sequestering of cell-surface GLUT2 in
25 lipid rafts and this sequestered GLUT2 exhibits a reduction in glucose transport activity by
26 approximately 25% (Ohtsubo et al. 2013). This reduction in transport is thought to occur through
27 interaction of the GLUT with lipid raft-bound stomatin (Ohtsubo et al. 2013; Zhang et al. 2001).
28 In mammals, GLUT2 serves a glucose-sensing role in the pancreatic β cells and is required for
29 the regulation of blood glucose through insulin and glucagon (Thorens and Mueckler 2010). The
30 lack of N-glycosylation of GLUT2 may contribute to observed high blood glucose concentration
31 in hummingbirds (Beuchat and Chong 1998).
32
33
34
35
36
37

38 **Conclusions**

39
40 Our results have leveraged cutting-edge technology to provide a compelling first direct look at
41 the transcriptome of this incredible organism. By using PacBio sequencing, we have been able
42 to generate full length cDNA transcripts from the hummingbird liver. We subjected this data to a
43 new and innovative pipeline, including steps to generate ORFs, merge transcripts to reconstruct
44 gene families (COGENT), determine orthologs (OrthoMCL) and even measure evolutionary
45 conservation and analyze genes.
46
47
48

49 Our data is remarkably representative of the transcriptome, capturing ~75% of metazoan
50 universal orthologs, and ~50% of aves universal orthologs. This is an amazing diversity from a
51 single tissue, given that many tissue specific genes are expressed at low levels if at all in liver.
52 By subsetting our data, we show a saturation of the number of unique genes detected,
53 suggesting that our data represents a nearly complete transcriptome of the hummingbird liver.
54 Our data show a high degree of alignment to the existing *Calypte anna* dataset; even the reads
55 which did not align with GMAP were subsequently mostly found to have a best match to *C. anna*
56 with BLAST.
57
58
59
60
61
62
63
64
65

1
2
3
4 *Long read cDNA sequencing allows for assembly-free, less ambiguous isoform detection*
5

6 Previous studies have shown that in a non-model organism, long read RNA sequencing
7 improves detection of alternative splicing events nearly five-fold (Abdel-Ghany et al. 2016).
8 Using full-length transcript data, we found alignment unnecessary to generate clear pictures of
9 the gene isoforms and COGENT reconstructed gene contigs. The long reads negate the need
10 for transcript assembly, a precarious analysis in the absence of a genome. The longest read
11 identified was acetyl-coA carboxylase 1, with read length greater than 7kb, capable of spanning
12 the length of the coding sequence. We have annotated our COGENT de novo transcriptome
13 with the closest equivalent gene from NCBI nr/nt database (O'Leary et al. 2016) which, in
14 combination with GO and OrthoMCL analyses, has improved transcriptome characterization.
15 Although absolute functional assignments using GO annotation would not be ideal for this
16 non-model species, comparing relative abundance of characterized orthologs between different
17 datasets is useful for establishing a framework for forthcoming investigations.
18
19
20
21
22

23 *Transcriptome data provides insight into biological function of hummingbird liver*
24

25 Often non-model organism protein sequences available in public repositories are translated
26 from coding sequences predicted from whole genome data, not transcriptome data. While
27 powerful, this approach for elucidating a transcriptome does not provide information regarding
28 tissue-specific transcription, relative transcript abundance or isoform prevalence. In contrast,
29 whole transcript mRNA sequencing from the liver tissue of *Archilochus colubris* has allowed us
30 to clearly annotate gene families and the specific isoforms expressed in liver. In follow up
31 studies, we will be able to compare the isoforms expressed between the liver and other tissues,
32 e.g. pectoralis muscle.
33
34
35
36

37 *Orthology and gene ontology analysis give clues underlying uniqueness of extreme metabolism*
38

39 Using our transcriptome, we have identified genes with unique orthologs in hummingbird as
40 compared to other bird species, reptiles or even mammals. These genes showed a specific
41 enrichment for pathways involved in lipid metabolism - suggesting that the hummingbird has
42 evolved variants of these genes to achieve its high levels of metabolic efficiency.
43
44

45 *Polished transcriptome provides basis for future genomic and biological studies*
46

47 Transcriptome data generated using the Iso-seq methodology, when coupled to sophisticated
48 recently developed gene synthesis techniques (Kosuri and Church 2014), allows for simple
49 generation of relevant isoforms for biochemical experiments. Some of the key metabolic
50 enzymes identified from our work as being unique to either *A. colubris* or at most common to *C.*
51 *anna* and *A. colubris* could be quickly cloned and expressed. Follow up studies will allow for
52 biochemical studies of proteins generated directly from our transcriptome data, measuring their
53 enzymatic properties, e.g. k_{cat} or V_{max} , as compared to other avian or mammalian analogues
54 (Suarez et al. 2009; FernándezM.J. et al. 2011; Suarez et al. 1988). Expressed proteins may
55 also be used for structural biology studies, applying either x-ray crystallography or cryoEM to
56 generate structural maps of the proteins, and examine how the structure compares to other
57
58
59
60
61
62
63
64
65

1
2
3
4 analogues. Importantly, most of the glycolytic (hexokinase, phosphofructokinase) and lipogenic
5 (acetyl-CoA carboxylase and pyruvate carboxylase) enzymes are highly evolutionarily
6 conserved, and structural information exists for most (Lasso et al. 2014; Xiang and Tong 2008;
7 Zhang et al. 2003; Aleshin et al. 1998; Kamata et al. 2004; Mulichak et al. 1998).
8
9

10 **Methods**

11 *Sacrifice and RNA extraction*

12
13
14
15 Wild adult male ruby-throated hummingbirds (*Archilochus colubris*) were captured at the
16 University of Toronto Scarborough using modified box traps. Birds were housed in the
17 University of Toronto Scarborough vivarium and fed NEKTON-Nectar-Plus (Nekton, Tarpon
18 Springs, FL, USA) ad libitum. Birds were sacrificed after ad libitum feeding, and tissues were
19 sampled immediately after euthanization using RNase-free tools. A hepatopancreas tissue
20 sample was collected from one bird. Tissue was homogenized at 4°C in 1 ml cold Tri Reagent
21 using an RNase free glass tissue homogenizer and RNase free syringes of increasing needle
22 gauge. We used 100 mg of tissue per 1 ml of Tri Reagent (Sigma-Aldrich, St. Louis, MO, USA),
23 and chloroform extraction was performed twice to ensure quality. RNA was precipitated,
24 centrifuged down, washed with ethanol, vacuum dried and eluted in RNase free water. DNase I
25 digestion and spin column cleanup were performed. RNA concentration and RIN were
26 determined with RNA Bioanalyzer (Agilent).
27
28
29
30

31 *Sample preparation and sequencing*

32
33
34 Pacific Bioscience's Iso-Seq sequencing protocol was followed to generate sequencing libraries
35 (Thomas et al. 2014). The Clontech SMARTER cDNA synthesis kit with Oligo-dT primers was
36 used to generate first and second-strand cDNA from polyA mRNA. After a round of PCR
37 amplification, the amplified cDNA was size selected into 4 size fractions (1-2kb, 2-3kb, 3-6kb,
38 and 5-10kb) to prevent preferential small template sequencing, using the Blue Pippin (Sage
39 Sciences). Additional PCR cycles were used post size-selection to generate adequate starting
40 material, and then SMRTbell hairpin adapters were ligated onto size-selected templates. Each
41 of the 4 size fractions was sequenced on 10 SMRT Cells, for a total of 40 SMRT Cells.
42 Sequencing was performed by the JHU HiT Center using P6-C4 chemistry on the RSII
43 sequencer.
44
45
46
47

48 **Analysis Methods**

49 *Data processing, isoform clustering and sorting using the DNANexus Pipeline*

50
51
52
53 The SMRTanalysis 3.1 software (<https://github.com/PacificBiosciences/SMRT-Analysis>) and
54 IsoSeq pipeline were employed using a DNANexus interface. Raw sequence files produced
55 from the Pacbio RSII (bax.h5, bas) were converted into BAM files using bax2bam,
56 zmws_per_split 3. Circular consensus sequence (CCS) was generated from subread BAM files,
57 parameters: min_length 300, max_drop_fraction 0.8, no_polish TRUE, min_zscore -9999,
58 min_passes 1, min_predicted_accuracy 0.8, max_length 15000. CCS.BAM files were output,
59
60
61
62
63
64
65

1
2
3
4 which were then classified into full length and non-full length reads using pbclassify.py,
5 ignorepolyA false, minSeqLength 300. Non-full length and full-length fasta files produced were
6 then fed into the cluster step, which does isoform-level clustering (ICE), followed by final Arrow
7 polishing, hq_quiver_min_accuracy 0.99, bin_by_primer false, bin_size_kb 1, qv_trim_5p 100,
8 qv_trim_3p 30.
9

10 11 *Aligning to reference using GMAP*

12 We aligned with GMAP (Wu and Watanabe 2005) version 2016-09-23 with parameters -f samse
13 -n 0 -z senseforce against Calypte anna genome (Zhang et al. 2014).
14

15 16 17 *Assessing transcriptome completion using BUSCO*

18 BUSCO v2.0 BETA (<https://gitlab.com/ezlab/busco>, (Simão et al. 2015), accessed 9/8/2016)
19 was used to benchmark transcriptome completion, by checking for essential single copy
20 orthologs which should be present in a whole transcriptome dataset for any member of the
21 given lineage. We used both Metazoan and Aves lineages (ortholog sets) to examine
22 transcriptome completion.
23
24
25
26

27 28 *Gene family prediction and reducing transcript redundancy using COGENT*

29 COGENT (coding genome reconstruction tool) v1.3 (<https://github.com/Magdoll/Cogent>) uses
30 k-mer similarity profiles in order to partition full length coding sequences into gene families, after
31 which it reconstructs contigs containing the full coding region.
32
33
34

35 36 *Orthologous gene prediction using OrthoMCL*

37 To examine protein sequence similarity and divergence between *Archilochus colubris* and other
38 avian species, we used OrthoMCL (<http://genome.cshlp.org/content/13/9/2178.full>). OrthoMCL
39 works by performing an all-vs-all BlastP comparison of input sequences, and determining
40 reciprocal best hit of all input pairings (cut off at e-5 and >50% match). Putative orthologs are
41 the reciprocal best hits between species, while putative paralogs constitute reciprocal better hits
42 within species. A normalized similarity matrix, followed by Markov clustering, produces ortholog
43 groups. We compared our ruby-throated hummingbird *Archilochus colubris* to five other birds-
44 *Calypte anna* (Anna's hummingbird, release date 2014-04-24,
45 ftp://climb.genomics.cn/pub/10.5524/101001_102000/101004/Calypte_anna.pep.gz),
46 *Gallus gallus* (chicken, Galgal5, ftp://ftp.ncbi.nih.gov/genomes/Gallus_gallus/protein/protein.fa.gz),
47 *Chaetura pelagica* (chimney swift, release date
48 2014-04-24,ftp://ftp.ncbi.nih.gov/genomes/Chaetura_pelagica/protein/protein.fa.gz),
49 *Taeniopygia guttata* (zebra finch, taeGut3.2.4,
50 ftp://ftp.ensembl.org/pub/release-85/fast/taeniopygia_guttata/pep/Taeniopygia_guttata.taeGut3
51 .2.4.pep.all.fa.gz), and *Melopsittacus undulatus* (budgeriger, melUnd1,
52 ftp://climb.genomics.cn/pub/10.5524/100001_101000/100059/BGIMUN1.120628.gene.withUTR
53 pep), as well as human (*Homo sapiens*, hg38,
54 ftp://ftp.ensembl.org/pub/release-85/fast/homo_sapiens/pep/Homo_sapiens.GRCh38.pep.all.fa
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 .gz) and American alligator (*Alligator mississippiensis*, allMis0.2/1,
5 ftp://ftp.ncbi.nih.gov/genomes/Alligator_mississippiensis/protein/protein.fa.gz).
6
7

8 *GO analysis using PANTHER*

9

10 Orthologs shared between *Calypte anna* and *Archilochus colubris*, but not with other birds
11 included in OrthoMCL analysis, were examined more closely using gene ontology (GO)
12 analysis. BlastP with default alignment settings was used to determine the top three putative
13 hits for each ortholog. Genbank accession numbers were converted in a universal gene symbol
14 using BioDB (<https://biodbnet-abcc.ncifcrf.gov/db/db2db.php>, (Mudunuri et al. 2009). Gene
15 symbols were then run through Panther (<http://pantherdb.org/>, (Mi et al. 2017), producing GO
16 terms for the input orthologs.
17
18
19

20 *Conservation analysis using PAML*

21

22 PAML4.9c (Yang 2007) was used in order to estimate pairwise conservation between
23 ruby-throated hummingbird and the five species (*Archilochus colubris*, *Calypte anna*, *Gallus*
24 *Gallus*, *Chaetura pelagica*, *Alligator mississippiensis* and *Homo sapiens*) used as comparators
25 in pathway analysis. Pairwise alignment was performed using CLUSTALX2.1 with PHYLIP
26 output. Codeml module was used, runmode= -2, and default settings. dN/dS was estimated
27 using Nei and Gojobori (Nei and Gojobori 1986). The mRNA sequences were then translated to
28 protein using ExPASy, and these proteins were aligned using CLUSTAL, alignment scores were
29 recorded.
30
31
32
33

34 *ORF prediction and protein translation using ANGEL*

35

36 The ANGEL pipeline (<https://github.com/PacificBiosciences/ANGEL>), a long read
37 implementation of ANGLE (Shimizu et al. 2006) was used in order to determine protein coding
38 sequences from cDNAs. ANGEL consists of three primary stages: dumb ORF prediction,
39 classifier training, and prediction. Dumb open reading frame (ORF) prediction, which produces
40 all six possible open reading frames per given transcript, was run for all transcripts with a
41 minimum length of 300 amino acids. Training involves the creation of a random subset of
42 non-redundant transcripts, which is then used to create a classifier pickle file implemented in the
43 prediction stage. Prediction then outputs the most likely ORF (minimum peptide length of 50
44 amino acids) based on length and coding potential of each given sequence. We performed this
45 analysis on both our HQ polished isoform (HQD) and all sequences datasets (ASD).
46
47
48
49

50 This resulted in 119,292 HQD and 1,061,147 ASD peptide sequences, with a size distribution
51 comparable to human, chicken, swift and alligator, with a mean amino acid length of roughly
52 500 AA, and a long tail.
53
54

55 **Data Accession**

56

57 Filtered fastq of clustered CCS reads deposited in SRA accession number SRP099041.
58 Predicted coding sequence and annotations, peptide and untranslated region data are available
59
60
61
62
63
64
65

1
2
3
4 at Zenodo 10.5281/zenodo.311651. Genbank submission is in progress. All other data available
5 upon request.
6

7 **Acknowledgments**

8
9 Pacific Biosciences for reagents and SMRTcells as well as technical support. M. Schatz, E.
10 Jarvis, J. Korlach, Y. Guo for discussion. HFSP grant #RGP0062/2016. Natural Sciences and
11 Engineering Research Council of Canada Discovery Grant (#386466) to KCW.
12
13

14 **Disclosure Declaration**

15
16 W.T. and R.W. have received travel funds to speak at symposia organized by Pacific
17 Biosciences. Bulk of reagents for IsoSeq were provided by Pacific Biosciences.
18
19

20 **Figure legends**

21
22 **Figure 1.** Analysis pipeline. **A** Raw sequence reads from a Pacbio RSII sequencer (bax.h5,
23 bas.h5) were sorted into full and non-full length reads using a classification algorithm that
24 identified full length reads with forward and reverse primers, as well as a poly-A tail. Iterative
25 clustering for isoforms (ICE) was performed on full length reads, and non-full length reads were
26 recruited to perform ARROW polished on the consensus isoforms. Polishing sorted reads into
27 high and low-quality bins, and either high quality data (HQD), all sequence data (ASD) or both
28 sets of data, were carried on to further applications (**B**).
29
30
31
32

33 **Figure 2.** Transcriptome dataset quality control. Average read lengths and isoform counts for 4
34 sequenced size fractions given in **A**, and read length for all sequence data (ASD, HQ + LQ)
35 plotted in **B**, with black line representing Mb data greater than read length. For example, at
36 2000bp, 5000Mb of sequence data was larger than 2000bp. **C.** BUSCO transcriptome
37 assessment results for *Archilochus colubris* (ruby-throated hummingbird, all sequence data
38 ASD, high quality sequence data HQD), COGENT-collapsed data (CCD), *Calypte anna* (Anna's
39 hummingbird), *Gallus gallus* (chicken) Thomas (single-tissue transcriptome).
40
41
42

43 **Figure 3.** Reducing transcript redundancy and predicting gene families with COGENT.
44 COGENT gene families predicted and classified by relationship to *Calypte anna* genome
45 assembly **A**. Number of full-length reads which support each isoform reduced **B**. **C** IGV view of
46 the MATR3 gene, which was reduced from 11 redundant reads to 3 unique isoforms using this
47 pipeline.
48
49

50
51 **Figure 4.** Orthology analysis. The transcriptomes of five birds (*Calypte anna*, *Tinamus guttatus*,
52 *Gallus gallus*, *Chaetura pelagica*, and *Melopsitticus undulatus*), one mammal (*Homo sapiens*)
53 and one reptile (*Alligator mississippiensis*) were compared against *Archilochus colubris* using
54 OrthoMCL to detect and compare similar sequences. A Venn diagram illustrating sequences
55 with reciprocal blast hits between the given species and *A. colubris* is shown in **A**. Ortholog
56 pairs unique to hummingbirds *Calypte anna* and *Archilochus colubris* were selected for gene
57 orthology (GO) annotation analysis, which revealed enzymes of many biological functions (**B**).
58
59
60
61
62
63
64
65

1
2
3
4 **Figure 5.** Pathway analysis of key enzymes in hepatic lipogenesis. **A** An overview of the
5 relationship between the investigated genes and their roles in triacylglycerol, phospholipid and
6 fatty acid synthesis, **B** a heat map illustrating percent identity of these proteins relative to
7 *Archilochus colubris* predicted sequences, and nucleotide coding sequence conservation scores
8 predicted using PAML (dN/dS), with genes dN/dS>1 starred.
9

10
11 **References:**

- 12
13
14 Abdel-Ghany SE, Hamilton M, Jacobi JL, Ngam P, Devitt N, Schilkey F, Ben-Hur A, Reddy ASN.
15 2016. A survey of the sorghum transcriptome using single-molecule long reads. *Nat*
16 *Commun* **7**: 11706.
17
18 Aleshin AE, Zeng C, Bourenkov GP, Bartunik HD, Fromm HJ, Honzatko RB. 1998. The
19 mechanism of regulation of hexokinase: new insights from the crystal structure of
20 recombinant human brain hexokinase complexed with glucose and glucose-6-phosphate.
21 *Structure* **6**: 39–50.
22
23 Baker HG. 1975. Sugar Concentrations in Nectars from Hummingbird Flowers. *Biotropica* **7**:
24 37–41.
25
26 Beuchat CA, Chong CR. 1998. Hyperglycemia in hummingbirds and its consequences for
27 hemoglobin glycation. *Comp Biochem Physiol A Mol Integr Physiol* **120**: 409–416.
28
29 Braun EJ, Sweazea KL. 2008. Glucose regulation in birds. *Comp Biochem Physiol B Biochem*
30 *Mol Biol* **151**: 1–9.
31
32 Carpenter FL, Hixon MA, Beuchat CA, Russell RW, Paton DC. 1993. Biphasic Mass Gain in
33 Migrant Hummingbirds: Body Composition Changes, Torpor, and Ecological Significance.
34 *Ecology* **74**: 1173–1182.
35
36 Chai P, Dudley R. 1996. Limits to flight energetics of hummingbirds hovering in hypodense and
37 hypoxic gas mixtures. *J Exp Biol* **199**: 2285–2295.
38
39 Chen CCW, Welch KC. 2014. Hummingbirds can fuel expensive hovering flight completely with
40 either exogenous glucose or fructose. *Funct Ecol* **28**: 589–600.
41
42 Fan L, Gardner P, Chan SJ, Steiner DF. 1993. Cloning and analysis of the gene encoding
43 hummingbird proinsulin. *Gen Comp Endocrinol* **91**: 25–30.
44
45 FernándezM.J., BozinovicF., SuarezR.K. 2011. Enzymatic flux capacities in hummingbird flight
46 muscles: a “one size fits all” hypothesis. *Can J Zool* **89**: 985–991.
47
48 Gordon SP, Tseng E, Salamov A, Zhang J, Meng X, Zhao Z, Kang D, Underwood J, Grigoriev
49 IV, Figueroa M, et al. 2015. Widespread Polycistronic Transcripts in Fungi Revealed by
50 Single-Molecule mRNA Sequencing. *PLoS One* **10**: e0132628.
51
52 Gregory TR, Andrews CB, McGuire JA, Witt CC. 2009. The smallest avian genomes are found
53 in hummingbirds. *Proc Biol Sci* **276**: 3753–3757.
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 Hermier D. 1997. Lipoprotein metabolism and fattening in poultry. *J Nutr* **127**: 805S–808S.
5
6 Hou L, Verdirame M, Welch KC. 2015. Automated tracking of wild hummingbird mass and
7 energetics over multiple time scales using radio frequency identification (RFID) technology.
8 *J Avian Biol* **46**: 1–8.
9
10 Hou L, Welch KC Jr. 2016. Premigratory ruby-throated hummingbirds, *Archilochus colubris*,
11 exhibit multiple strategies for fuelling migration. *Anim Behav* **121**: 87–99.
12
13 Hughes AL, Hughes MK. 1995. Small genomes for better flyers. *Nature* **377**: 391.
14
15 Jacob E, Unger R. 2007. A tale of two tails: why are terminal residues of proteins exposed?
16 *Bioinformatics* **23**: e225–30.
17
18 Jarvis ED, Mirarab S, Aberer AJ, Li B, Houde P, Li C, Ho SYW, Faircloth BC, Nabholz B,
19 Howard JT, et al. 2015. Phylogenomic analyses data of the avian phylogenomics project.
20 *Gigascience* **4**: 4.
21
22 Joost HG, Thorens B. 2001. The extended GLUT-family of sugar/polyol transport facilitators:
23 nomenclature, sequence characteristics, and potential function of its novel members
24 (review). *Mol Membr Biol* **18**: 247–256.
25
26 Kamata K, Mitsuya M, Nishimura T, Eiki J-I, Nagata Y. 2004. Structural basis for allosteric
27 regulation of the monomeric allosteric enzyme human glucokinase. *Structure* **12**: 429–438.
28
29 Korlach J, Gedman G, Kingan S, Chin J, Howard J, Cantin L, Jarvis ED. 2017. De Novo PacBio
30 long-read and phased avian genome assemblies correct and add to genes important in
31 neuroscience research. *bioRxiv* 103911. <http://biorxiv.org/content/early/2017/02/02/103911>
32 (Accessed February 9, 2017).
33
34 Kosuri S, Church GM. 2014. Large-scale de novo DNA synthesis: technologies and
35 applications. *Nat Methods* **11**: 499–507.
36
37 Kryazhimskiy S, Plotkin JB. 2008. The population genetics of dN/dS. *PLoS Genet* **4**: e1000304.
38
39 Lasso G, Yu LPC, Gil D, Lázaro M, Tong L, Valle M. 2014. Functional conformations for
40 pyruvate carboxylase during catalysis explored by cryoelectron microscopy. *Structure* **22**:
41 911–922.
42
43 Li L, Stoeckert CJ Jr, Roos DS. 2003. OrthoMCL: identification of ortholog groups for eukaryotic
44 genomes. *Genome Res* **13**: 2178–2189.
45
46 McGuire JA, Witt CC, Altshuler DL, Remsen JV Jr. 2007. Phylogenetic systematics and
47 biogeography of hummingbirds: Bayesian and maximum likelihood analyses of partitioned
48 data and selection of an appropriate partitioning strategy. *Syst Biol* **56**: 837–856.
49
50 McGuire JA, Witt CC, Remsen JV Jr, Corl A, Rabosky DL, Altshuler DL, Dudley R. 2014.
51 Molecular phylogenetics and the diversification of hummingbirds. *Curr Biol* **24**: 910–916.
52
53 Mi H, Huang X, Muruganujan A, Tang H, Mills C, Kang D, Thomas PD. 2017. PANTHER version
54 11: expanded annotation data from Gene Ontology and Reactome pathways, and data
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 analysis tool enhancements. *Nucleic Acids Res* **45**: D183–D189.
5
6 Mi H, Muruganujan A, Casagrande JT, Thomas PD. 2013. Large-scale gene function analysis
7 with the PANTHER classification system. *Nat Protoc* **8**: 1551–1566.
8
9 Minoche AE, Dohm JC, Schneider J, Holtgräwe D, Viehöver P, Montfort M, Sörensen TR,
10 Weisshaar B, Himmelbauer H. 2015. Exploiting single-molecule transcript sequencing for
11 eukaryotic gene prediction. *Genome Biol* **16**: 184.
12
13 Mudunuri U, Che A, Yi M, Stephens RM. 2009. bioDBnet: the biological database network.
14 *Bioinformatics* **25**: 555–556.
15
16 Mulichak AM, Wilson JE, Padmanabhan K, Garavito RM. 1998. The structure of mammalian
17 hexokinase-1. *Nat Struct Biol* **5**: 555–560.
18
19 Nei M, Gojobori T. 1986. Simple methods for estimating the numbers of synonymous and
20 nonsynonymous nucleotide substitutions. *Mol Biol Evol* **3**: 418–426.
21
22 Ohtsubo K, Takamatsu S, Gao C, Korekane H, Kurosawa TM, Taniguchi N. 2013.
23 N-Glycosylation modulates the membrane sub-domain distribution and activity of glucose
24 transporter 2 in pancreatic beta cells. *Biochem Biophys Res Commun* **434**: 346–351.
25
26 O’Leary NA, Wright MW, Brister JR, Ciufo S, Haddad D, McVeigh R, Rajput B, Robbertse B,
27 Smith-White B, Ako-Adjei D, et al. 2016. Reference sequence (RefSeq) database at NCBI:
28 current status, taxonomic expansion, and functional annotation. *Nucleic Acids Res* **44**:
29 D733–45.
30
31 Polakof S, Mommsen TP, Soengas JL. 2011. Glucosensing and glucose homeostasis: from fish
32 to mammals. *Comp Biochem Physiol B Biochem Mol Biol* **160**: 123–149.
33
34 Powers DR, Brown AR, Van Hook JA. 2003. Influence of normal daytime fat deposition on
35 laboratory measurements of torpor use in territorial versus nonterritorial hummingbirds.
36 *Physiol Biochem Zool* **76**: 389–397.
37
38 Shimizu K, Adachi J, Muraoka Y. 2006. ANGLE: A SEQUENCING ERRORS RESISTANT
39 PROGRAM FOR PREDICTING PROTEIN CODING REGIONS IN UNFINISHED cDNA. *J*
40 *Bioinform Comput Biol* **04**: 649–664.
41
42 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO:
43 assessing genome assembly and annotation completeness with single-copy orthologs.
44 *Bioinformatics* **31**: 3210–3212.
45
46 Suarez RK. 1992. Hummingbird flight: sustaining the highest mass-specific metabolic rates
47 among vertebrates. *Experientia* **48**: 565–570.
48
49 Suarez RK, Brownsey RW, Vogl W, Brown GS, Hochachka PW. 1988. Biosynthetic capacity of
50 hummingbird liver. *Am J Physiol* **255**: R699–702.
51
52 Suarez RK, Lighton JR, Moyes CD, Brown GS, Gass CL, Hochachka PW. 1990. Fuel selection
53 in rufous hummingbirds: ecological implications of metabolic biochemistry. *Proc Natl Acad*
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 *Sci U S A* **87**: 9207–9210.
5

6 Suarez RK, Welch KC Jr, Hanna SK, Herrera M LG. 2009. Flight muscle enzymes and
7 metabolic flux rates during hovering flight of the nectar bat, *Glossophaga soricina*: further
8 evidence of convergence with hummingbirds. *Comp Biochem Physiol A Mol Integr Physiol*
9 **153**: 136–140.
10

11 ter Kuile BH, Westerhoff HV. 2001. Transcriptome meets metabolome: hierarchical and
12 metabolic regulation of the glycolytic pathway. *FEBS Lett* **500**: 169–171.
13

14
15 Thomas S, Underwood JG, Tseng E, Holloway AK, Bench To Basinet CvDC Informatics
16 Subcommittee. 2014. Long-read sequencing of chicken transcripts and identification of new
17 transcript isoforms. *PLoS One* **9**: e94650.
18

19
20 Thorens B, Mueckler M. 2010. Glucose transporters in the 21st Century. *Am J Physiol*
21 *Endocrinol Metab* **298**: E141–5.
22

23 Vianna CR, Hagen T, Zhang CY, Bachman E, Boss O, Gereben B, Moriscot AS, Lowell BB,
24 Bicudo JE, Bianco AC. 2001. Cloning and functional characterization of an uncoupling
25 protein homolog in hummingbirds. *Physiol Genomics* **5**: 137–145.
26

27 Wang B, Tseng E, Regulski M, Clark TA, Hon T, Jiao Y, Lu Z, Olson A, Stein JC, Ware D. 2016.
28 Unveiling the complexity of the maize transcriptome by single-molecule long-read
29 sequencing. *Nat Commun* **7**: 11708.
30

31
32 Weidensaul S, Robinson TR, Sargent RR, Sargent MB, Poole A. 2013. Ruby-throated
33 hummingbird (*Archilochus colubris*). *Birds of North America Online (A Poole, Editor) Cornell*
34 *Lab of Ornithology, Ithaca, NY, USA* <http://bna.birds.cornell.edu/bna/species/204>.
35

36
37 Welch KC Jr, Allalou A, Sehgal P, Cheng J, Ashok A. 2013. Glucose transporter expression in
38 an avian nectarivore: the ruby-throated hummingbird (*Archilochus colubris*). *PLoS One* **8**:
39 e77003.
40

41 Welch KC Jr, Altshuler DL, Suarez RK. 2007. Oxygen consumption rates in hovering
42 hummingbirds reflect substrate-dependent differences in P/O ratios: carbohydrate as a
43 “premium fuel.” *J Exp Biol* **210**: 2146–2153.
44

45
46 Welch KC Jr, Chen CCW. 2014. Sugar flux through the flight muscles of hovering vertebrate
47 nectarivores: a review. *J Comp Physiol B* **184**: 945–959.
48

49 Wu TD, Watanabe CK. 2005. GMAP: a genomic mapping and alignment program for mRNA
50 and EST sequences. *Bioinformatics* **21**: 1859–1875.
51

52 Xiang S, Tong L. 2008. Crystal structures of human and *Staphylococcus aureus* pyruvate
53 carboxylase and molecular insights into the carboxyltransfer reaction. *Nat Struct Mol Biol*
54 **15**: 295–302.
55

56
57 Yang Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. *Mol Biol Evol* **24**:
58 1586–1591.
59

60 Zhang G, Li B, Li C, Gilbert MTP, Mello CV, Jarvis ED, Consortium TAG, Wang J. 2014.
61
62
63
64
65

1
2
3
4 Genomic data of the Anna's Hummingbird (*Calypte anna*). <https://doi.org/10.5524/101004>.

5
6 Zhang G, Rahbek C, Graves GR, Lei F, Jarvis ED, Gilbert MTP. 2015. Genomics: Bird
7 sequencing project takes off. *Nature* **522**: 34.

8
9
10 Zhang H, Yang Z, Shen Y, Tong L. 2003. Crystal structure of the carboxyltransferase domain of
11 acetyl-coenzyme A carboxylase. *Science* **299**: 2064–2067.

12
13 Zhang JZ, Abbud W, Prohaska R, Ismail-Beigi F. 2001. Overexpression of stomatin depresses
14 GLUT-1 glucose transporter activity. *Am J Physiol Cell Physiol* **280**: C1277–83.

15
16 Joost HG, Thorens B. 2001. The extended GLUT-family of sugar/polyol transport facilitators:
17 nomenclature, sequence characteristics, and potential function of its novel members
18 (review). *Mol Membr Biol* **18**: 247-256.

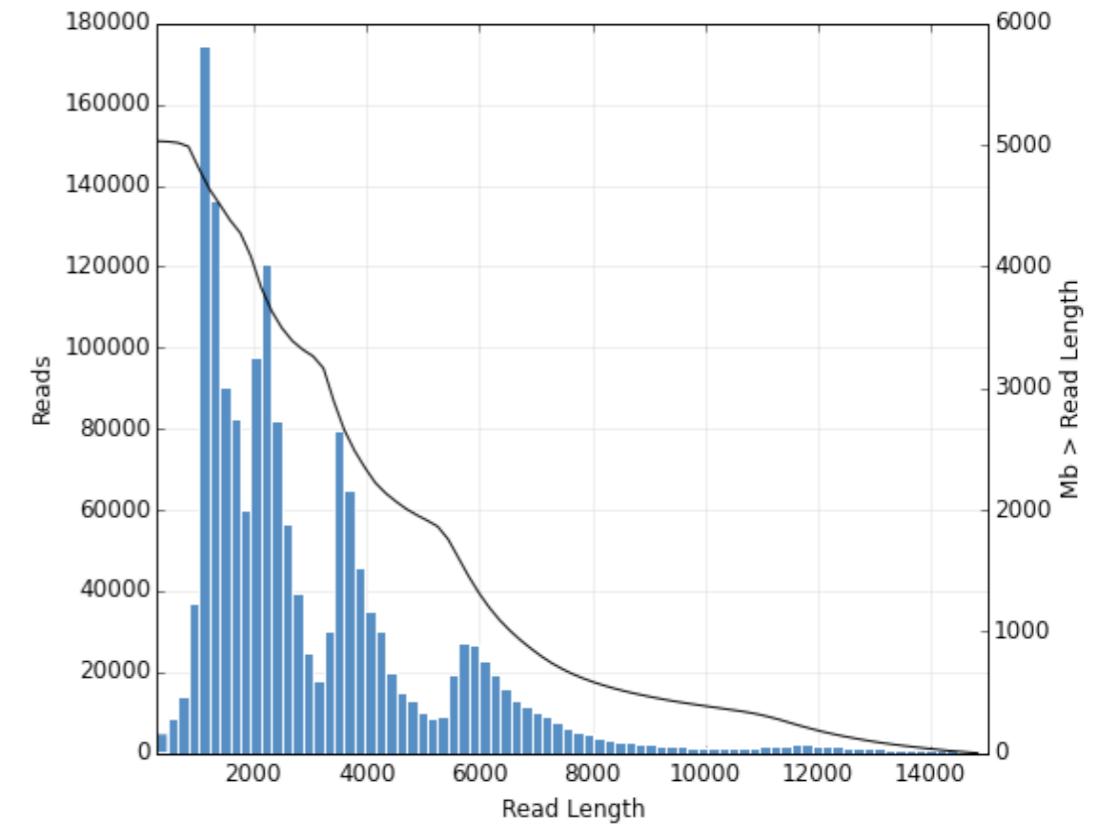
19
20
21 Kryazhimskiy S, Plotkin JB. 2008. The population genetics of of dN/dS. *PLoS Genet* **4**:
22 e1000304.

23
24 McGuire JA, Witt CC, Remsen JV Jr., Corl A, Rabosky DL, Altshuler DL, Dudley R. 2014.
25 Molecular Phylogenetics and Diversification of Hummingbirds. *Curr Biol* **24**: 910-916.
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]

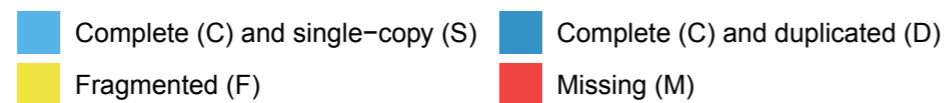
Size Fraction	1-2kb	2-3kb	3-6kb	5-10kb	Total
# of cells	10	10	10	10	40
Reads of Insert (ROI)	688,069	591,050	735,670	625,194	2,639,983
Avg length ROI (bp)	1533	2464	3650	5444	
ROI Yield (Mbp)	1055	1457	2685	3404	8601
Filtered Reads (FLNC)	430,381	306,841	272,781	193,906	1,203,909
# Consensus Isoforms	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]

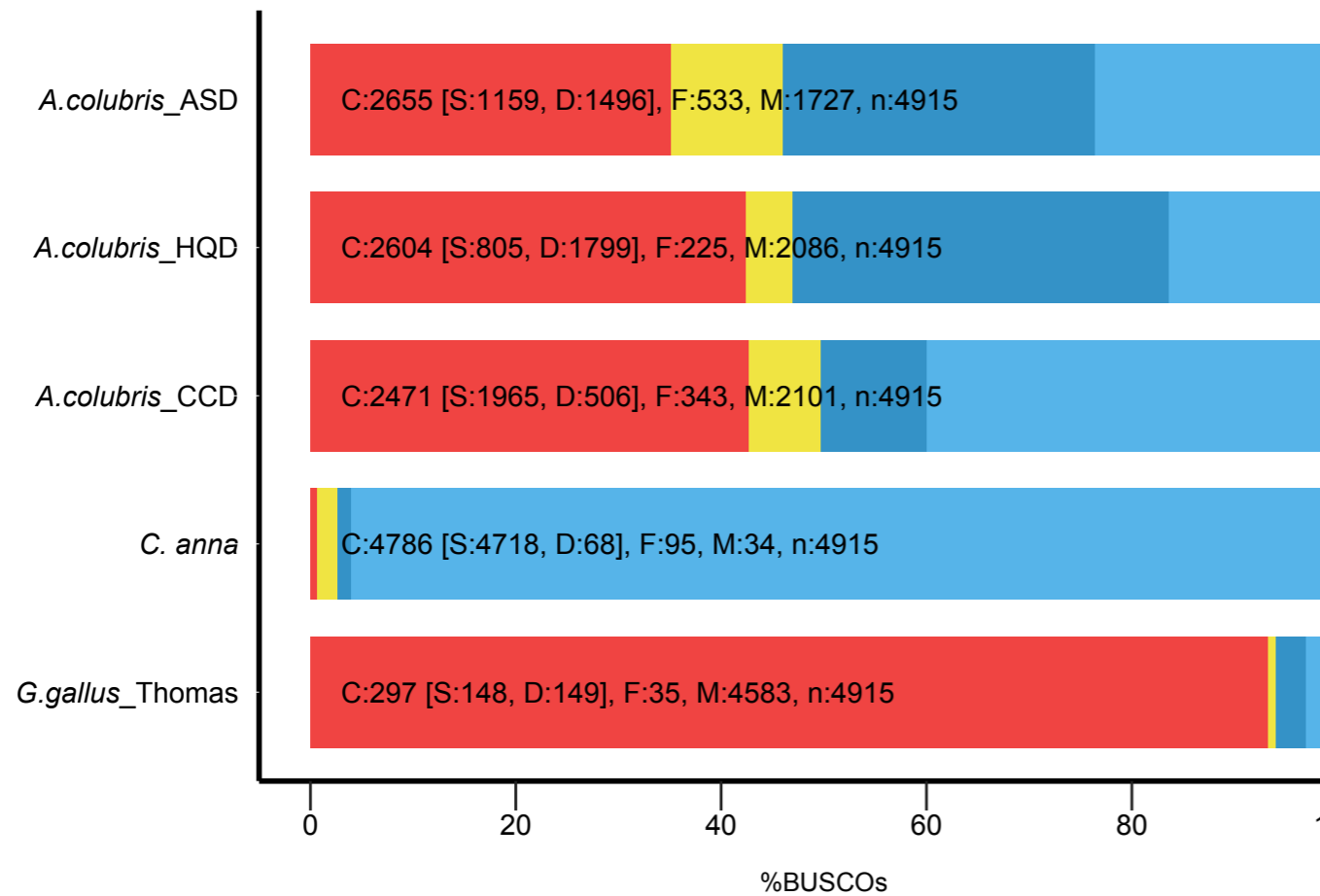


C]

BUSCO ASSESSMENT RESULTS



Aves



Metazoan

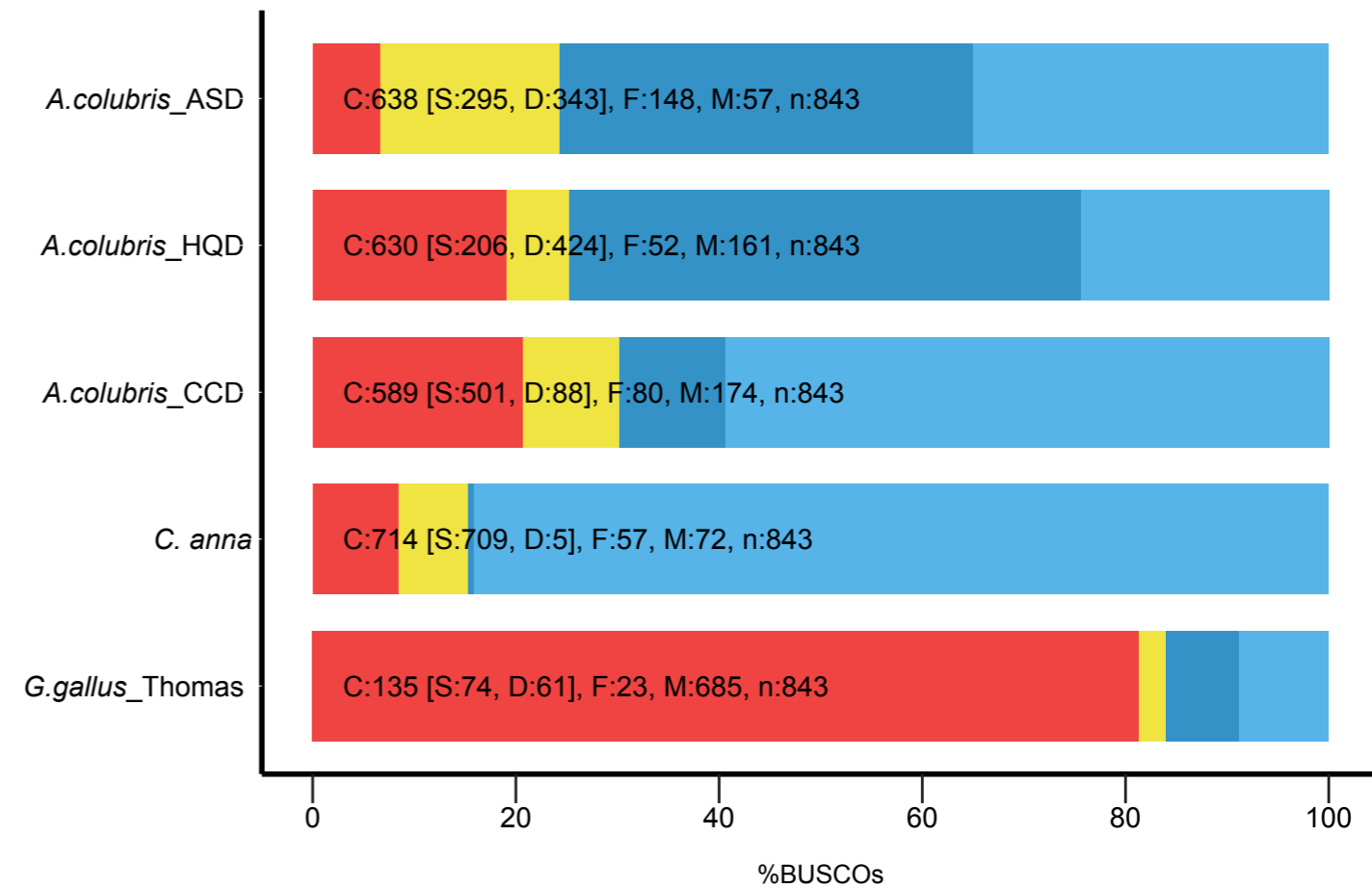
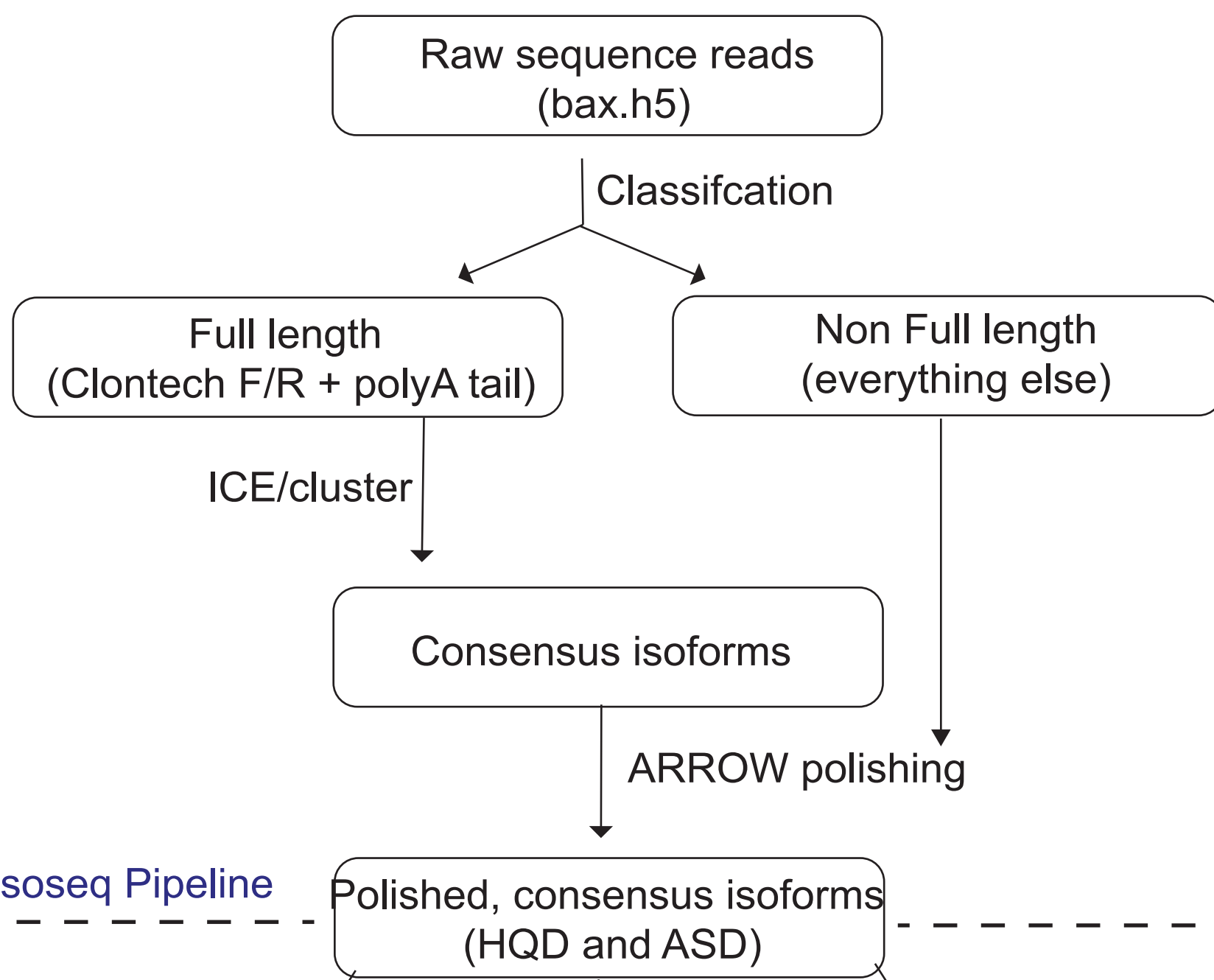


Figure 1. Transcriptome dataset quality control. Average read lengths and isoform counts for 4 sequenced size fractions given in **A**, and read length for all sequence data (ASD, HQ + LQ) plotted in **B**, with black line representing Mb data greater than read length. For example, at 2000bp, 5000Mb of sequence data was larger than 2000bp. **C**. BUSCO transcriptome assessment results for *Archilochus colubris* (ruby-throated hummingbird, all sequence data ASD, high quality sequence data HQD), Cogent-collapsed data (CCD), *Calypte anna* (Anna's hummingbird), *Gallus gallus* (chicken) Thomas (single-tissue transcriptome).

A]



Pacbio SMRT Analysis Isoseq Pipeline

Applications

B]

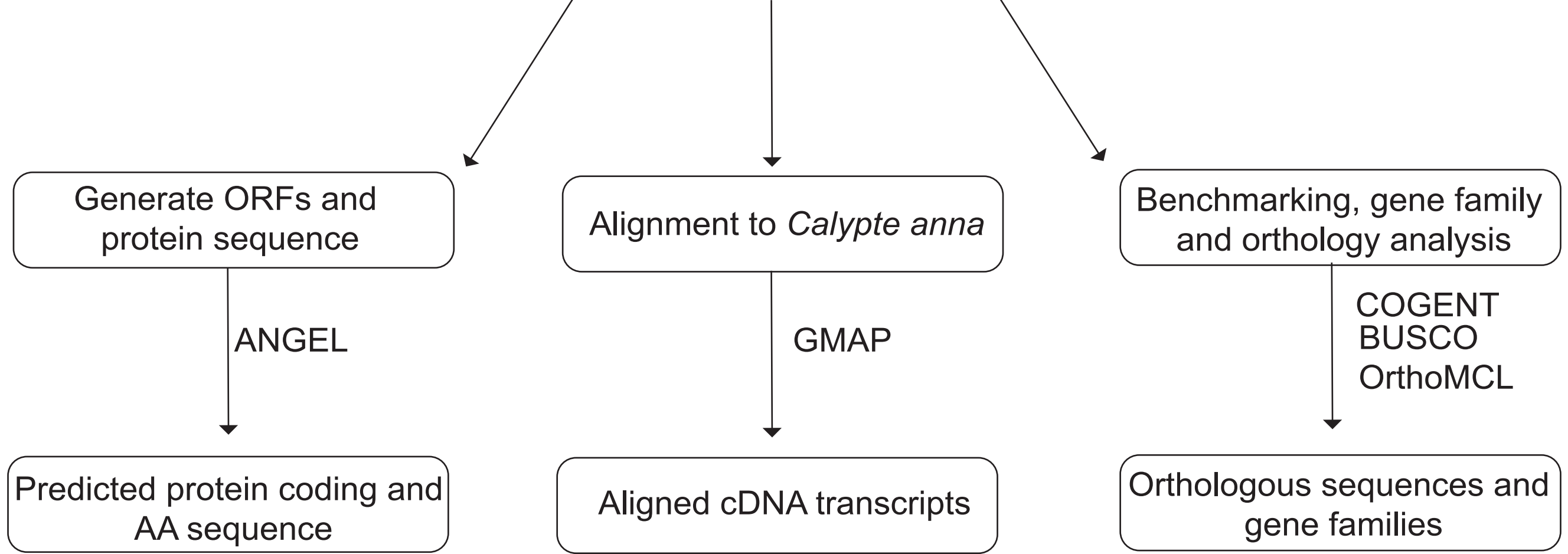


Figure 2. Analysis pipeline. **A** Raw sequence reads from a Pacbio RSII sequencer (bax.h5, bas.h5) were sorted into full and non-full length reads using a classification algorithm that identified full length reads with forward and reverse primers, as well as a poly-A tail. Iterative clustering for isoforms (ICE) was performed on full length reads, and non-full length reads were recruited to perform ARROW polished on the consensus isoforms. Polished sorted reads into high and low-quality bins, and either high quality data (HQD), all sequence data (ASD) or both sets of data, were carried on to further applications (**B**).

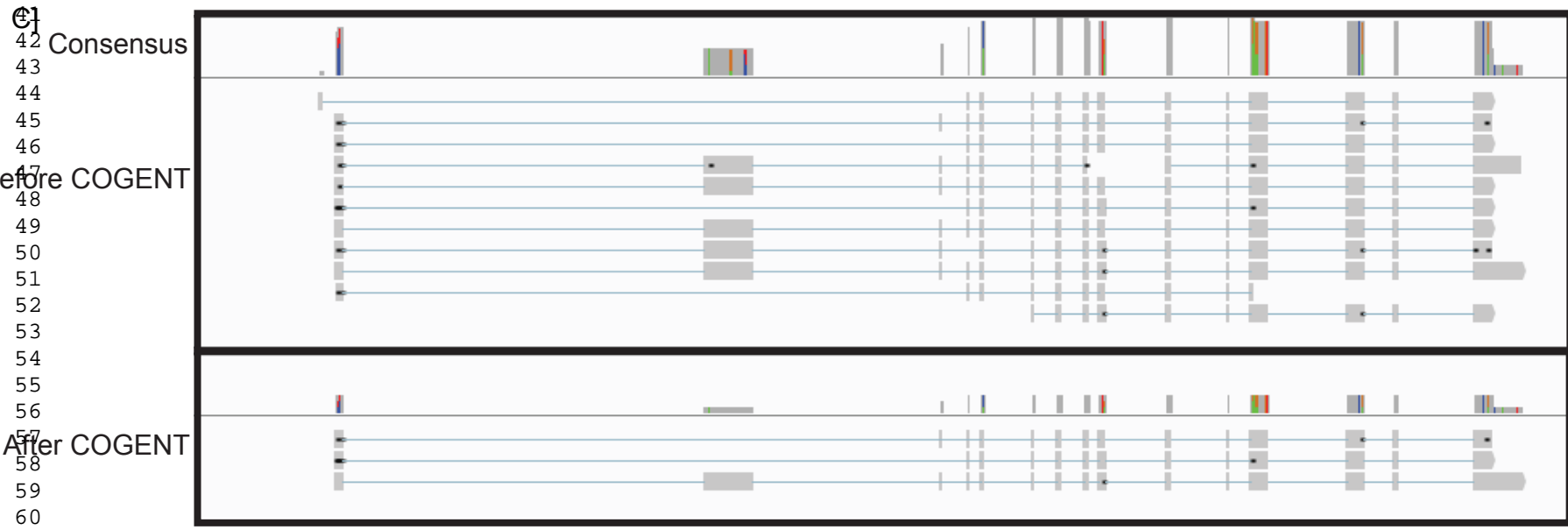
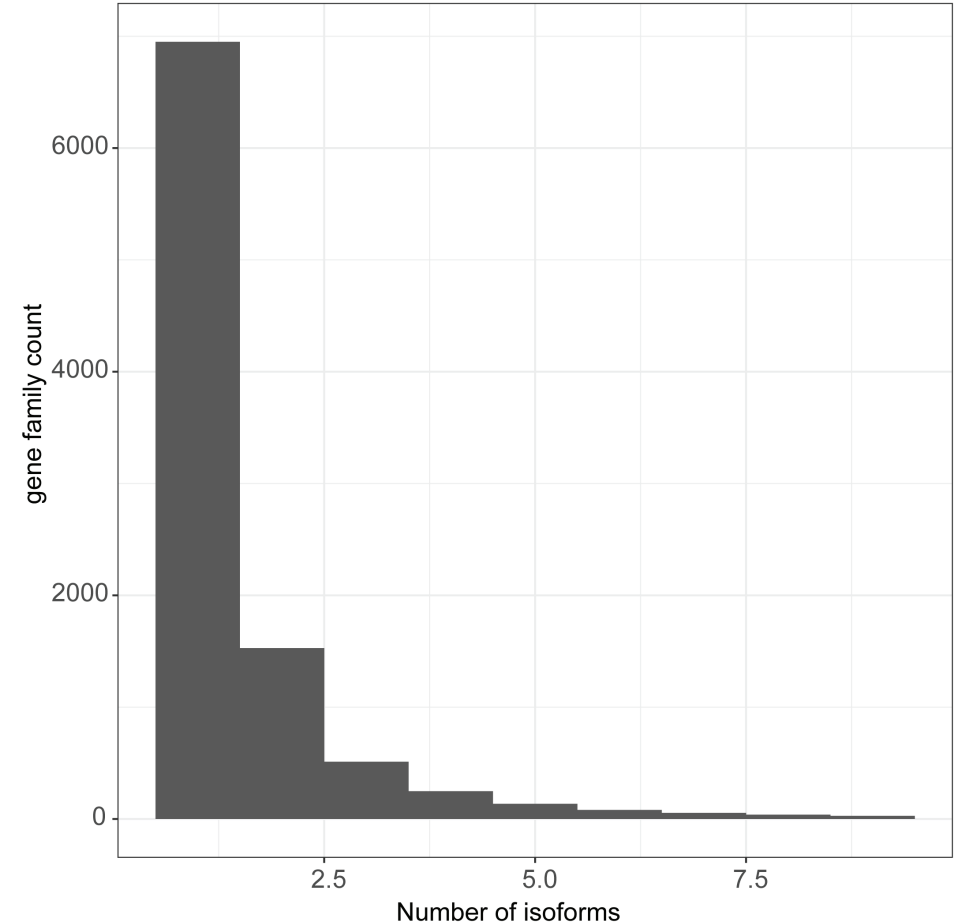
1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

9
10
11 **A)**

Clustering results	Counts		
Total HQ Arrow Isoforms	94724		
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

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

B)

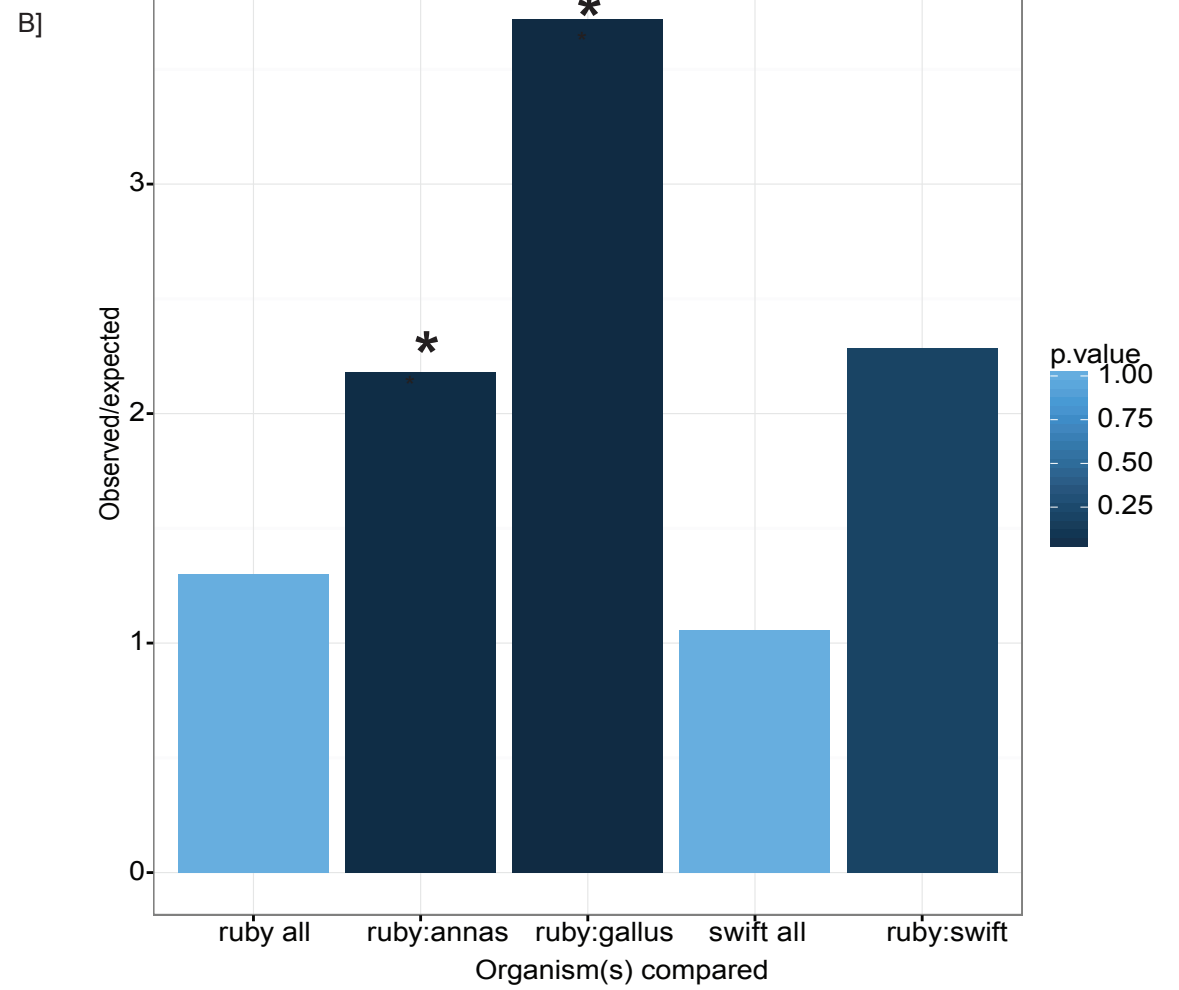
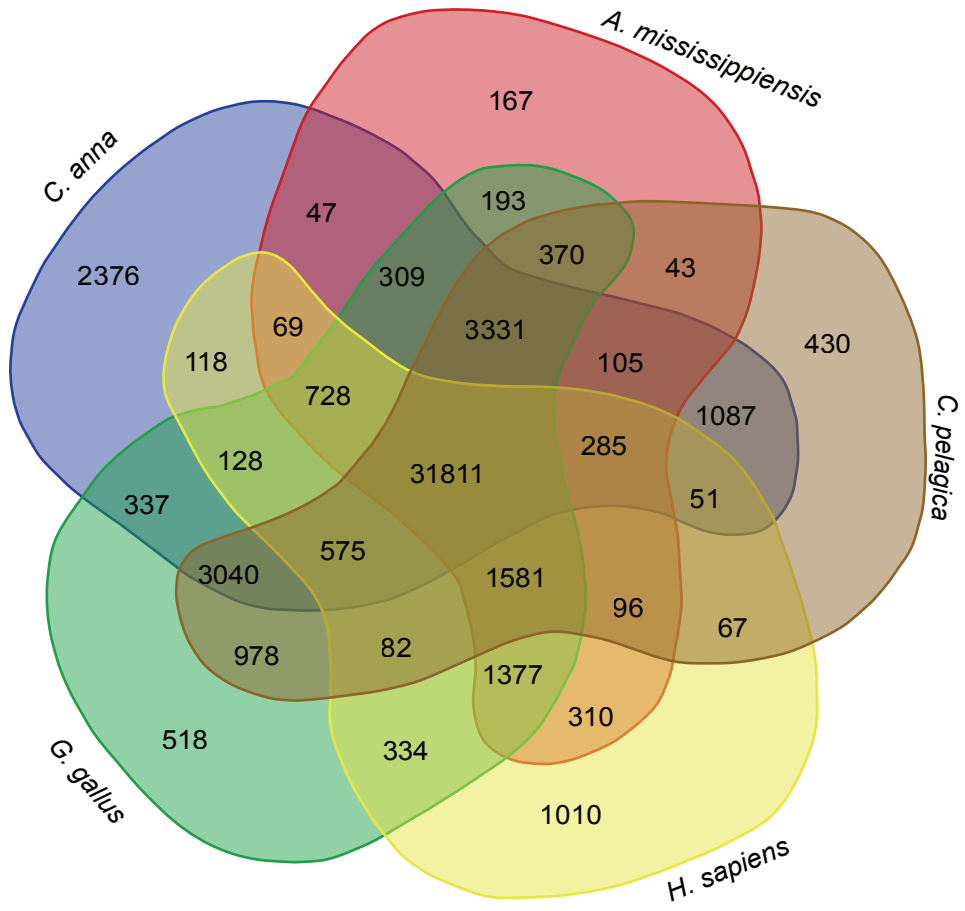


Cogent family 14912
MATR3 gene
FALCON assembly
scaffold 000168F
860,227-921,621

62
63
64

Figure 3. Reducing transcript redundancy and predicting gene families with COGENT. Cogent gene families predicted and classified by relationship to *Calypte anna* genome assembly **A**. Number of full-length reads which support each isoform reduced **B**. **C** IGV view of the MATR3 gene, which was reduced from 11 redundant reads to 3 unique isoforms using this pipeline.

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 Figure 4. Orthology analysis. The transcriptomes of three birds (*Calypte anna*, *Gallus gallus*, *Chaetura pelagica*), one mammal (Homo
61 sapiens) and one reptile (*Alligator mississippiensis*) were compared against *Archilochus colubris* using OrthoMCL to detect and compare
62 similar sequences. **A** Venn diagram illustrating sequences with reciprocal blast hits between the given species and *A. colubris* is shown in
63 **A**. Ortholog pairs unique to hummingbirds *Calypte anna* and *Archilochus colubris* were selected for gene orthology (GO) annotation
64 analysis, which revealed enzymes of many biological functions (**B**).
65

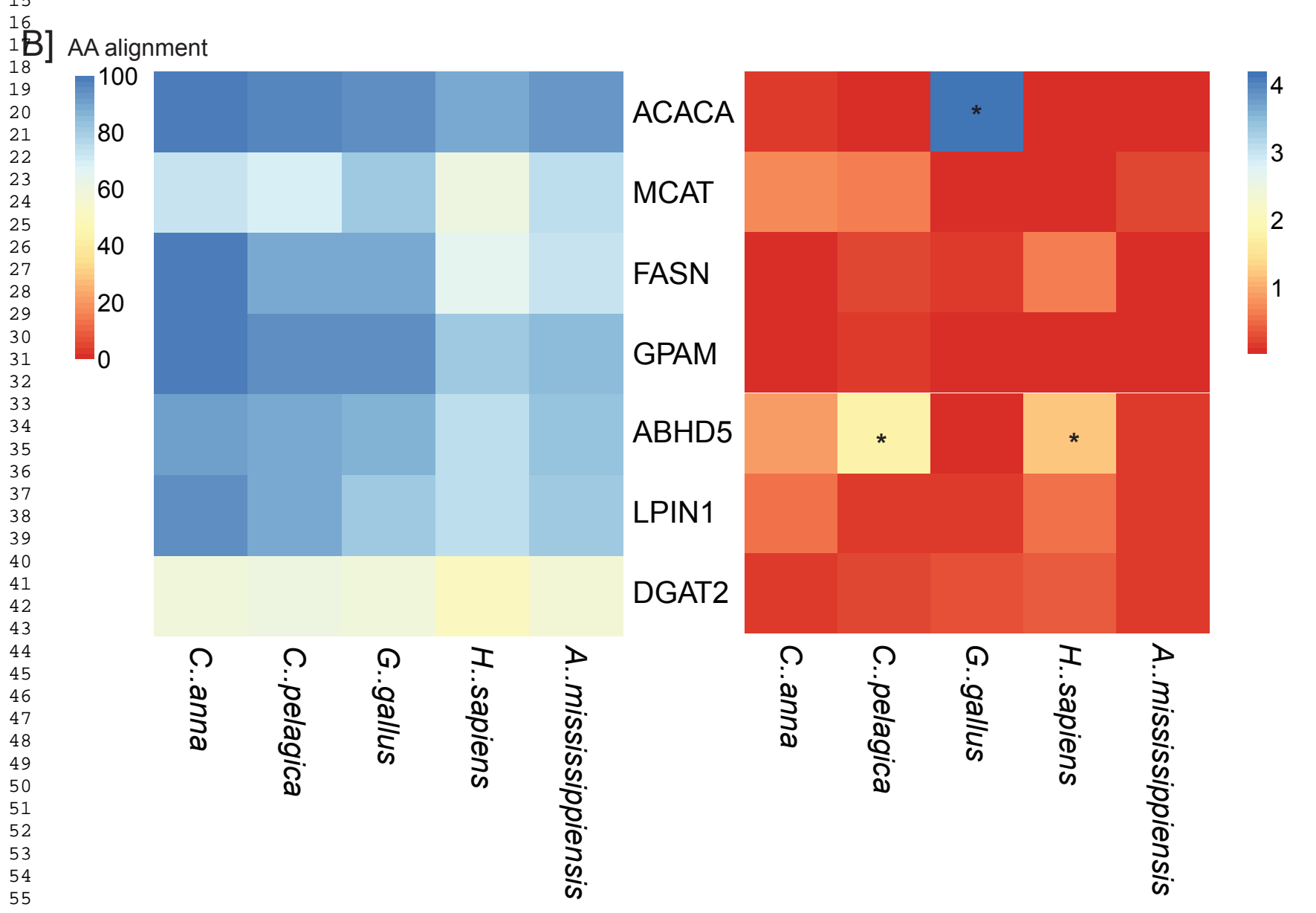
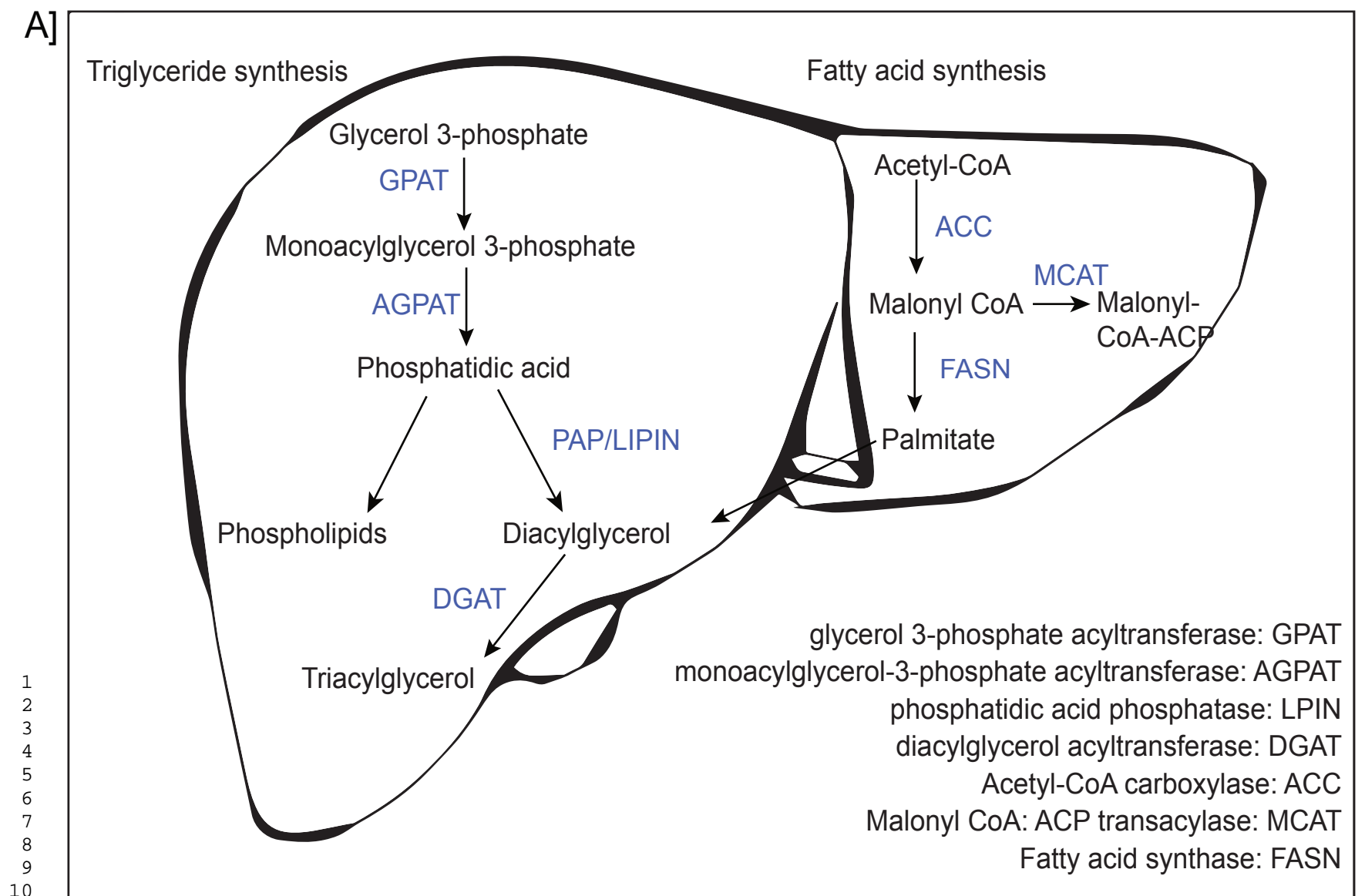


Figure 5. Pathway analysis of key enzymes in hepatic lipogenesis. **A]** An overview of the relationship between the investigated genes and their roles in triacylglycerol, phospholipid and fatty acid synthesis, **B** a heat map illustrating percent identity of these proteins relative to *Archilochus colubris* predicted sequences, and nucleotide coding sequence conservation scores predicted using PAML (dN/dS), with genes dN/dS>1 starred.

A]

Size Fraction	1-2kb	2-3kb	3-6kb	5-10kb	Total
# of cells	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 Reads (FLNC)	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	

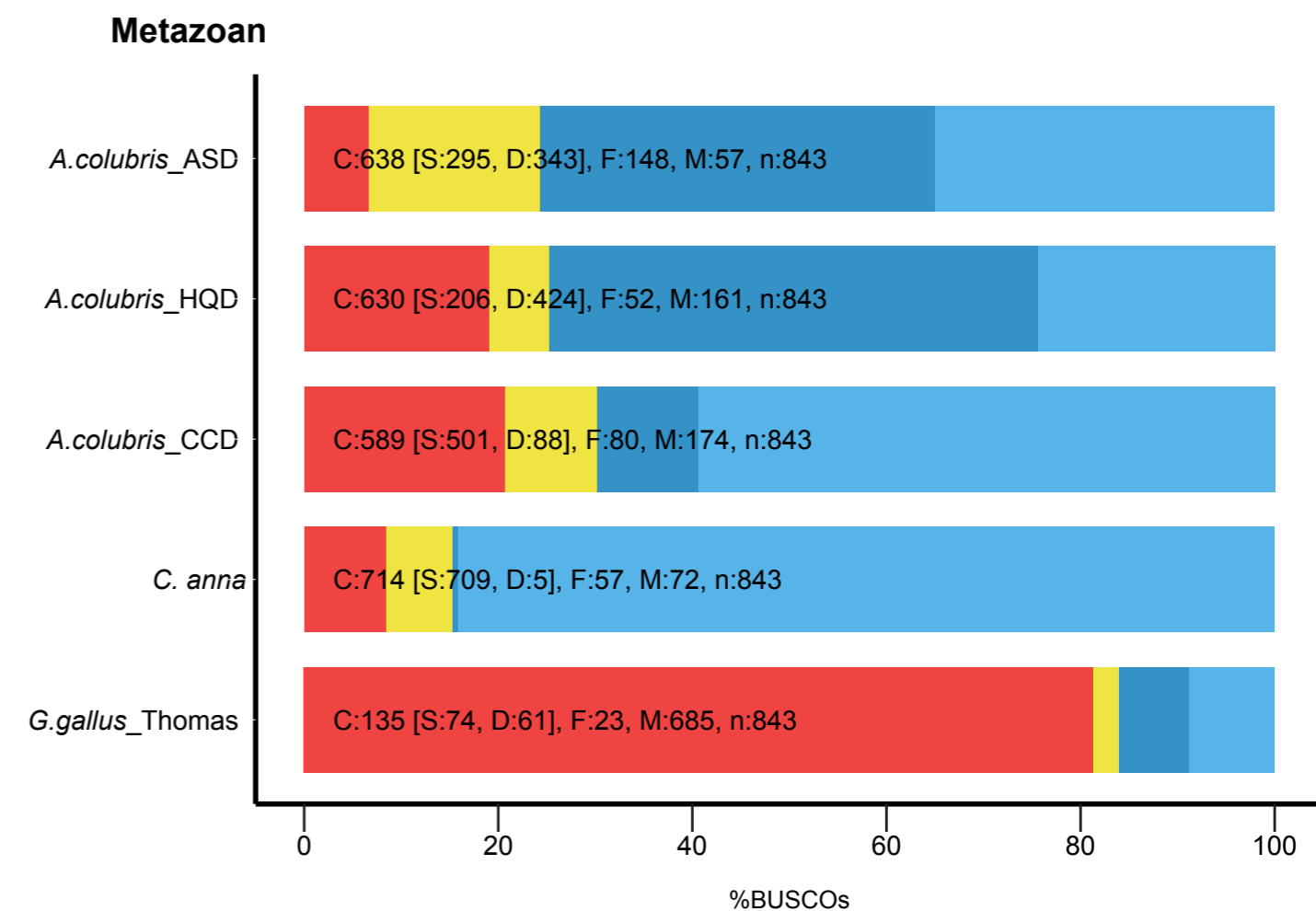
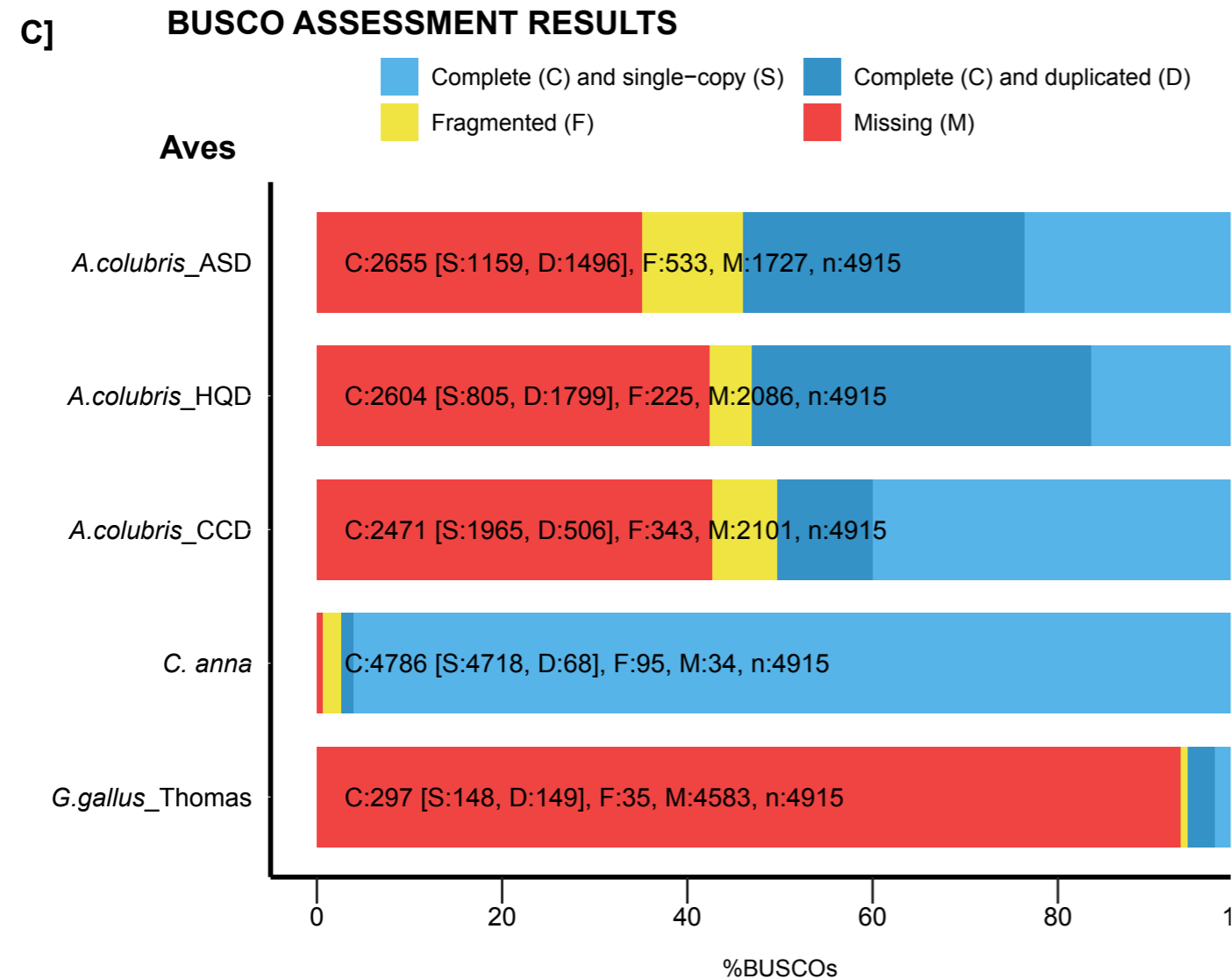
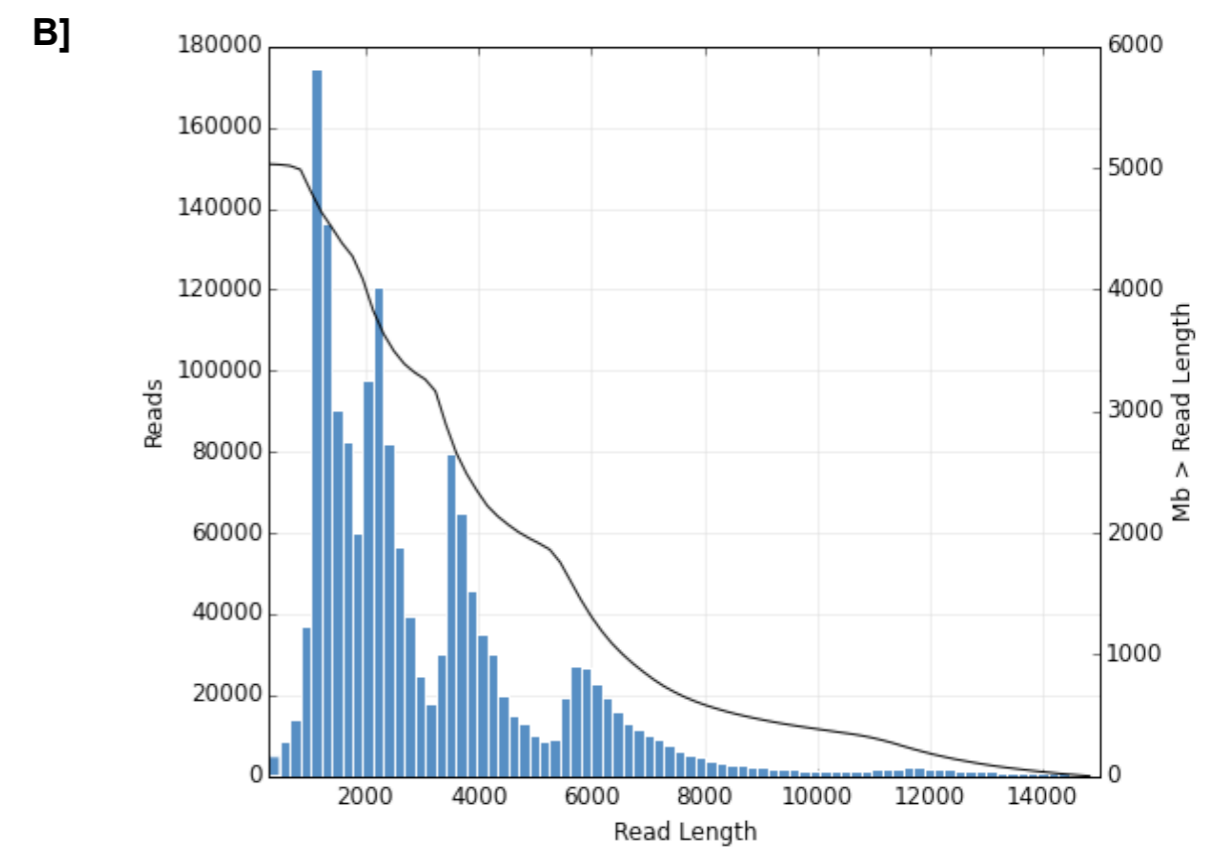
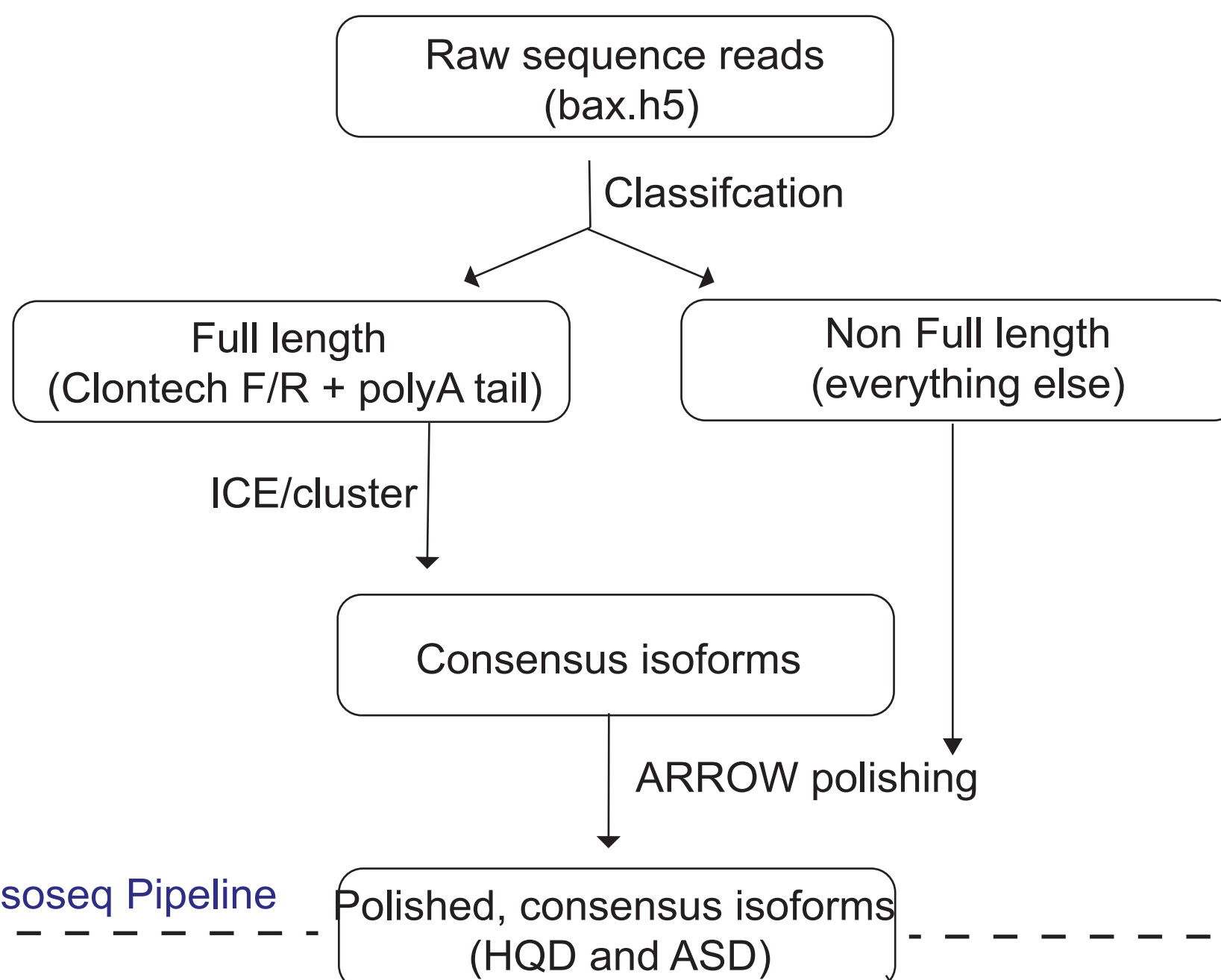


Figure 1. Transcriptome dataset quality control. Average read lengths and isoform counts for 4 sequenced size fractions given in **A**, and read length for all sequence data (ASD, HQ + LQ) plotted in **B**, with black line representing Mb data greater than read length. For example, at 2000bp, 5000Mb of sequence data was larger than 2000bp. **C**. BUSCO transcriptome assessment results for *Archilochus colubris* (ruby-throated hummingbird, all sequence data ASD, high quality sequence data HQD), Cogent-collapsed data (CCD), *Calypte anna* (Anna's hummingbird), *Gallus gallus* (chicken) Thomas (single-tissue transcriptome).

A]



Pacbio SMRT Analysis Isoseq Pipeline

Applications

B]

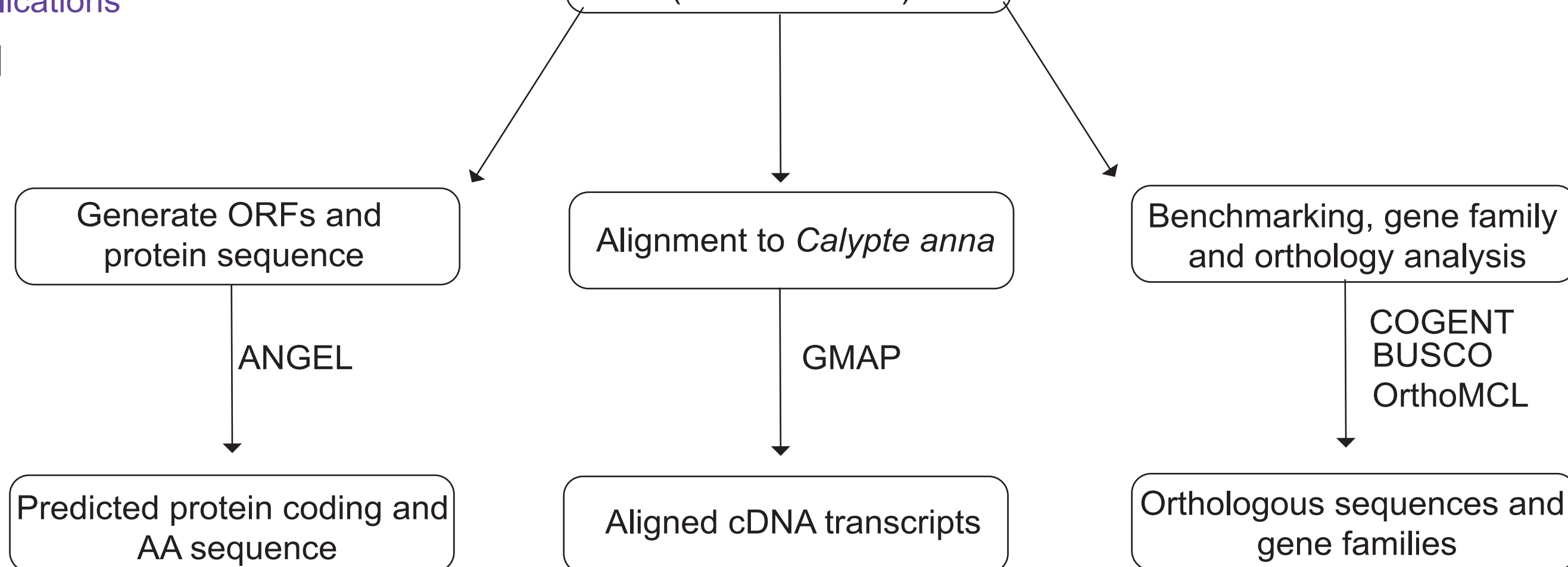
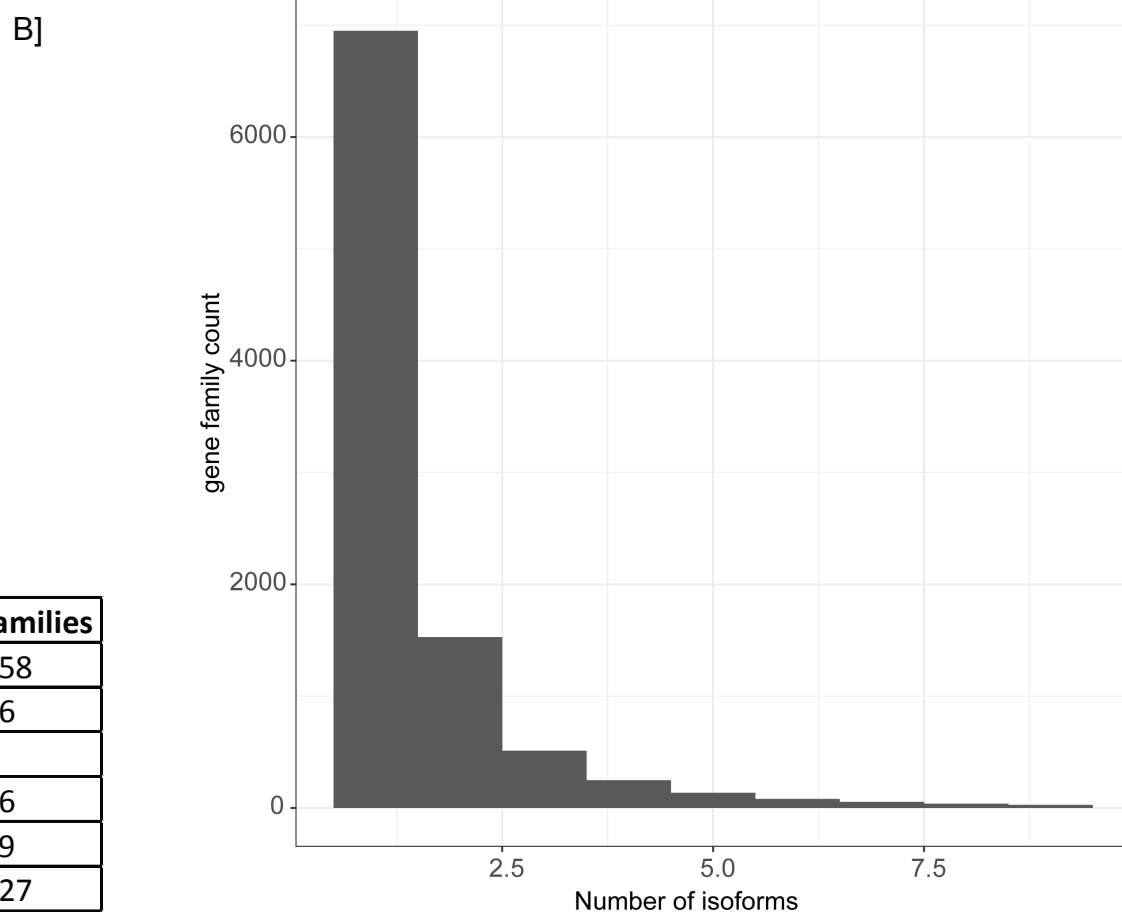


Figure 2. Analysis pipeline. **A** Raw sequence reads from a Pacbio RSII sequencer (bax.h5, bas.h5) were sorted into full and non-full length reads using a classification algorithm that identified full length reads with forward and reverse primers, as well as a poly-A tail. Iterative clustering for isoforms (ICE) was performed on full length reads, and non-full length reads were recruited to perform ARROW polished on the consensus isoforms. Polished sorted reads into high and low-quality bins, and either high quality data (HQD), all sequence data (ASD) or both sets of data, were carried on to further applications (**B**).

A)

Clustering results	Counts		
Total HQ Arrow Isoforms	94724		
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

Figure 3. Reducing transcript redundancy and predicting gene families with COGENT. Cogent gene families predicted and classified by relationship to *Calypte anna* genome assembly **A**. Number of full-length reads which support each isoform reduced **B**. C IGV view of the MATR3 gene, which was reduced from 11 redundant reads to 3 unique isoforms using this pipeline.

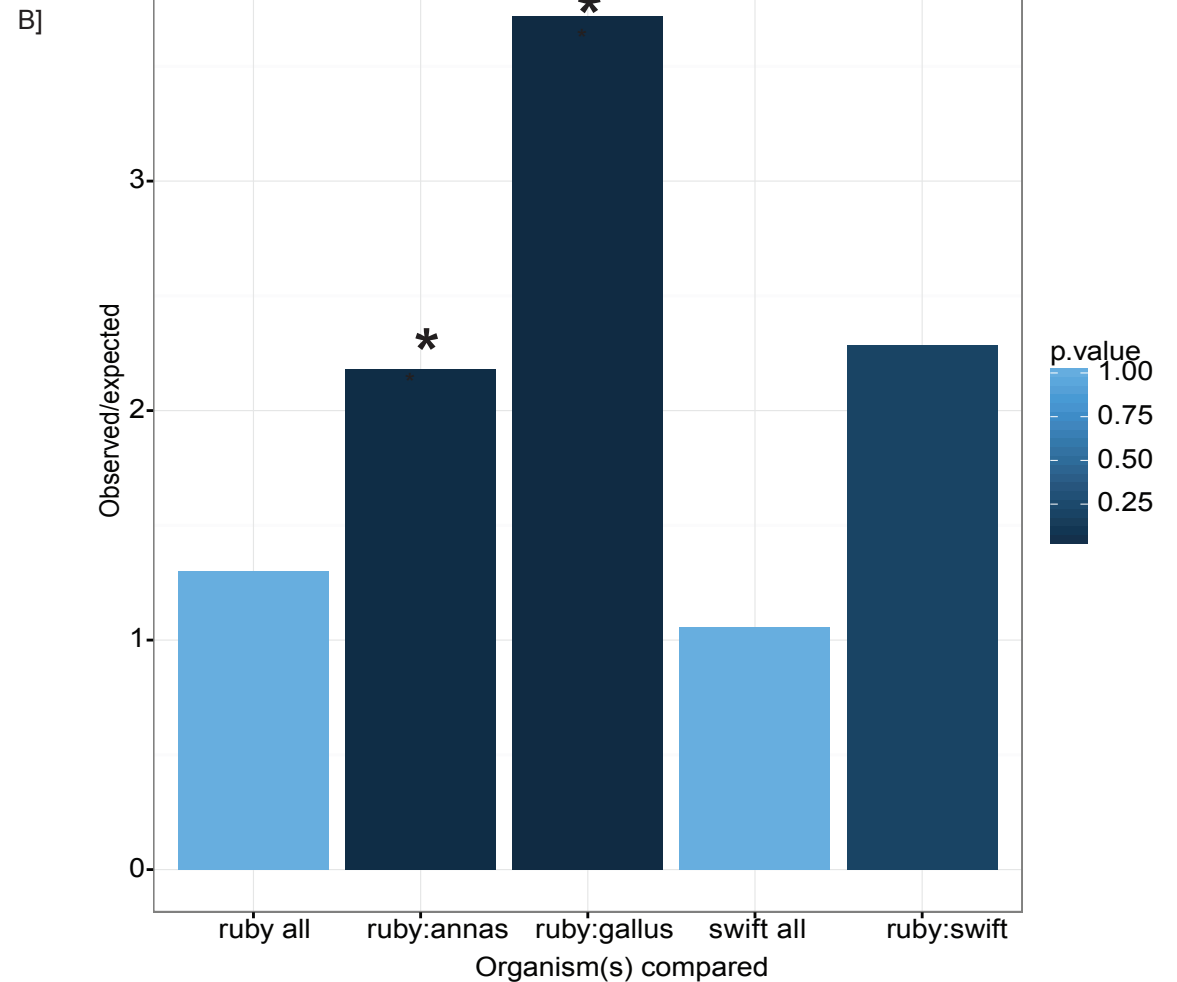
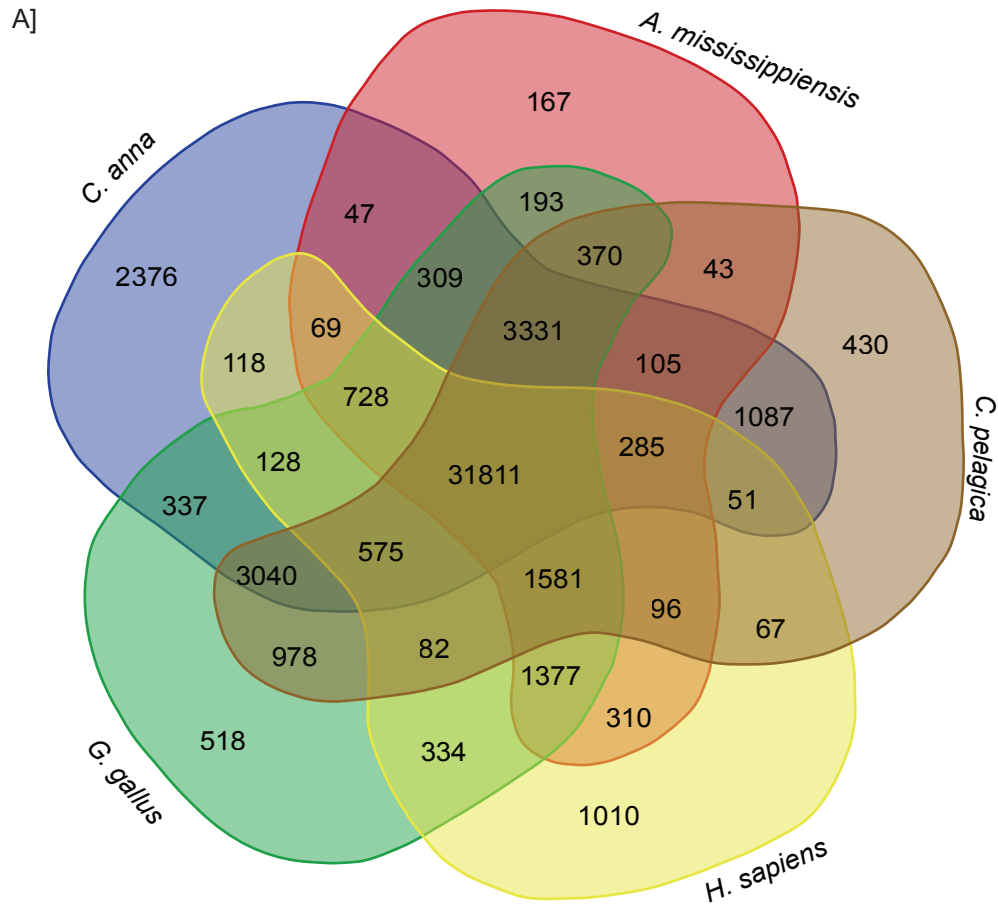


Figure 4. Orthology analysis. The transcriptomes of three birds (*Calypte anna*, *Gallus gallus*, *Chaetura pelagica*), one mammal (Homo sapiens) and one reptile (*Alligator mississippiensis*) were compared against *Archilochus colubris* using OrthoMCL to detect and compare similar sequences. **A** Venn diagram illustrating sequences with reciprocal blast hits between the given species and *A. colubris* is shown in A. Ortholog pairs unique to hummingbirds *Calypte anna* and *Archilochus colubris* were selected for gene orthology (GO) annotation analysis, which revealed enzymes of many biological functions (**B**).

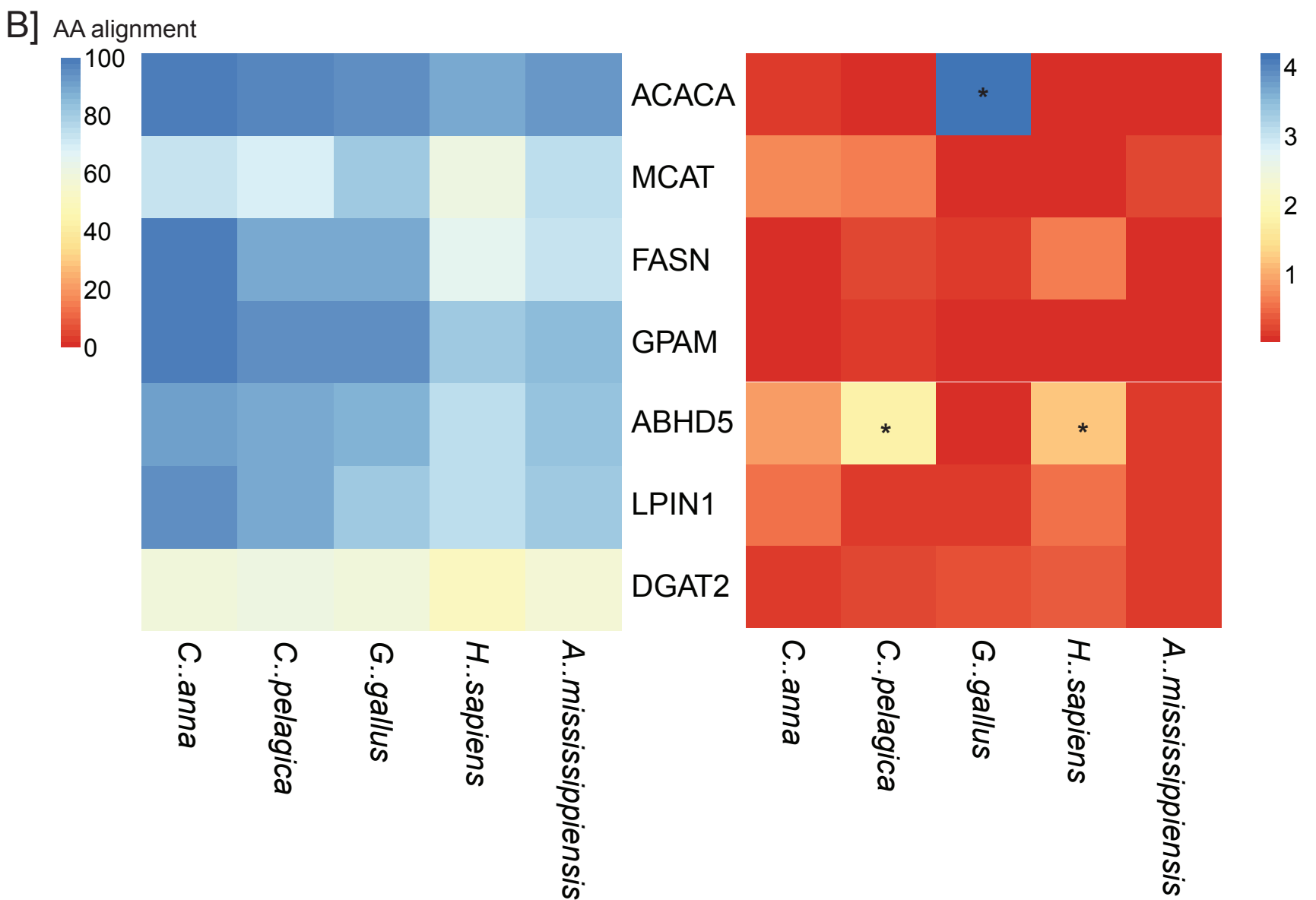
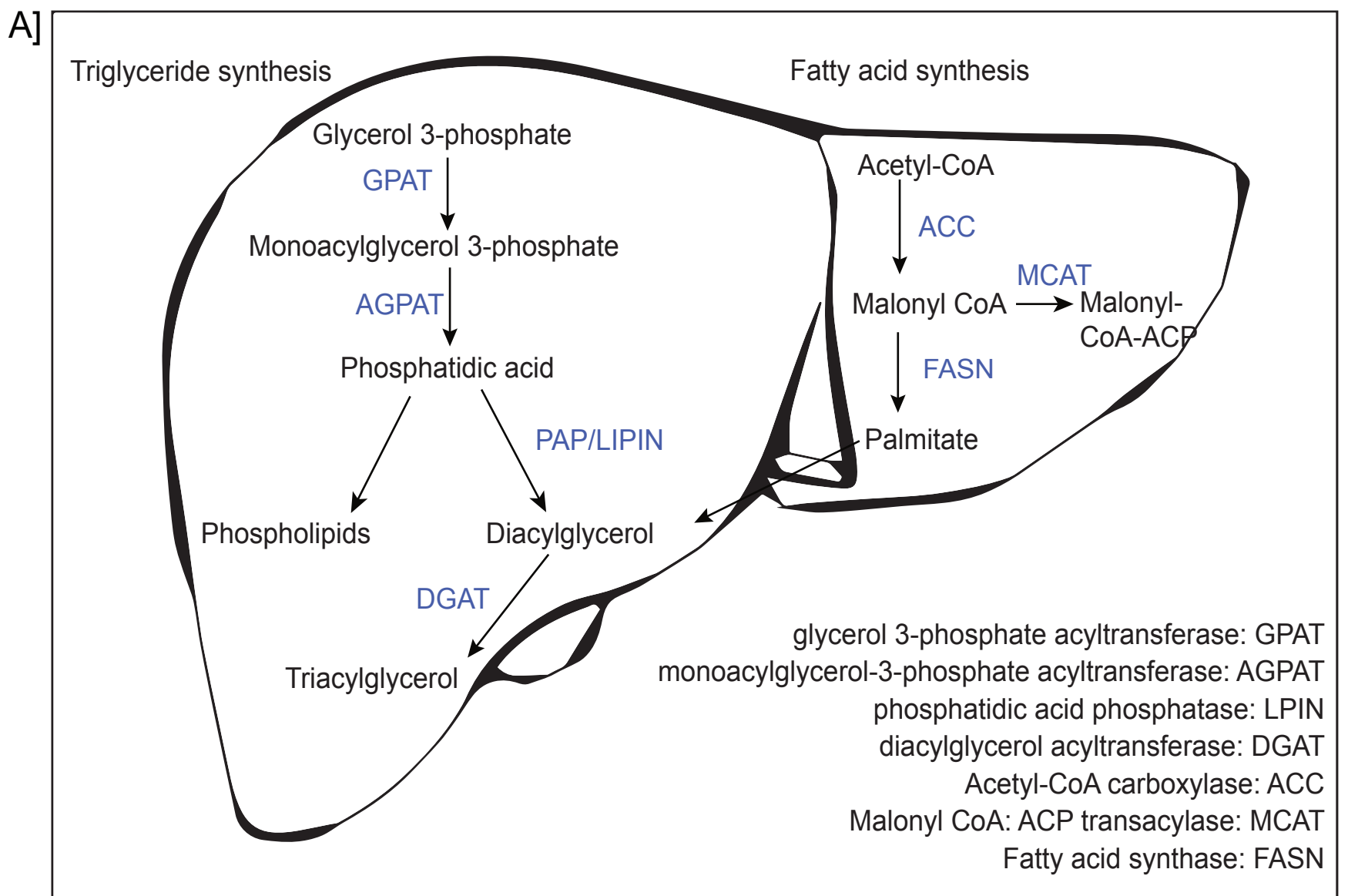


Figure 5. Pathway analysis of key enzymes in hepatic lipogenesis. **A)** An overview of the relationship between the investigated genes and their roles in triacylglycerol, phospholipid and fatty acid synthesis, **B)** a heat map illustrating percent identity of these proteins relative to *Archilochus colubris* predicted sequences, and nucleotide coding sequence conservation scores predicted using PAML (dN/dS), with genes dN/dS > 1 starred.



[Click here to access/download](#)

Supplementary Material

[MS_isoseq_hummingbird-cleancopy.pdf](#)





Click here to access/download

Supplementary Material

Suppdata_Workmanetal_7.pdf

