
Lab Notebook

Postdoctoral Research

Sonal Singhal

sonal.singhal1@gmail.com

Beginning 7 July 2014

Contents

Monday, 7 July 2014	1
1 Variant calling.	1
2 Data management.	1
Tuesday, 8 July 2014	2
1 Data management.	2
2 Variant calling.	2
3 Making masked genomes.	2
4 Phasing variants.	3
Wednesday, 9 July 2014	5
1 Making masked genomes.	5
Thursday, 10 July 2014	6
1 Phasing variants.	6
Friday, 11 July 2014	8
1 Data management.	8
2 Phasing variants.	8
Weekend and Monday, 14 July 2014	10
1 Phasing variants.	10
2 Generating recombination map.	10
Tuesday, 15 July 2014	11
1 Phasing variants.	11
2 Generating recombination map.	11
Wednesday, 16 July 2014	12
1 Phasing variants.	12
2 Generating recombination map.	12
Thursday, 17 July 2014	14
1 Generating recombination map.	14
2 Phasing variants.	14
3 Data management.	15

Monday, 7 July 2014

1 Variant calling.

After talking with Molly, I looked at the weird peaky behavior of VQSR filter scores for long-tailed and double barred finch. All the weird peaks could be attributed to variants that were fixed or nearly fixed.

Also, I wrote code `count_triallelic_sites.py` to look at multi-allelic sites in long-tailed finch. It turns out that a significant portion of these sites are actually bi-allelic because the reference allele (or, the zebrafish allele) is not at all represented.

2 Data management.

Issues with the /KG/ drive continued. Started to migrate data off and explored ways to make one centralized repository for data. See emails with John Zekos for information on what happened and how to avoid it.

Tuesday, 8 July 2014

1 Data management.

All data but a bare minimum of final VCF files was moved off the Columbia server. BAM files for doublebarred finch and longtailed finch were moved off the /KG/drive and onto my more secure home drive.

2 Variant calling.

Looking at the multiallelic sites in longtailed finch found the following:

- 3.52% of LTF variable sites have 2 or more alternate alleles reported.
- 1.34% of LTF variable sites have 2 or more alternate alleles reported, but only 2 of the alleles were genotyped individuals – i.e., no individuals had the reference allele.
- 2.18% of LTF variable sites are sites that are truly multiallelic.

3 Making masked genomes.

Each run of the script `make_masked_genomes.py` produces two files, the masked genome and a summary file that tells the user how many of each site type was produced. This information is best summarized in the README, a portion of which is duplicated here.

The masked genome is represented as a FASTA file with each site given as a numeric, mutua

The sites are coded as following:

- 0: coverage within acceptable bounds; no variation
- 1: coverage within acceptable bounds; variant that is HQ and no evidence of Mendelian di
- 2: coverage within acceptable bounds; variant that is HQ and has evidence of Mendelian d
- 3: coverage within acceptable bounds; variant that is low quality
- 4: unacceptable coverage; no variation
- 5: unacceptable coverage; variant that is HQ and no evidence of Mendelian distortion
- 6: unacceptable coverage; variant that is HQ and has evidence of Mendelian distortion
- 7: unacceptable coverage; variant that is low quality

Tuesday, 8 July 2014

Summing sites falling in categories 0 to 3 gives an estimate of total callable sites. Sum

I had made the doublebarred finch masked genome earlier last week. These were the files used.

```
vcf_file = '/mnt/lustre/home/sonal.singhal1/DBF/after_vqsr/DB.allchrs.vqsr.snps.indels.vcf'
cov_summary = '/mnt/lustre/home/sonal.singhal1/DBF/masked_genome/doublebarred_depth_summary.txt'
cov_data = '/mnt/lustre/home/sonal.singhal1/DBF/masked_genome/doublebarred_avg_depth.txt'
mendel_file = '\\"'
```

These are the files used to make the zebrafinch masked genome.

```
vcf_file = '/mnt/gluster/home/emleffler/genotype_callsets/zebrafinch/zf_unrels/unified_genotype.vcf'
cov_summary = '/mnt/lustre/home/sonal.singhal1/ZF/masked_genome/zebrafinch_depth_summary.txt'
cov_data = '/mnt/lustre/home/sonal.singhal1/ZF/masked_genome/zebrafinch_avg_depth.txt'
mendel_file = '/mnt/gluster/home/emleffler/genotype_callsets/zebrafinch/zf_family/unified_genotype.vcf'
```

Note that in zebrafinch, coverage was based on unrelated individuals only, because related individuals were sequenced to higher depth. Also, in the same vein, only variants from the unrelated were considered. Is this going to be a problem?

These are the files used to make the longtailed masked genome.

```
vcf_file = '/mnt/gluster/home/emleffler/genotype_callsets/longtailedfinch/after_vqsr/gatk.vcf'
cov_data = '/mnt/lustre/home/sonal.singhal1/LTF/masked_genome/longtailed_avg_depth.txt'
cov_summary = '/mnt/lustre/home/sonal.singhal1/LTF/masked_genome/longtailed_depth_summary.txt'
mendel_file = ''
```

4 Phasing variants.

To phase these chromosomes, we are going to use ShapeIt, first using their feature that lets you determine phase-informative reads (PIR) – these are mate-pairs that span two heterozygous sites. They tell you phase. Cool idea. However, to do this, we need VCFs that include SNPs and indels and that are by chromosome, and in this case, are also filtered by callable sites.

To do this:

1. Frustratingly, some VCFs are ordered by chromosome number (Chr1, Chr1A, Chr1B, etc) whereas others are ordered alphabetically (Chr10, Chr11, etc). This makes GATK and a number of programs unhappy, so the first thing I did was take a genome ordered alphabetically and create sequence dict and index files, using Picard and samtools respectively.

Tuesday, 8 July 2014

2. Then, it turns out that we didn't have a filtered VCF file for LTF that combined across both SNPs and indels and across the chromosomes. So, I wrote a little script `filtered_variants.py` that I used to pick out variable passing sites from the full genomic, all quality VCF from Ellen.
3. Then, I wrote a script that borrows from `make_masked_genome.py` and creates a VCF without inappropriate coverage (higher than $2\times$, lower than $0.5\times$ average genomic coverage) and splits it across chromosomes.

Wednesday, 9 July 2014

1 Making masked genomes.

It looks like all my masking of genomes has worked and that it is finally completed. Hurray! I put the README in each of the `masked_genome` folders and sent out the info to the group. Also, I copied the LTF directory to Columbia for Alva. Note that adding categories 0 through 3 gives the effective sequence length, which will be crucial for all other analyses.

To determine the final call sets, I took the (until that point) most final call set and removed sites with very low or very high coverage using the script `make_vcf_filtered_for_coverage_by_ch`. The VCFs used were the following:

- LTF: `/mnt/lustre/home/sonal.singhal1/LTF/after_vqsr/gatk.ug.ltf.allchrs.allvar.filtered`
 - Note: I had to create this VCF, because for whatever reason, it did not exist already.
 - I created it using `filtered_variants.py` on `/mnt/gluster/home/emleffler/genotype_callse`
- ZF: `/mnt/gluster/home/emleffler/genotype_callsets/zebrafinch/zf_unrels/unified_genotyp`
`/gatk.ug.unrelzf.allchrs.snps.indels.vqsr2.filtered.nomendel.recode.vcf.gz`

Thursday, 10 July 2014

1 Phasing variants.

I started my phasing experiments. There are four main ways that I could phase.

1. by using PIRs in ShapeIt and then LDhelmet
2. by using family information in ShapeIt and then LDhelmet
3. by using PIRs and family information in ShapeIt and then LDhelmet (is this even possible??)
4. by using LDhat straight away

I am going to pursue approaches 1 and 4 first, because I am still trying to figure out how to implement approach 3 and approach 2 seems like it would be less informative. Also, SNP calling for approach 2 was weird, and I am still getting my head around that.

In order to do approach 1, I got messed up by the extractPIRs program, but after some trial and error, I found:

- the pre-compiled version of extractPIRs works great; compiling on servers is hard because their version of g++ etc is so outdated
- you cannot have indels in the VCF file
- you cannot have non-biallelic SNPs in the VCF file
- VCFs have to be separated by chromosome
- although some of our multiallelic SNPs are really biallelic, I decided to just drop them, because keeping them would require recoding the VCF, which seems inadvisable
- I filtered the VCF files by using the script `remove_multialleles_indels.py`
- I really wanted to have proper BAM files for PIR calling, because it uses mate pair info. Due to the data problems, I didn't have one for LTF sample G118. I replicated KS's and EL's commands on LTF and moved forward with that.
- It turns out that ZF bam files and SNP calls are given two separate IDs, the intersection of which is in my local dir `/Users/singhal/zebrafinch/samples/zf_ids.txt`. Why would you do this?

Thursday, 10 July 2014

ShapeIT needs the BAM files to be indexed, so I did that, but it looks like G294 is truncated. So, copied that over again and indexed it, too.

I also realized that the ShapeIT program automatically defaults to N_e and ρ values for humans, which don't seem advisable for these birdies. So, I am going to do some work to approximate these values for birds.

Let's start with N_e . The easiest way to get a proxy estimate for N_e is to infer θ and to hope that the mutation rate estimate is reasonable. I looked through some papers, and the zebrafish mutation rate has been reported as $2.21 \times 10^{-9} \frac{\text{mutations}}{\text{site} \cdot \text{year}}$ (Nam et al 2010; doi:10.1186/gb-2010-11-6-r68) and $2.95 \times 10^{-9} \frac{\text{mutations}}{\text{site} \cdot \text{year}}$ (Balakrishnan and Edwards; doi: 10.1534/genetics.108.094250). Zebrafish have 3 - 4 generations a year, which means the per generation mutation rate is incredibly low – on the order of $7 \times 10^{-10} \frac{\text{mutations}}{\text{bp} \cdot \text{generation}}$. This paper (Ksepka et al. 2014; dx.doi.org/10.1098/rspb.2014.0677) suggests these estimates tend to predict much older divergence times than fossils, which makes me wonder if these mutation rates are downwardly biased. Plus, they are all based on the Zink calibration for mitochondrial rates, which I really wouldn't trust.

Anyways, I am going to go from θ estimates to N_e using Watterson's theta, so I need to calculate the number of segregating sites. (The VCF includes fixed sites and indels, which aren't appropriate for inclusion.) I wrote a script called `calculate_segregating_sites.py`.

Friday, 11 July 2014

1 Data management.

I lost most of the day to moving files around. I have 1 TB of space on the /mnt/lustre/ drive and 1 TB of space on the /mnt/glustre/ drive. That should be enough for everything, though it is definitely not ideal to have things spread out.

This is getting super annoying. This lab notebook is also starting to sound like a teenage diary.

2 Phasing variants.

I started my phasing experiments. There are four main ways that I could phase.

1. by using PIRs in ShapeIt and then LDhelmet
2. by using family information in ShapeIt and then LDhelmet
 - a) I really cannot make sense of how this was done already. Should I just go ahead and use it, or redo it in a way that makes more sense?
3. by using PIRs and family information in ShapeIt and then LDhelmet
 - a) An e-mail from Oliver Deleneau suggests that no, this is not possible. Bummer!
 - b) I suppose a kludge-y way to do this would be to phase family and then use that as a reference
4. by using LDhat straight away

Honestly, I expect these results are going to be robust across analysis types. I should probably do some sort of power analysis to check and see what my PIR power is going to be – calculating the average distance between biallelic SNPs should tell me a lot. Also, should I look at some kind of four gamete test like they did in Mimulus? A lot of this gets output by ShapeIT, so I should just look at that. It is non-intuitive to understand the reporting results, but it is a decent proxy.

Back to calculating N_e . For zebrafish, the total number of segregating sites is: 48726579. The total sequence length is: 859594837 bp. So, the segregating sites measure (S_n) is: $\frac{48726579}{859594837} = 0.0567$. Remember that $\theta = \frac{S_n}{a_n}$, where a_n is $\frac{1}{1 + \frac{1}{2} + \frac{1}{3} + \dots + \frac{1}{n-1}}$. Here, $n = 38$

Friday, 11 July 2014

because we have 19 diploid individuals. $\alpha_{38} = 0.2317$, which means that $\theta = 0.2447$. $\theta = 4N_e\mu$, and let's take $\mu = 1 \times 10^{-9} \frac{\text{mutations}}{\text{site} \cdot \text{gen}}$, which means $N_e = 61$ million. That's absurd. If mutation rates are more on the order of humans, then we get a more reasonable $N_e = 6$ million. This actually concords quite well with the Balakrishnan and Edwards value of $N_e = 7$ million. Similarly, for longtailed finch, $S_n = \frac{27995773}{897015175} = 0.0312$ and $\alpha_{40} = 0.229$, which means that $\theta = 0.136$. Assuming the same human-ish μ , $N_e = 3.4$ million.

Now for calculating ρ . I didn't know much about how ρ was calculated, or really what it meant. First, a centimorgan (c) is the distance between genes for which 1 out of 100 meiosis products are recombinant. Then, ρ is the expected number of crossover events per gene per base, and $\rho = 4N_e c$. The main tricky part of this equation is that c is the numerator of a fraction (the denominator is 100), so we need to take care of that and make it actual number before moving forward. In this case, the mean recombination rate in zebrafishes has been reported as $1.3 \frac{\text{cM}}{\text{Mb}}$ (Backstrom et al 2009; 0.1101/gr.101410.109). Converting to bp and turning it back into a fraction, $c = 1.3 \times 10^{-8} \text{meiotic products per bp}$. For zebrafish, $\rho = 4 \cdot 6e6 \cdot 1.2e-8 = 0.312$, and for longtailedfinch, $\rho = 4 \cdot 3.4e6 \cdot 1.2e-8 = 0.1768$.

These numbers might all be balderdash, but they are certainly quite different from the defaults, so I'll go with them for now.

So went ahead and got started, and got the dreaded underflow in sequencing error. I am trying four things:

1. rerunning with more RAM (20g) – still failed
2. getting rid of window size – still failed
3. getting rid of window size and theta / rho
4. just getting rid of theta / rho

I am redoing the same file (chr23) that worked before, so something is probably off.

This redos showed that for all that work looking at theta and rho, that is what was bugging the program! Will run with the incredibly wrong default values. This seems like a bad idea. Well, I am going to go forth with the bad idea, but I e-mailed Deleneau first. Another user had seen the same error, so I posted an update to that message.

IN GENERAL I AM NOT HANDLING THE Z CHROMOSOME WELL. WILL NEED TO RECONSIDER THIS.

Weekend and Monday, 14 July 2014

1 Phasing variants.

This weekend, started running ShapeIt based on PIRs for LTF and finished up the runs that failed for ZF. Some runs for ZF failed because they had insufficient memory. It is worth noting that the cluster does not generate the standard memory heap error when there is no more memory; rather, it just sorta shuts down.

2 Generating recombination map.

In order to run LDhelmet, the file formats I have now need to be parsed in many ways. The first is that I need a mutation matrix, which shows the stationary distribution of mutation rates between different bases. To do this, I followed Chan et al. 2012 (10.1371/journal.pgen.1003090) and started to implement the method by writing the script `get_mutation_matrix.py`.

Tuesday, 15 July 2014

1 Phasing variants.

Some of the initial ZF ShapeIt runs failed – for a random subset of the smaller chromosomes. I reran those with larger window sizes, hoping it was an issue with window being too small.

2 Generating recombination map.

I finished the `get_mutation_matrix.py` script – it took some finagling of the reference genome so that it was in the same format as my masked genome files – the new reference genome is `reference/taeGut1.60.bamorder.fasta`. I am currently running it on the ZF genome. If that works, I will repeat with the LTF genome. It worked, so I did it with the LTF genome.

I spent a lot of time trying to figure out how to use the De Maio, Schlotterer, and Koisel (2013) method (PoMo) to get ancestral allele states. Along the way, I had to install Python and finagle with that and associated issues with libraries and such. It ended up wasting most of the day, because it turns out that PoMo does the ancestral allele reconstruction, but there is no way to get that information out of the program. Bummer.

I then went back to our original idea of using the outgroups to get the ancestral state, though an informal counting scheme, essentially. I started to implement this in `simple_ancestral_allele`. There seems to be a lot of ancestral polymorphism, so it will be hard for this to be as exacting as I'd like. We'll see! Maybe it is fine given that everything is a prior, anyways.

Wednesday, 16 July 2014

1 Phasing variants.

Good progress on the phasing. Everything worked but for chromosome 20 in ZF. Pretty good! Need to figure out what happened there, but worth moving forward, at least.

2 Generating recombination map.

Another thought – maybe use the Darwin’s finch to try phasing? It is about 10 mya diverged, which is pretty far, but is better than nothing. Will look at mapping efficiency, and go from there. Downloaded the reads from the SRA using the sra-toolkit – wow, this is much better than just FTPing the reads to the server. Will definitely want to use it in the future. Used the reads in Project PRJNA178982 in SRA binary format and then dumped into FASTQ format using fastq-dump.

In the meantime, I prepared the genome using Stampy. I am using Stampy rather than my favorite (bowtie2) because it can handle divergent reads well, which I imagine will be relevant here.

```
~/bin/stampy-1.0.23/stampy.py --species=zebrafinch --assembly=taeGut1 -G taeGut1.bamorder  
~/bin/stampy-1.0.23/stampy.py -g ~/reference/taeGut1.bamorder -H ~/reference/taeGut1.bamorder  
/mnt/lustre/home/sonal.singhal1/bin/stampy-1.0.23/stampy.py -g ~/reference/taeGut1.bamorder
```

Here are the long-awaited mutation matrices! I calculated these as followed by Chan et al 2013. These are both reported as A, C, G, T for both rows and columns, and the directionality is from row to column. These look a lot like the Drosophila matrices published in the Chan paper, though slightly different.

Here is the matrix for the zebrafinch.

$$\begin{pmatrix} 0.455 & 0.104 & 0.322 & 0.119 \\ 0.206 & 0.001 & 0.135 & 0.659 \\ 0.659 & 0.135 & 0 & 0.206 \\ 0.119 & 0.322 & 0.103 & 0.455 \end{pmatrix}$$

Here is the matrix for the longtailed finch.

$$\begin{pmatrix} 0.437 & 0.103 & 0.344 & 0.117 \\ 0.205 & 0 & 0.151 & 0.644 \\ 0.644 & 0.151 & 0 & 0.205 \\ 0.117 & 0.344 & 0.103 & 0.436 \end{pmatrix}$$

Wednesday, 16 July 2014

These two matrices look pretty similar, as I would think.

Thursday, 17 July 2014

1 Generating recombination map.

Aligning the ground finch genome sequences to the zebrafish genome was going very slowly, and Stampy doesn't work in parallel unless you are running in BWA mode, so I created my own stupid parallel Stampy by splitting up the read file into 22 subfiles, and then running Stampy on 22 processes. It is still going pretty slowly, but at least this is a 22x speedup. I will then combine all the SAM files to get one master result, whenever it is finished.

2 Phasing variants.

I am not sure why chromosome 20 for zebrafish keeps failing, so I am going to try rerunning extractPIRs with greater stringency for quality, and see if that makes a difference. New command:

```
~/bin/extractPIRs.v1.r68.x86_64/extractPIRs --bam /mnt/lustre/home/sonal.singhal1/ZF/phas
```

Although it didn't influence the number of PIRs found or used, this run finished. Doesn't make sense, but I suppose it need not to.

I worked on getting the filtered for quality and coverage VCF files for longtailed finch into LDhat format. The main constraints of the program is that it can only handle variable, biallelic positions, so I had to get rid of any fixed or true polyallelic positions. Another fussy thing about LDhat – it cannot handle multiple variants at a site, so I had to dump any such variants. All of these tend to be when GATK calls a SNP and indel for the same position. I did this using the script `convert_vcf_to_ldhat.py`, which works on the files (all variants, not just biallelic) in the `/mnt/lustre/home/sonal.singhal1/LTF/after_vqsr/by-`

For LDhat, I need to get a likelihood lookup table, which is contingent on the θ for the species and the N (number of chromosomes) sampled for the species. To do that, I ran the complete program that is part of LDhat. The command I ran was:

```
~/bin/LDhat_v2.2/complete -n 40 -rhomax 100 -n_pts 101 -theta 0.136 -prefix LTF
```

The θ value I used was the same one I calculated using the number of segregating sites (Watterson's θ).

Thursday, 17 July 2014

I also wanted to try running ShapeIt using the duoHMM option, which allows you to use any level of pedigree information to phase chromosomes. Ellen had used an older version of ShapeIt that only allowed one trio or one duo. So, to do that, I had to first properly merge the unrel_zf and rel_zf files, because for some reason, they were called independently of each other. I am not sure why. To do this, I did the following:

1. Took my working ZF files to be the `*filtered.coverage.vqsr*` files in `/mnt/lustre/home/sonal.singhal`. These files are ideal because they only contain variable sites, they have both SNPs and indels (ShapeIt can handle both, when not run in PIR mode), and they do not contain any sites that fail coverage guidelines and Mendelian errors.
2. Took as the starting VCF for the un_rel_zf to be `/mnt/gluster/home/emleffler/genotype_callsets`. Note that this includes all sites; I need to be able to distinguish between non-variable sites and sites with missing data, and this is the only way to do that.
3. Took the un_rel_zf VCF and separated it by chromosome and filtered it for coverage and non-mendel sites using `make_vcf_filtered_for_coverage_by_chr_no_mendel.py`.
4. I then started writing a script that will take these two VCFs and merge them into a PED Plink output. I am not using VCF format because it would require me to recode the VCF to better handle indels and fakepolyallelic sites, and I don't want to do that.
5. VCF to ShapeIt PED format
 - a) family id (unique for every individual)
 - b) individual id
 - c) father
 - d) mother
 - e) sex
 - f) phenotype??
 - g) Genotypes
 - 1 - indel 1
 - 2 - indel 2
 - A,T,C,G, as expected
 - 0 = missing

3 Data management.

I set up my education GitHub account, which allows me to have private repos. I set up two repos, one for my scripts (<https://github.com/singhal/postdoc>) and the other for my lab notebook (<https://github.com/singhal/labnotebook>).