

# **Phylogenetic Resolution with mtDNA D-loop vs. HVS 1: Methodological Approaches in Anthropological Genetics Utilizing Four Siberian Populations**

By

Stephen M Johnson

Submitted to the graduate degree program in Anthropology and the Graduate Faculty of the University of Kansas in partial fulfillment of the requirements for the degree of Master of Arts.

---

Chairperson: Michael H Crawford, PhD

---

Bartholomew Dean, PhD

---

James H Mielke, PhD

Date Defended: 5 July 2013

The Thesis Committee for Stephen M Johnson  
certifies that this is the approved version of the following thesis:

**Phylogenetic Resolution with mtDNA D-loop vs. HVS 1: Methodological  
Approaches in Anthropological Genetics Utilizing Four Siberian Populations**

---

Chairperson: Dr. Michael H Crawford, PhD

Date approved: 5 July 2013

## ABSTRACT

Mitochondrial DNA is a useful genetic marker for answering evolutionary questions due to its high copy number, maternal mode of inheritance, and its high rate of evolution (Stoneking and Soodyall, 1996). The vast majority of research on mitochondrial DNA in anthropological studies has utilized the hypervariable segment 1 (HVS 1) to reconstruct population history and structure, explore population ancestry, construct phylogenies, and answer questions about the origins of prehistoric populations. A common debate in this field is whether better phylogenetic resolution can be obtained by the use of additional sequence data or genomic regions. If only the hypervariable segment is under scrutiny, does adding all three regions provide the same results or does it provide a deeper resolution, conveying a better understanding of populations of inquiry?

Sequence data from the D-loop of four Siberian populations: Altai, Evenki, Yakut, and Udehe have been analyzed using multivariate statistics in order to gain insight on evolutionary questions about these populations, and to investigate the utility and efficacy of sequencing the entire D-loop (HVS 1,2,and3 combined) versus sequencing solely the HVS 1. By comparing sequence data from the HVS1 to the whole D-Loop, this project investigated: 1.) if the increase in the number of SNPs sequenced revealed different phylogenetic relationships between Siberian populations; 2.) if additional genetic variation can be revealed by the addition of more genomic regions; and 3.) whether additional SNPs reveal stronger relationships between genetics, linguistics, and geography than using the HVS1 alone. Results of these statistics are consistent with previously reported findings for these populations based on HVS1 sequences data. The Altai, Evenki, and Yakut are predominately characterized by mtDNA lineages C and D, with various other Eurasian and East Asian lineages influencing their gene pool, whereas the Udehe from this study are solely characterized by the Native Siberian haplogroup C, and the Eastern Asian lineages of M, N, and Y. This is not surprising based their close geographical proximity to East Asian cultures. Results from this study indicate that the addition of the second and third hypervariable segment only sometimes aids in further characterizing the populations. In three of the four populations (Altai, Evenki, and Yakut), all of the haplotypes that were not characterized by HVS1 alone were resolved by adding HVS 2 and 3. However, no tangible differences were reported between the two sets of data for gene and nucleotide diversity. Likewise, no significant difference in results is gained from adding the HVS 2 and 3 data for AMOVA, Neighbor-joining trees based on population FSTs, neutrality tests or mismatch analysis. The addition of these extra genomic regions did, however, allow for a better resolution of population relatedness in multidimensional scaling, NJTs built directly from DNA sequences with bootstrapping (1000 replicates), and although it did not give higher correlation coefficients for mantel randomization tests, it did yield greater statistical significance.

This study and the conflicting findings therein suggest that the debate over whether or not resources should be used to analyze the entire mitochondrial D-Loop or simply the seemingly standard HVS1 remains a valid query. It is not necessarily as dichotomous a question as whether one is better than the other, but it depends on the research questions, the types of analyses conducted, and the sample sizes of the populations of study.

## **ACKNOWLEDGEMENTS**

I would like to first thank the indigenous people of the various Siberian populations that participated in this study, for without them, none of this research would be possible. Thank you, as well, to Dr. Larissa Nichols for her collaboration in collecting the Yakut samples, to Dr. Rem Sukernick for his work in collecting the Udehe samples, and to Dr. Michael Crawford and team for the collection and use of the Evenki and Altai samples.

I must extend my appreciation to Dr. Jodi Irwin at the Armed Forces DNA Identification Laboratory for managing the sequencing of the samples used for this project, and to Dr. Moses Schanfield at George Washington University for facilitating this collaboration.

I would also like to thank my committee members, Dr. Michael Crawford, Dr. Bartholomew Dean, and Dr. Jim Mielke, for their mentorship, advice and support throughout my time working on this project.

My boundless gratitude goes out to my friends and fellow (current and former) colleagues at the Laboratory of Biological Anthropology and the Department of Anthropology, specifically: Orion Graf and Dr. Norberto Baldi-Salas, Dr. Christine Phillips-Krawczak, Dr. Anne Justice, Jacob Boyd and Randy David for their collaboration, support, patience with my endless questions, and invaluable help along the way.

Thank you to the Department of Anthropology office staff: Carol Archinal, Le-Thu Erasmus, and Kathleen Womack for making all of our lives easier and for being so knowledgeable, kind, and helpful, to the Department of Anthropology faculty for providing me with the knowledge and background needed to advance my career, and to Dr. Crawford especially, my adviser and mentor.

I wish to thank my family, friends, and loved ones, for always encouraging me to dream big and helping me to achieve my goals in life. Your love and continued support mean the world to me.

I could not have done this without any of you.

## TABLE OF CONTENTS

Chapter 1: Introduction .....	1
Chapter 2: Background.....	5
2.1 Molecular Genetics.....	5
2.1.1 Mitochondrial DNA.....	6
2.1.2 Siberian mtDNA Review.....	8
2.2 Population Background.....	12
2.2.1 Gorno Altai.....	13
2.2.2 Evenki.....	14
2.2.3 Yakut.....	15
2.2.4 Udehe.....	17
Chapter 3: Materials and Methods.....	19
3.1 Materials.....	19
3.2 Lab Methods.....	20
3.3 Analytical Methods.....	20
3.3.1 Haplogroup Prediction.....	20
3.3.2 Within Population Variation.....	21
3.3.3 Among Population Variation.....	22
3.3.3.1 AMOVA.....	22
3.3.3.2 Multi-dimensional Scaling Plots .....	23
3.3.3.3 Neighbor-Joining Trees.....	24
3.3.4 Measures of Forces of Evolution.....	26
3.3.4.1 Neutrality Tests.....	26
3.3.4.2 Mismatch Analysis.....	27
3.3.5 Phylogeographic Methods.....	28
3.3.5.1 Mantel Randomization.....	29
Chapter 4: Results.....	31
4.1 Haplogroup and Haplotype Results.....	31

4.2 Hypervariable Segments Sequencing.....	37
4.3 Within Population Variation.....	39
4.3.1 Genetic Diversity.....	39
4.4 Among Population Variation.....	41
4.4.1 AMOVA .....	41
4.4.2 Multi-dimensional Scaling Plots .....	45
4.4.3 Neighbor-Joining Tree.....	48
4.5 Forces of Evolution.....	50
4.5.1 Neutrality Tests.....	50
4.5.2 Mismatch Analysis.....	52
4.6 Phylogeography.....	67
4.6.1 Mantel Randomization.....	67
4.6.2 Mantel Randomization and D-loop vs. HVS 1.....	68
Chapter 5: Discussion.....	69
5.1 mtDNA Lineages and Within Group Variation.....	69
5.2 Variation Among Populations.....	71
5.2.1 AMOVA.....	71
5.2.2 Multi-dimensional Scaling.....	71
5.2.3 Neighbor-Joining Trees.....	73
5.3 Forces of Evolution.....	75
5.3.1 Neutrality Tests.....	75
5.3.2 Mismatch Analysis.....	77
5.4 Phylogeography.....	79
5.4.1 Mantel Randomization.....	79
Chapter 6: Conclusions.....	81
Bibliography.....	85
Appendix.....	94

## CHAPTER 1: INTRODUCTION

Genetics can be defined as the study of heredity in organisms (Relethford, 2012). The field of population genetics extends these concerns of heredity to the population level, seeking to investigate the distribution of genes of a population and how its genetic composition changes throughout time due to forces of evolution. Anthropological genetics often employs a comparative approach to characterize the distribution of genes of multiple populations by explaining the variation within populations and between populations, all the while taking into consideration demographic events and the forces of evolution events (i.e., gene flow, natural selection, mutation and genetic drift).

In anthropological genetics biological markers are utilized to answer evolutionary questions about populations. These markers are “discrete segregating, genetic traits which can be used to characterize populations by virtue of their presence, absence, or high frequency in some populations and low frequency in others” (Crawford, 1973; Crawford 2007). After the advent of the molecular revolution in the 1980s, molecular markers became the focus of study for an increasing amount of research. Those markers based on blood groups and white blood cells are now often referred to as classical genetic markers (Crawford and Workman, 1973; Relethford, 2012). Molecular markers have become incredibly powerful tools in population studies because it allows the researcher to study patterns of inheritance of traits, both uni- and bi-parentally, from parents to offspring. Bi-parentally inherited autosomal markers are found in our nuclear DNA, or our chromosomes, of which we have 23 pair. Of these, the two sex chromosomes (X and Y), are transmitted uni-parentally, while the Y-chromosome is restricted to males and is passed solely from father to son, and the two X-chromosome are passed strictly from mother to offspring. This is highly useful for population studies in anthropology because it allows the researcher to

effectively trace inherited mutations and maternal and paternal lineages with greater detail (Crawford, 2007; Jobling et al., 2004).

In addition to our nuclear DNA, other molecular markers that are frequently utilized are found in the mitochondrial DNA (mtDNA). Like the Y-chromosome, this genetic material does not undergo recombination and is passed uni-parentally to offspring. However, in the case of this molecule, the DNA within the mitochondrion is passed maternally to all offspring (Francalacci et al., 1999; Crawford; 2007). There are hundreds of copies of mitochondria in each cell in our bodies and the mtDNA within them has a much higher mutation rate (some ten times higher) than autosomal DNA. In fact, there is a segment of the molecule called the D-loop that consists of three hypervariable segments (HVS 1, HVS2, and HVS3) that accumulate mutations at a faster rate than the rest of the mtDNA genome. Noting these features, mtDNA is a useful molecule for studying inheritance patterns and matrilineal ancestry information about a population. Since most nuclear DNA is too degraded in ancient samples, these features also make mtDNA highly desirable for ancient DNA studies where contamination and degradation are important factors to overcome (Crawford; 2007; Fancalacci et al., 1999; Stoneking and Soodyall, 1996; Stoneking, 2000).

Sequencing the entire mtDNA genome is an increasingly common practice in population genetic studies, especially as new technology helps to make this methodology more affordable. For well-funded laboratories, the high cost is typically not a problem. However, the costs associated with conducting DNA sequencing still remains high, especially when dealing with large sample sizes or long segments of DNA. The larger the sequence, the more segments need to be analyzed, utilizing more resources. The vast majority of research on the mitochondrial DNA molecule has been done on the hypervariable segment 1 (HVS 1). However, there exists the idea that the more genomic data you have, the more information you will be able to gain about your samples and the population(s). A study conducted by

Non et al. (2006) suggested that the more loci you have the better accuracy you will obtain with your results and the higher degree of phylogenetic resolution. If only the hypervariable region is under scrutiny, would similar results be obtained if the second and third hypervariable segment were also investigated? With resource limitation a factor, is the seemingly standard HVS 1 sufficient, or should one strive to include all three hypervariable segments (the whole D-loop) when using DNA sequencing to answering evolutionary questions about populations? The implications of this question could potentially affect financial and material resources of a laboratory, as well as lead to a better phylogenetic representation of the populations of interest.

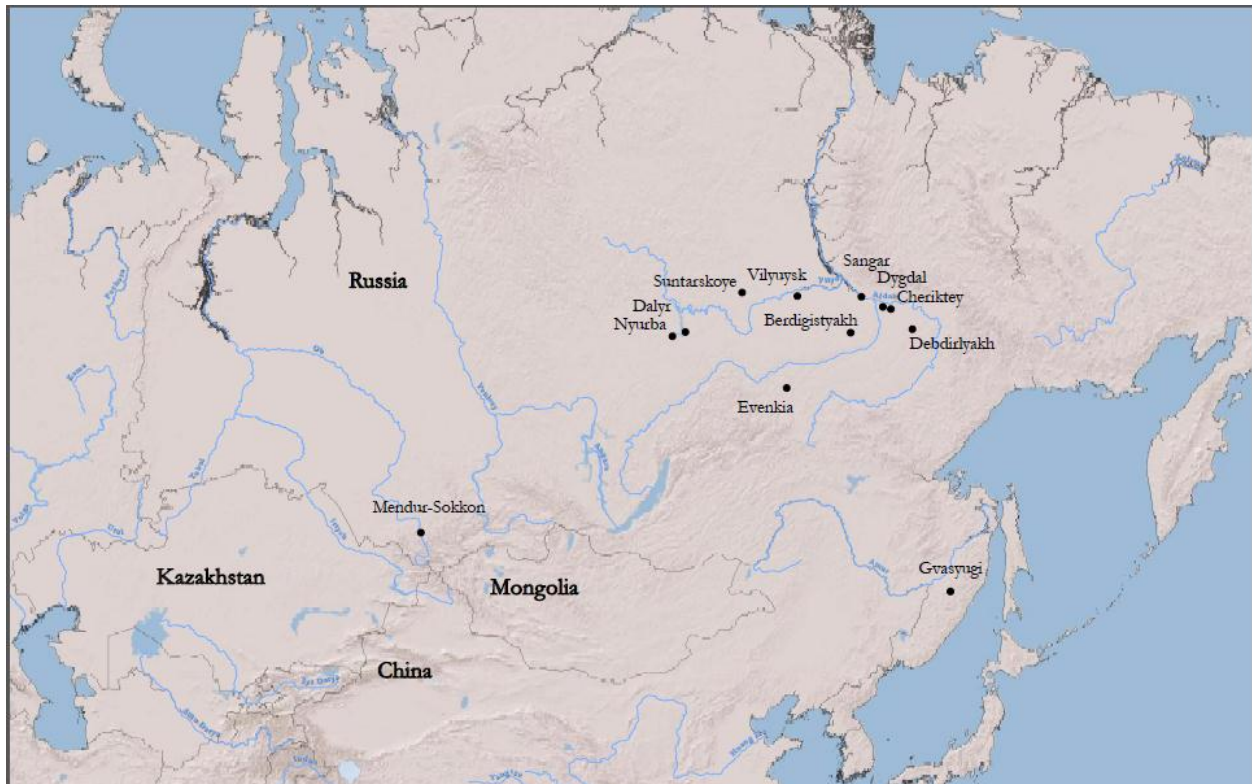


Figure 1.1 Map of Siberia and locations of populations in this study. This map was made with ArcGIS and edited with Adobe Illustrator CS6.

This project analyzes mtDNA sequence data from the D-loop of four Siberian populations: Altai, Evenki, Yakut, and Udehe, (*Figure 1.1*) and seeks to answer the question of whether it is better for scientific inquiry as well as a better use of resources to sequence the entire D-Loop (HVS 1, 2, and 3) or if the HVS 1 alone is sufficient for population studies analyzing mtDNA. Multivariate statistics were used to investigate the following questions:

- 1.) Whether the increase in the number of Single Nucleotide Polymorphisms (SNPs) sequenced reveals different phylogenetic relationships between Siberian populations;
- 2.) If additional genetic variation can be revealed by the addition of more genomic regions; and
- 3.) Whether additional SNPs reveal stronger relationships between genetics, linguistics, and geography than using the HVS1 alone.

The subsequent chapters provide relevant background information in order to put this research in perspective, and present the results in a discussion of the aforementioned questions. Chapter two contains a review of literature of mitochondrial DNA and its uses in population genetics, the history of the populations being compared in this study, and an overview of previous mtDNA research on the populations. The third chapter highlights the sample collection methods, laboratory and analytical methods used in this study. The results of the project are presented in Chapter four and discussed thereafter in Chapter five. The conclusions of this study are presented in Chapter six.

## **CHAPTER 2: BACKGROUND**

Chapter 2 summarizes the literature on the application of mitochondrial markers in anthropological population genetics, and introduces the populations utilized in this study. This information prepares the reader for understanding the analytical methodologies while providing a framework on which they are based.

### **2.1 Molecular Genetics**

A typical somatic cell houses 23 pair of chromosomes in its nucleus, which are collectively called nuclear DNA (Crawford et al., 2007). There is the one set (of the 23 pair) per nucleus, and only one nucleus per cell. Historically, working with this DNA molecule has been challenging because of the ease at which it degrades, the propensity for contamination, and the minute quantities in which it is found. Scientists today still deal with these hurdles but are better able to overcome them by the progress and contributions other researchers have made along the way. When James Watson and Francis Crick discovered that the structure of the deoxyribonucleic acid (DNA) molecule was composed of two helical chains that were each coiled around the same axis (1953), they laid the groundwork for every molecular biologist and geneticist who would follow in their footsteps. In the subsequent decades, researchers discovered: the concept of the Y-chromosome being the sex-determining chromosome, demonstrated that 85% of genetic diversity in human beings is within a population as opposed to between populations, and the development of DNA fingerprinting techniques using Variable Number Tandem Repeats (VNTRs) for forensic work (Crawford, 2007). In 1985, another groundbreaking innovation powered the molecular revolution when Kary Mullis et al. conceived a technique called Polymerase Chain Reaction (PCR). This method exponentially amplifies the amount of DNA through a cyclical process of denaturing, hybridization and extension of the template DNA strand. Each phase is catalyzed by different

temperatures, allowing the polymerase enzyme and oligonucleotides to replicate a target sequence of DNA, making downstream analysis more efficient and easier to carry out (1985). PCR has made it possible to examine the sequences of the smallest quantities of DNA and other molecules. Innovations and discoveries like these have allowed science to delve into deeper, more specific questions about human biology and origins. Without them, the research for this thesis would not be possible.

### **2.1.1 Mitochondrial DNA**

In modern population genetics research, studies based on mitochondrial DNA (mtDNA) and Y-chromosome DNA are an excellent way of illustrating population structure while tracing uni-parental inheritance and ancestry—mtDNA is maternally inherited while the Y-chromosome is paternally inherited. Mitochondrial DNA is considered highly valuable in answering evolutionary questions about populations due to its high copy number (around 2,000 mitochondria per cell on average) it's strictly maternal mode of inheritance, as well as its high rate of evolution (Stoneking and Soodyall, 1996). The mitochondrion is an organelle that occurs in vast quantities in the majority of cells in the body and provides cellular energy in the form of adenosine triphosphate (ATP). Unlike the double-helix form of the nuclear DNA, the mtDNA is a circular molecule with only 16,569 base pairs, and is present in the cytoplasm of unfertilized ovum during the reproduction cycle. Due to this pattern of maternal inheritance, mtDNA does not undergo recombination, upon fertilization the zygote will contain solely the mother's mtDNA, making it ideal for analysis of the evolutionary history of the maternal line of the individual (1996). This lack of recombination provides a predictable rate of mutations that can be used to establish a chronology for divergence of individuals and populations. This utility has led mtDNA to become a staple in the anthropological geneticists' toolkit for: reconstructing population history and structure, conducting ancestry research, constructing phylogenies, as well as working with ancient DNA.



by sequencing the mtDNA or through the use of restriction fragment length polymorphisms (RFLPs) (Young, 2009). RFLPs utilize restriction enzymes that can recognize the presence or absence of specific polymorphic DNA regions, and cut sites in the coding region of the mtDNA. These restriction sites correspond to the presence or absence of single nucleotide polymorphisms (SNPs). SNPs are then in turn used to characterize diversity in a population by categorizing individuals into mtDNA-haplogroups (Crawford, 2007; Young, 2009). A haplogroup is a package of inherited mutations. In this case, mtDNA-haplogroups are based on SNPs that are passed maternally to offspring, inherited as a unit, creating traceable mitochondrial lineages. These units of heritable mutations are labeled as mitochondrial haplogroups A-Y, and are regionally or geographically derived, as we can trace the origins of each haplogroup, through maternal lineages, all the way to their inception within specific parts of the world. The mt-haplogroups that are found in the New World are A, B, C, D and X. Mitochondrial haplogroups C and D occur at the highest rate in South America. Interestingly, these founding units of mutations occur in Northeastern Asia, such as Siberia, which would support theory that Native Americans entered the New World via the Land bridge during the Last Glacial Maxima (LGM), (Crawford et al., 2007; Fagundes et al., 2008; Rubicz et al., 2001).

### **2.1.2 Siberian mtDNA Review**

The major haplogroups that are found in Siberia are: A, B, C, D, F, G, and X with varying frequencies of other Eurasian haplogroups. *Table 2.2* below highlights the various diagnostic sites for each of these, respectively. All data in the table below are referenced in [www.mitomap.org](http://www.mitomap.org), and relative to the Revised Cambridge Sequence (Anderson et al., 1981; Andrews et al., 1999).

**Table 2.2:**  
**Diagnostic Polymorphisms of Major mtDNA Haplogroups of Siberian Populations**

Haplogroup	Diagnostic HVS1 Transition	Diagnostic RFLP sites with <i>Endonucleases</i>
A	16223, 16290, 16319	+663 <i>HaeIII</i>
B	16189, 16219	COII-tRNA(Lys) intergenic 9 bp $\Delta$
C	16223, 16298, 16327	+10394 <i>DdeI</i> , +10397 <i>AluI</i> , -13259 <i>HincII</i> , +13262 <i>AluI</i>
D	16223, 16362	-5176 <i>AluI</i> , +10394 <i>DdeII</i> , +10397 <i>AluI</i>
F	16304	-12406 <i>HpaI</i>
G	16223, 16227, 16278	+4830 <i>HaeII</i> , +4831 <i>HhaI</i> , +10394 <i>DdeI</i> , +10397 <i>AluI</i>
X	16189, 16223, 16278	-1715 <i>DdeI</i> , +144665 <i>AccI</i>

Torroni et al. (1993) were the first investigators of Siberian mtDNA diversity. The authors of the study investigated 411 individuals from ten different Siberian populations and analyzed them for RFLPs with endonucleases and mtDNA control region sequencing of the hypervariable region 1. Two of these populations (Evenki n=51 and Udehe n=46) are also studied in this thesis. Results indicated that of the four haplogroups present in Native American populations (A, B, C, and D), three of the four (A, C, and D) were found at high frequencies in the Siberian populations of study, collectively at >64%. For the Evenki, 3.9% of the sample was represented by haplogroup A, 84.3% by haplogroup C, and 9.8% by haplogroup D, whereas the Udehe showed much different frequencies (19.6% haplogroup C and 80.4% "other." Of these classified as "other," 63% were characterized by the indel *DdeII* np +10394 and over half of these containing the indel *AluI* np +10397, which have been shown to suggest other common East and Southeast Asian sub-haplogroups (Ballinger et al., 1992).

Another early mtDNA study of variation in Siberian populations was conducted by Derenko and Shields (1997 and 1998). These authors investigated a sample of Yakut, sequencing the HVS1 and Region V of the coding region. Results revealed Yakut-specific haplotypes, high sequence diversity and

presence of 9bp deletion in region V, common to Haplogroup B (4.6%). However, the sample size was only 22 individuals, so the study likely does not reflect accurate diversity measures within the population. Similar studies have also shown that haplogroups A, C and D are present at high frequencies in Northeastern Siberian populations, whereas haplogroup B is common in Central Asia and southern Siberia such as in the Altai (also in this study), Tibetans, Mongolians, and the Buryat with percentages >8%(Ballinger et al., 1992; Kolman et al., 1996; Torroni et al., 1993; Schurr and Wallace, 1999; Shields, 1992, 1993; Sambuughin et al., 1991; Starikovskaya et al., 1998). This has been a significant anthropological discovery, as it supports the notion that Siberian populations were the ancestral founding sources for Native American populations.

Haplogroup B does not appear to be restricted to Southern Siberia. In 1993 Petrishchev et al. discovered the 9bp deletion, characteristic of this haplogroup at low frequencies in the Northeastern Siberian populations of the Nanai and Evenki. This marker has also been found at low frequencies in the Yakut by various studies (Derenko and Shields, 1998; Puzyrev et al., 2003; and Federova et al., 2003).

The magnitude of mtDNA diversity appears to be higher in central and Eastern Siberian populations compared to those of the North. In 2003, Derenko et al. investigated 480 Southern Siberian samples from seven populations: Altai, Khakassians, Buryats, Sojots, Tuvinians, Todjins, and Tofalars. RFLP and HVS1 sequencing revealed that 81% of the haplogroups represented were classified as Asian (C, D, G, Z, A, B, F, N, Y), whereas, 17% belonged to Western Eurasia (H, U, J, T, I). This pattern reflects the diversity seen in populations from East and Central Asia. On the other hand, Northern Siberian populations tend to represent high frequencies of haplogroups C and D, with smaller, varying amounts of Western Eurasian haplogroups (Derbeneva et al., 2002).

The Yakut, one of the populations studied in this thesis, have shown similar levels of diversity as those in Northern Siberia. Sequencing and RFLP of 83 Yakut revealed most of the major Siberian

haplogroups (A, B, C, D, G, F, and M) at a frequency of 94%. The remaining 6% are represented by the Eurasian H, J, and U. Of the 94%, haplogroups C and D represented 80% of the diversity (Puzyrev et al. 2003). It may be interesting to note that of the 83 samples tested, only 22 haplotypes were characterized. This may be indicative of overrepresentation of family lineages, inbreeding, or due to sample error. A similar study, published the same year (2003) revealed similar results. HVS1 Sequencing of 117 Yakut individuals showed that haplogroups C and D represented 69.2% of the sample (Pakendorf et al.). The author notes that this showed close genetic affinities with Evenki and Tuva.

Further evidence for these patterns in the Yakut come from a study conducted by Federova et al. (2003), in which HVS1 sequences were constructed for 191 Yakut individuals from a DNA bank at the Department of Molecular Genetics at the Russian Academy of Medical Sciences in Yakutsk. The results indicated that 91% of the samples belonged to Eurasian haplogroups, 74% of which came from C and D. The others belonged to haplogroups H, J, T, U and W, which is a similar pattern of variation seen in Mongolian and Central Asian ethnic groups.

Though haplogroup X is one of the founding Native American haplogroups, it is documented at relatively low frequencies in the Americas, compared to the other four (A-D). It was originally thought to have been brought to the Americas from Europe—Greenland has a high frequency of haplogroup X—as it has not been found in any Asian population (Brown et al., 1998). In 2001, however, Derenko et al. conducted a study of approximately 790 individuals from ten different Siberian populations, 202 of which were Altai. The authors noted that haplogroup X had been present in the Altai at a frequency of 3.5%. It is important to note, however, that sequence analysis of those individuals who belonged to haplogroup X from the Altai, when compared to haplogroup X sequences from Native Americans, and Europeans, are all strikingly different phylogenetically, and are characterized as being from distinct subhaplogroups, leaving the exact origin of the Native American X yet to be fully understood (Bandelt et

al., 2003, and Reidla et al., 2003). The presence of haplogroup X was confirmed in the Altai yet again, however, when Phillips-Krawczak et al. (2006) used mtDNA RFLP markers and HVS 1 sequencing to identify one individual from those sampled (n = 98) as Haplogroup X. The haplogroup assignment of this individual was validated by results of the present study, as these 98 individuals from Mendur-Sokkon formed the majority of the Altai samples that were analyzed for this work.

With improvements in methodologies and better sequencing techniques of HVS 1 and 2, further subhaplogrouping has been done on Siberian and Native American populations, and as a result, the Native American founder subhaplogroups are now recognized as A2, B2, C1, D1, and X2a (Forster et al., 1996; Hernstadt et al., 2002; and Bandelt et al., 2003). Further, subhaplogroup, D2, has been discovered as another founder lineage from Siberia and not derived from the American D1 (Starikovskaya et al., 1998; Derbeneva et al., 2002; Rubicz et al., 2003; Zlojutro et al., 2006; Saillard et al., 2000). In review, central Siberian populations, such as the Altai, seem to experience a greater degree of genetic diversity, and show more Western Eurasian haplogroups, whereas northeastern populations like the Evenki and Yakut tend to be more representative of the Asian haplogroups A-D, with the highest frequencies falling in C and D. The Udehe are more disparate, as they contain lower frequencies of the major Siberian haplogroups A-D, as well as fewer Western Eurasian haplogroups, and a significant number of East Asian haplogroups, which is likely due to contact and close geographical proximity.

## **2.2 Population Background**

The following information will provide a concise ethnographic background on the populations used in this study. Information on their history, language, and culture and ecological territory will be highlighted so as to provide a cultural framework that can be used to aid interpretation of results of this study.

### 2.2.1 Gorno Altai

The Gorno Altai is a genetically diverse people indigenous to southern Siberia, and live in what is now the Altaic Republic, an autonomous region within the Russian Republic. Their homeland is situated in Eurasia, at the confluence of China, Mongolia, Kazakhstan and Siberia. Today the Gorno Altai region is inhabited by approximately 60,000 native peoples, composed of various Turkic speaking groups. Bowles (1977) hypothesized that the Turkic-speaking people originated in the Altai region, later spreading rapidly to adjacent geographical areas. This has been argued, however, suggesting that the Altaic language family more likely developed in the steppe region of western Siberia more than seven thousand years ago, and only later spread to the Altai region of today (Miller, 1991).

In the 12<sup>th</sup> century, Mongol hordes invaded the Altai, taking political control of the region. This dominion lasted into the 18<sup>th</sup> century with the fall of the Oirat Khanate (Chiefdom). At this point, the Chinese saw an opportunity, and began invading the Southern Altaic region. In the latter part of the century, Russians forced the Chinese out of area, assuming political control of the region. In the 1930s, the Soviets forced the Altai into farming and herding collectives (Forsyth, 1992). Today, the Kizhi of the Gorno Altai autonomous republic are still semi-nomadic, pastoralists. The Altaic valleys between surrounding mountains of the Altai-Sayan range provide exceptional grazing rounds for cattle (Crawford et al., 2002).

The tribes of the Altai region today are often divided into two groups based on language, culture, geographic distribution and genetics: Northern and Southern. The Northern tribes consist of: Chelkans, Kumadins, Tubalars, and Maimalars. Southern Ethnic groups include: Altai-Kizhi, the Teleuts, and the Telenghits (Derenko et al., 2003; Levin and Potapov 1964; Phillips-Krawczak et al., 2006; Crawford et al., 2002). The Kizhi of the Gorno Altai (used in the present study), are speakers of a Turkic language that is distantly related to Tungusic, spoken by the Evenki. The samples used for this study

were collected in Mendur-Sokkon, which is part of the Southern Altaic Region. Mendur-Sokkon is home to approximately 1000 individuals who are subdivided into three major patrilineal clans (Irkit, Kipchak, and Todosh), in addition to various smaller clans.

Because of the situation of the Altaic region, a convergence of migration routes from Europe and Asia, this area has experienced a significant population movements as well as invasions from Mongolia, China, and Russia, not to mention trade and cultural exchanges with these regions. This dynamic history suggests that the Altai display a mosaic of genetic and cultural motifs from Western Europe and Eastern Asia (Derenko et al., 2001, 2003; Phillips-Krawczak et al., 2006). The Altai have been of great significance to the field of Anthropology because it is the only place outside of the new world that all five of the founding Native American mtDNA haplogroups are located (A, B, C, D and X) (Derenko et al., 2001).

### **2.2.2 Evenki**

The Evenki whose name reportedly means “he who runs swifter than a reindeer” (Oaks and Riewe, 1998) are an indigenous people of Siberia that subsist as reindeer herders and breeders and are widely distributed throughout the taiga and boreal forests of central Siberia all the way to the Amur region of the East. Their territory is so widespread in fact, they were once thought to occupy one quarter of all of Siberia (Vasilevich and Smolyak, 1964). Reindeer are extremely important to the Evenki way of life. They are used as pack and riding animals and also supply food, clothing, and tents. In addition to herding reindeer, the Evenki hunt wild reindeer and moose, fish, trap, and for some in the southern sections of their territory, cattle herding is not uncommon (Oaks and Riewe, 1998; Torroni et al., 1993).

The Evenki speak a Tungusic language, which is a branch of the Altaic language family. Other languages that are considered Altaic are Turkic, Mongolian, and Manchu. Based on archaeological evidence, it is believed that the Evenki come from reindeer hunters from the Lake Baykal region who eventually expanded east, and in doing so began herding the animals (Okladnikov, 1964; Crawford, 2002). By the time of European contact in the 1600s, the Evenki occupied an expanse of territory from the Yenisey River to the eastern coast of Siberia, thousands of miles away. The Evenki distribution was bisected in the 16<sup>th</sup> century by the Yakut from South Siberia, as they began to migrate northward. (Crawford, 2002).

Traditionally the Evenki lived in patrilocal, patrilineal tribes that are headed by princely families. Until the 20<sup>th</sup> century, they lived in relative isolation, with the exception of trading with Russian merchants as well as the imposition of the fur tax. The Evenki may have been able to escape much of Russia's influence on their culture due to the mobile nature of their nomadic society. In the 1930s, however, the Soviets tried to consolidate herding groups, and in doing so many leaders and shamans were executed as a way of maintaining control over the tribes. Herds were collectivized and reindeer were assigned to brigades of young, unrelated Evenki males, replacing traditional family herding practices. These herds consisted of approximately 1500 animals. After the fall of the Soviet Union, extended herding families began reclaiming most of the reindeer herds (Crawford, 2002). Today, the Evenki number over thirty thousand individuals and their population continues to grow (Oaks and Riewe, 1998).

### **2.2.3 Yakut**

The Yakut of Northeastern Siberia are Turkic-speaking, horse and cattle breeders and herders, surrounded by Tungusic-speaking reindeer herders and hunter gatherers—the Evenki, for example (Zlojutro et al., 2009; Pakendorf et al., 2003; Balzer, 1994). The Saha, as they call themselves, are one of

the largest indigenous groups in Russia, numbering over 400,000 individuals (Tarskaia et al., 2006; Jordan et al., 2001). It is believed that the Yakut left the steppes of southern Siberia in the 1500s, traveling northward into the taiga, living among the Evenki around the Lena River. This region today is known as Saha, an independent republic in the Russian Federation. Most of this area lies above 60 degrees north, which makes it one of Siberia's colder environments (Oaks and Riewe, 1998).

Traditionally the Yakut were known to be horse breeders and cattle herders, and milk products were an important part of their subsistence. They are renowned for making the cows' milk into a type of soured condensed milk. In addition to being used as pack animals, the horse meat was also consumed. Because of their time among the Evenki, the Yakut have today adopted Evenki-style reindeer herding, wild moose and reindeer hunting, but certainly not to the extent of the Evenki (Armstrong, 1968; Forsyth, 1992; Okladnikov, 1970; Luick, 1978; Pakendorf et al., 1999).

The Yakut have a long history of metal working and smelting of iron from ore for forging. During the middle ages, Yakut warriors were armored like that of European knights. It is even said that they used armored war horses. Instead of chainmail like that of the Europeans, their armor was more fish scale-like, making it much more similar to East Asian armors. In ancient times, wealthy Yakut owned slaves that served as cowhands, servants, driver, butcher, or cooks. It was not an uncommon practice to bury the slave or wife alive with the slave owner (men), so as to ensure they would be served in the afterlife (Okladnikov, 1970; Oaks and Riewe, 1998).

Archaeological and ethnographic evidence suggest that the Saha originated from an ancestral population from the region around Lake Baikal and migrated northward, partially to evade the Mongol invasions (Schönig, 1990 Tarskaia et al., 2006; Forsyth 1992; Okladnikov, 1955). Rock paintings discovered on the northwest shores of Lake Baikal have been attributed to the Kurykan people, who were cattle and horse breeders. The decorations on horse bridles, riding dress, and style of pottery all

show remarkable similarities to that of the Yakut. This has been interpreted to suggest that the Kurykans were likely the ancestral population to the Yakut (Alekseev 1996; Kostanistinov 1975; Okladnikov 1970). Because of this migration from South to North and practice of exogamy, they are genetically related to both northern and southern Siberian populations. During their trek northward, the Yakut continued to expand their territory through intermarriage with other surrounding ethnic groups such as the Evens, Evenki, Yukagir, and Samoyed (Forsyth, 1992). Mitochondrial DNA evidence from HVS1 sequences have suggested that they may be closely related to the Tuvinians of the southern Siberia (Tungusic-speakers). Tarskaia et al. (2002) has shown that protein polymorphisms reveal a strong affinity between the Sakha and Altai, Mongols, and Buryats, as well as the Evenki, Even, and Chukchi. Pastoralism, clothing, festivals, and other aspects of the Yakut culture, as well as many Buryat and Mongolian words in the Yakut language, all support the idea of these ancestral ties to Southern Turkic peoples (Zlojutro et al., 2008; Tokarev and Gurvich, 1956).

#### **2.2.4 Udehe**

The Udehe (also spelled *Udegey* or *Udeghey*) are a small, isolated ethnic group of aboriginal hunters and fishermen of Southeastern Siberia who live along tributaries of the Amur, Rikin and Ussuri Rivers along the Amur River Basin (Torrioni et al., 1993). Until recent times, they occupied both slopes of the Sikhote-Alin mountain range south of the Amur River Basin. They are one of the least represented ethnic groups in the Russian Federation. According to the 2002 census, Russia was home to 2657 Udehe. Today, however, it is estimated that the population is no more than 1000 individuals (Han-Jun Jin et al., 2010). Though most of the Udehe today speak Russian, a Tungusic language, thought to be a subdivision of the Manchu Branch of the Altaic linguistic family, is traditionally spoken by this population (Torrioni et al., 1993; Jin et al., 2010; Levin and Potapov, 1964; Krauss 1988; Starikovskaya et al., 2005).

The Udehe mainly subsist by engaging in hunting, fishing, and ginseng collecting. Because of the heavy influence from surrounding Russians, most of them have adopted westernized lifestyles (Smolyak, 1975). The Udehe share many similarities, anthropologically, with the Ulchi, Nanay, and Orok, who are all self-assigned to be part of the Nani group (Petrova, 1967). Crawford et al., (2002) suggests that Udehe might be better classified as Tungusic and then viewed as the most easterly extension of the Evenki reindeer herders. It has been widely suggested that they have migrated to the territories of the Amur Basin from an area around Lake Baikal in Southern Siberia, which is also the prospective ancestral home of other Siberian groups such as the Evenki and Yakut (Smolyak, 1975; Alexeev and Gohman, 1984). Many connections between ancestral Udehe and East Asians exist. The direction and amount of genetic admixture between these groups remains unclear, however (Han, 2005). The Udehe have also been heavily influenced by migrations from further West. Migration networks from central Asia and Western Eurasia have been documented. The completion of the Trans-Siberian Railway in 1907 brought a large influx of Russian and Ukrainian settlers to the Udehe region (Alexeev, 1989; Vasiliev, 1993; Derevianko 1998). According to Ballinger et al. (1992), the Udehe appear to be more genetically diverse than expected, as they contain not only the major Siberian Haplogroups (A-D), and East Asian haplogroups, but a significant representation of haplogroups coming from Europe as well, which is likely due to geographical proximity, trade, and the Trans-Siberian railway. The results of this study, however, suggest that although they are predominately comprised of East Asian lineages (M, N, and Y), and only a small percentage Siberian, represented by haplogroup C.

## CHAPTER 3: MATERIALS AND METHODS

### 3.1 Materials

This chapter reviews the molecular and statistical methodologies employed in analyzing the samples used in this study. During various LBA field sessions (1980s-90s), DNA samples were collected in the form of whole blood from members of the Altai, Evenki, Udehe, and Yakut ethnic groups of Siberia. These samples, though originally used for anthropological genetic studies of population origin and structure (Torrioni et al., 1993; Shields et al., 1993; Phillips-Krawczak et al., 2006; Starikovskaya et al., 2003; Derenko et al., 2003; Kong et al., 2003; Zlojutro et al., 2008), have been used for the present study in order to elucidate whether there is a significant difference between results of various population studies using only HVS 1 sequences and those using all three hypervariable regions of the mtDNA D-Loop (HVS 1, 2, and 3). This study was approved by the Human Subjects Committee (HSCL) of the University of Kansas (20788#). Verbal Informed consent was given by participants for involvement in this research project. A sampling summary of the populations can be seen in *Table 3.1*.

*Table 3.1 Siberian Sample Collection Table*

<b>Ethnic Group</b>	<b>Locality/Village (Administrative Centre)</b>	<b>Sample Size (n)</b>	<b>Collected by</b>
Altai	Altai Mountains (Mendur-Sokon)	101	Dr. Michael Crawford
Evenki	Evenkia	53	Dr. Michael Crawford
Udegey	Gvasyugi	41	Dr. Rem Sukernick
Yakut(Ck)	Cheriktey	34	Dr. Larissa Nichols
Yakut(Da)	Dalyr	66	Dr. Larissa Nichols
Yakut(De)	Debdirge (Debdirlyakh)	61	Dr. Larissa Nichols
Yakut(Dy)	Dygdal	34	Dr. Larissa Nichols
Yakut(YE)	Egolja Nurbinsky (Nyrba)	26	Dr. Larissa Nichols
Yakut(YJ)	Jarkhan (Suntarskoye)	48	Dr. Larissa Nichols
Yakut(YK)	Kulyatsy Vilvyskiy (Vilyuysk)	103	Dr. Larissa Nichols
Yakut(YM)	Mukuchu Kobiyskiy (Sangar)	36	Dr. Larissa Nichols
Yakut(YO)	Orto-Surt Gorny (Berdigistyakh)	63	Dr. Larissa Nichols

### **3.2 Laboratory Methods**

DNA was extracted from samples using the phenyl chloroform extraction method and the *Super Quick Gene*® protocol, (LBA, University of Kansas). Extracted DNA was then sent to Jodi Irwin at the Armed Forces DNA Identification Laboratory (AFDIL) for sequencing of the hypervariable segments 1, 2, and 3 of the mitochondrial control region. Resulting sequence motifs were reported back to the Laboratory of Biological Anthropology at KU where full sequences were built with MEGA 5.10 (Tamura et al., 2007; Tamura et al., 2011), by altering the revised Cambridge Reference Sequence (rCRS) with the specified mutations for each sample (Anderson et al., 1980; Andrews et al., 1999).

### **3.3 Analytical Methods**

For each of the tests presented in this section, results were compared between HVS1 sequence data and that from the combined D-Loop consisting of HVS 1, 2, and 3 collectively, in order to investigate the question of whether a greater degree of information is gained from the availability of sequence data from all three hypervariable segments. The following combinations of populations and sub-populations were utilized in the subsequent analyses: 1.) Comparing the Kizhi Altai, Evenki, Udehe, and Combined Yakut populations; and 2.) Comparing the Altai, Evenki, Udehe, and the nine different Yakut populations separated into their respective geographical sub-populations, grouped by village. For the intra-population analyses, the Yakut were treated both as one single group, as well as individual groups, determined by collection site. Each Yakut sub-population was composed of similar sample sizes to that of the other three main populations of study.

#### **3.3.1 Haplogroup Prediction**

All samples were characterized by haplotypes based on HVS1 sequence data and again from D-loop data. Haplogroups were determined for each of the haplotypes present in all populations using the

mtDNA haplogroup prediction tool published by *Mitomap*

(<http://www.mitomap.org/MITOMASTER/WebHome>). This prediction tool compares sequences to the revised Cambridge Reference Sequence (rCRS) and predicts an mtDNA haplogroup assignment based on mutations present in the sample.

### 3.3.2 Within Population Variation

Nucleotide and gene diversity were calculated using mtDNA D-loop sequence data in Arlequin version 3.5.1.2 (Excoffier et al., 2005; Excoffier and Lischer, 2010) in order to determine the amount of variation within populations represented by both HVS1 alone and the D-loop sequence data (Nei, 1987; Nei and Li, 1979; Tajima, 1983). Both of these tests of intrapopulation variation are equivalent to estimating a population's average heterozygosity for haploid data. Gene diversity is, however, less likely to be affected by stochastic changes in allele frequency and recent evolutionary events such as demographic events or genetic drift than nucleotide diversity (Nicholson et al., 2002; Helgason et al., 2003). Gene diversity, or the probability that randomly chosen two haplotypes from the sample are different, is defined as:

$$\hat{H} = \frac{n}{n-1} \left( 1 - \sum_{i=1}^k p_i^2 \right) \quad (\text{Equation 3.1})$$

where  $n$  is the number of gene copies within the sample,  $k$  is the number of haplotypes, and  $p_i$  is the frequency of the  $i^{\text{th}}$  haplotype (Nei, 1987). Nucleotide diversity is defined as the probability that two randomly chosen homologous nucleotide sites are different. It is represented by the equation:

$$\hat{\pi}_n = \frac{(\sum_{i=1}^k \sum_{j < i} p_i p_j \hat{a}_{ij})}{L} \quad (\text{Equation 3.2})$$

where  $p_i$  is the probability of the  $i^{th}$  sequence in the population,  $p_j$  is the probability of the  $j^{th}$  sequence in the population,  $\hat{d}_{ij}$  is the number of nucleotide differences between the  $i^{th}$  and  $j^{th}$  sequence, and  $L$  is the number of loci present (Tajima, 1983; Nei, 1987).

### 3.3.3 Among Populations Variation

#### 3.3.3.1 AMOVA

Using Arlequin version 3.5.1.2 (Excoffier and Lischerl, 2010), Analysis of Molecular Variance (AMOVA) was used to characterize the amount of variation among the study populations based on the entire D-loop and HVS1 sequence data. Results were compared between the two data sets for each grouping of populations. AMOVA is an extension of the Analysis of Variance (ANOVA), where genetic variance is sequestered into several hierarchical levels, partitioning the total variance among groups into components of covariance. F-statistic analogs, or  $\Phi$ -statistics, are calculated and reflect the proportion of variation and haplotypic diversity attributable to various hierarchical sub-divisions of the populations of study. Kimura 2-parameter distances with a gamma correction of 0.26 were used with mtDNA sequence data (Kimura, 1980). Kimura 2P distance weighs transversions (T or C  $\leftrightarrow$  A or G) and transitions (T $\leftrightarrow$ C and A $\leftrightarrow$ G) differently in the calculation, as there are unequal substitution rates for transitions and transversions. This is highly relevant for Hypervariable segment data, which displaces a high ratio of transition:transversion (Ward et al., 1991). Kimura 2P distance is defined by the equation:

$$\delta_{jk}^2 = \frac{1}{2} \ln(1 - 2\hat{P} - \hat{Q}) - \frac{1}{4} \ln(1 - 2\hat{Q}), \quad (\text{Equation 3.3})$$

where  $\hat{P}$  is the frequency of transitions and  $\hat{Q}$  is the frequency of transversions between the sequence being compared. AMOVA compares the proportion of the Sum of Squared Differences (SSD) to the Mean Square Deviations (MSD) among group hierarchies. The total SSD is characterized by the following equation:

$$SSD_{(Total)} = \frac{1}{2N} \sum_{j=1}^N \sum_{k=1}^N \delta_{jk}^2, \quad (\text{Equation 3.4})$$

Where  $N$  equals the number of haplotypes and  $\delta_{jk}^2$  is the Euclidean distance between haplotypes  $j$  and  $k$ . This allows the haplotypes to be partitioned into covariance components: SSD among groups, SSD among populations within groups, and SSD within populations. Mean Squared Deviation (MSD) is then obtained by dividing the corresponding SSD by the appropriate degrees of freedom (Excoffier et al., 1992). The  $\Phi$ -statistics for AMOVA are tested for significance by bootstrapping the molecular data 1000 times. These F ratios are defined as follows:

$$\phi_{ST} = \frac{\sigma_a^2 + \sigma_b^2}{\sigma_T^2}, \quad \phi_{SC} = \frac{\sigma_b^2}{\sigma_b^2 + \sigma_c^2}, \quad \phi_{CT} = \frac{\sigma_a^2}{\sigma_T^2}, \quad (\text{Equation 3.5})$$

where  $\sigma_a^2$  is the covariance among groups,  $\sigma_b^2$  is the covariance among populations within groups,  $\sigma_c^2$  is the difference among individuals within a population, and  $\sigma_T^2$  is the total covariance of haplotypes (Excoffier and Lischer, 2010; Excoffier et al., 1992). These F ratios describe the proportions of total variation explained by that level of grouping.

AMOVA requires the researcher to define groups of populations, which creates a self-defined genetic structure tested by the analysis of variance. For the purpose of this study, the following AMOVAs were run for both HVS1 sequence data and D-loop sequence data: 1.) Language family: {Turkic = Yakut and Altai, and Tungusic = Evenki and Udehe}, 2.) Ethnicity: {Altai, Evenki, Udehe, and Yakut}, and 3.) Geography: { Southwest Siberia = Altai, East Siberia = Evenki and Yakut, Far East Siberia = Udehe}.

### 3.3.3.2 Multi-dimensional Scaling

Multi-dimensional scaling (MDS) is an ordination method often used in population studies to infer relationships between variables. This is done by displaying (geometrical) patterns of similarities and differences between multivariate data in dimensional space. This analysis reduces the variation

between groups as much as possible, so that the interpoint distances between groups correspond to the original, observed distance matrix (Kruskal, 1964). For this study, Kimura 2-parameter distance matrices were calculated using Arlequin version 3.5.1.2 (Excoffier and Lischerl, 2010) with a gamma correction value of 0.26 (Meyer et al., 1999). This value takes into account the differing substitution rates of transversions and transitions and allows for multiple substitutions per site. The goodness of fit test employed by this analysis is called *stress* and is represented by the following formula:

$$stress = \sqrt{\frac{\sum(d_{ij}^* - d_{ij}^f)^2}{\sum d_{ij}^{*2}}} \quad (\text{Equation 3.6})$$

where  $d_{ij}^*$  is the distances between the pairs of points  $ij$ . These are then compared to the original distances  $d_{ij}$ . A monotone function,  $d_{ij}^f$  (a regression line) is run between the two distance matrices, to assess the goodness of fit of the matrices. A high stress value suggests a poor fit and indicates that the chosen number of dimensions is not the most accurate way of representing the relationship among groups (Sturrock and Rocha, 2000). The closer the stress value gets to zero, the better the fit of the MDS plot to the original distance matrix. Stress values below 0.10 are considered excellent goodness of fit values, whereas stress values less than 0.20 are considered intermediate (Kruskal, 1964). MDS plots for this study were performed with NTSYSpc version 2.2 (Rohlf, 2008).

### 3.3.3.3 Neighbor-Joining

Neighbor-Joining trees (NJT) are used to explain relationships between haplotypes or populations by creating a topology of genetic distances and minimizing the total branch length of the dendrograms created. This provides a representation of the tree based on a matrix of genetic distances with the shortest evolutionary times between neighbors. Evolutionary relationships of the most likely

closest operational taxonomic unit (OTU), or “neighbor,” can then be inferred (Saito and Nei, 1987). For the purposes of this study, NJT were constructed using Kimura 2P distances, with a gamma correction of 0.26, for mtDNA HVS1 and D-loop distances. Kimura 2p distances were created by bootstrapping 1000 replicates to form the distance matrix. Bootstrapping is a technique used to estimate the distribution of a sample by resampling the data, creating, in this case, the most likely distance matrix. The Yakut were treated as nine distinct populations based on village of origin. Cophenetic distance matrices were generated from the NJTs and tested against the original FSTs distance matrix with Mantel randomizations to assess the goodness of fit of the tree to the original data. Cophenetic distances are constructed by a hierarchical clustering technique that produces values of similarity or dissimilarity between, in this case, populations (Rohlf and Sokal, 1981). The Cophenetic matrix and Mantel randomization were both performed in NTSYSpc version 2.5 (Rohlf, 2000; 2008).

A second set of NJTs were also constructed in MEGA ver 5.10 (Tamura et al., 2007; 2011) based on sequence data, using 1000 bootstrapped trees. Again, using bootstrapping involves resampling the data a number of times (1000), creating a new tree with each set of sampled data. The topology of the tree is compared to the original tree, and each time a branch is the same, it receives a value of 1, whereas differences in branches receives a value of 0. The percentage of the number of times a branch is given the number 1 is provided at the branch point of the final tree. This bootstrap value is generally accepted as correct if it is 95% or higher (Felsenstein, 1985; Nei and Kumar, 2000; Efron, 1982; Tamura et al., 2007, 2011). Comparing the bootstrap values of each tree will provide an opportunity to make inferences about the result of adding more than SNPs (the remainder of the D-loop).

### 3.3.4 Measures of Forces of Evolution

Various analyses can be used to test the effects of the forces of evolution on a population, which can lead to better inferences of how the (genetic) differences and variation have come within and between populations. For this study, the demographic indices: Tajima's D and Fu's  $F_s$ , and Mismatch Analysis have been applied to both HVS1 and the whole D-loop. These analyses were computed using Arlequin version 3.5.1.2 (Excoffier and Lischer, 2010).

#### 3.3.4.1 Neutrality Tests

Fu's  $F_s$  and Tajima's D, two neutrality test statistics, were applied to mtDNA sequence data in order to detect departures from the null evolutionary model of constant population size (fluctuations in population size), as well as infer possible effects of natural selection on the population. The neutrality test Tajima's D (Tajima, 1989), is useful for haploid data such as mtDNA sequence data because it is based on the infinite-sites model without recombination. Tajima's D uses mtDNA pairwise sequence data and operates under the assumption of a constant mutation rate to test the neutral model. It is represented by:

$$D = \frac{\hat{\theta}_\pi - \hat{\theta}_S}{\sqrt{\text{var}(\hat{\theta}_\pi - \hat{\theta}_S)}} \quad (\text{Equation 3.7})$$

where  $\hat{\theta}_\pi$  represents the average number of nucleotide differences present in a population and  $\hat{\theta}_S$  is the number of segregating sites. Tajima's D compares the proportion of average number of pairwise differences ( $\hat{\theta}_\pi$ ) to the total number of nucleotide differences ( $\hat{\theta}_S$ ). Populations that experience a large number of low-frequency mutations will have  $\hat{\theta}_S$  values that are larger than  $\hat{\theta}_\pi$ , resulting in negative  $D$  values. This is indicative of population expansion. Conversely, positive  $D$  values result from larger  $\hat{\theta}_\pi$ , which suggests genetic drift, such as a population bottleneck. This latter scenario would occur when a

population has a greater number of intermediate- and high-frequency mutations. Stable populations, under no evolutionary pressures, should be represented by  $D$  values close to zero (1989). It should be noted, however, that significant  $D$  values can be caused by other factors such as mutation rate heterogeneity and selection (Tajima, 1989; Aris-Brosou and Excoffier, 1996).

The test statistic Fu's  $F_s$  (Fu, 1997), is based on the infinite-sites model without recombination. Instead of using pairwise differences like Tajima's  $D$ , Fu's  $F_s$  utilizes haplotypic distributions. This neutrality statistic tests for the same evolutionary processes' effects on a population, but compared to Tajima's  $D$ , this test is more sensitive to population growth (expansion) than genetic drift (Zlojutro et al., 2006). Fu's  $F_s$  can be calculated as:

$$F_s = \ln \left( \frac{s'}{1-s'} \right) \quad \text{(Equation 3.8)}$$

and where

$$s' = \Pr \langle K \geq k_{obs} \mid \theta = \theta_\pi \rangle \quad \text{(Equation 3.9)}$$

In other words,  $F_s$  is equal to the natural log of the probability of a random neutral sample having more haplotypes ( $K$ ) present in the data than the observed haplotypes ( $k_{obs}$ ), given  $\hat{\theta}_\pi$ , divided by the probability of  $K$  not being present. Negative  $F_s$  values are caused by the existence of more haplotypes than expected in the population. Positive  $F_s$  values occur when there are fewer haplotypes than expected. Large negative values then are indicative of population expansion, whereas large positive values suggest genetic drift.  $F_s$  values are considered significant at  $p < 0.02$  (Fu, 1997).

### 3.3.4.2 Mismatch Analysis

In addition to neutrality tests, the effects of evolutionary forces on populations can be detected through certain characteristics of mismatch distributions of (mitochondrial) sequence data, or the

distribution of the observed number of pairwise differences between populations' mtDNA haplotypes (Rogers and Harpending, 1992). Such characteristics may be used to infer population expansion, bottlenecks, or stability. Under the infinite-sites model with no recombination, a population at drift-mutation equilibrium will display a unimodal distribution of mismatches with a peak at zero. If the distribution is unimodal, but with a peak at greater than zero, it suggests that the population has undergone significant expansion. A multimodal distribution with a peak at zero indicates that the current population is stable. If the distribution is multimodal but with a significant peak at some point greater than zero, it is indicative of the population's expansion in the remnant past. Additionally, multimodal distributions with peaks greater than zero, can be caused by population bottlenecks or selection. The raggedness of the distributions can also be used to make inferences about the evolutionary history of a population. The mismatch distribution of populations that are unadmixed will display smooth distribution, whereas recently admixed populations will display a more ragged form. For this study, mismatch distributions were constructed for both HVS1 and the whole D-loop for each population. Further, distributions will be given a raggedness index ranging from 0.0 to 1.0, where the closer to 1.0, the more ragged is the distribution. Values close to 0.0 are indicative of population expansion, where high raggedness values signify a stable population (Hudson and Slatkin, 1991; Rogers and Harpending, 1992; Jobling et al., 2004). Mismatch distributions were created from mtDNA sequences in Arlequin version 3.5.1.2 (Excoffier and Lischer, 2010) and then plotted with Microsoft Excel®.

### **3.3.5 Phylogeographic Methods**

Phylogeography is a field that uses genetic data and geographical coordination in order to make inferences about the relationship between the genes and the ecology of a species. For population genetics, this is useful because phylogeography uses changes in both gene frequencies and geographic

spaces of a species to infer evolutionary events influencing the species (Avise et al., 1987; Avise, 1998). Another utility of these methods for anthropological genetics is that the distance parameters of study are not limited to strictly genes and geography, but linguistic, temporal, and cultural distances can also be applied.

### 3.3.5.1 Mantel Randomization

The mantel randomization test is used to examine the significance of the correlation between two matrices,  $X = \{x_{ij}\}$  and  $Y = \{y_{ij}\}$ . The correlation value is represented by:

$$r_{xy} = \frac{SP(X,Y)}{\sqrt{SS(X)SS(Y)}}, \quad (\text{Equation 3.10})$$

Where  $SP$  is the sum of the products for  $X$  and  $Y$  and  $SS$  is the sum of squares for the individual matrices,  $X$  and  $Y$  (Smouse et al., 1986). For the present study, a mantel test was utilized to evaluate: 1.) the relationship between phenotypic variation of HVS1 data and phenotypic variation of the whole D-loop, and 2.) phenotypic variation of the HVS1 and D-loop sequence data (separately) for each population, to geography. In order to create the matrices needed for this test, standardized genetic (pairwise) distances for mtDNA data as well as geographic distances were employed. The randomization requires the holding of one matrix constant while creating a randomly configured matrix from the second, for which the correlation statistic  $r_{xy}$  is then calculated. This test uses 1000 randomizations to compute  $r_{xy}$ . The correlation from the original distance matrices is then compared to the randomized correlation in order to check the significance of the statistic. If the results are in fact significant, and there is no true relationship between the two matrices in comparison, then the correlation between the two ( $r_{xy}$ ) would be low, or similar to the correlation value among the randomized and that which was held as a constant. If the results are significant and there is, in fact a relationship between the two variables, genes and geography for example, then one would see a high correlation value ( $r_{xy}$ ) when compared to

the randomized matrix. For this study, pairwise distance matrices (FSTs) for sequence data were created with Arlequin version 3.5.1.2 (Excoffier and Lischer, 2010) and the geographic distances were created using Latitude/Longitude Distance Calculator from NOAA/National Weather Service. In addition to comparing genetic and geographic distance matrices, the HVS1 FST matrix was compared to that of the whole D-loop (HVS1, 2, and3) to test the significance of the correlation between the two matrices. Mantel Randomization tests were run with Mantel version 3.1 (Relethford, 2003).

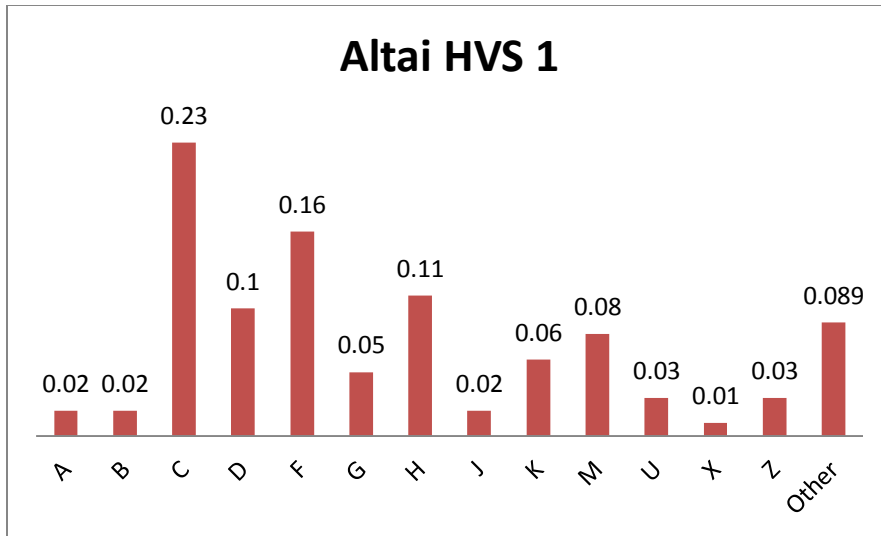
## CHAPTER 4: RESULTS

The analytical methods described in the previous chapter were applied to the four Siberian populations of study: the Altai, Evenki, Udehe, and Yakut. Since the Yakut sample size was extremely large, they were subdivided into their nine respective collection sites or villages, resulting in a total of 12 populations for the entire project. This chapter presents the results of these analyses based on sequencing results of the mtDNA Hypervariable segments 1, 2, and 3. These analyses include: haplotype ( $h$ ) and gene diversity ( $\pi$ ), neutrality test statistics (Tajima's  $D$  and Fu's  $F_s$ ), AMOVA, Mantel tests, neighbor-joining tree (NJT), mismatch analysis, and multidimensional scaling.

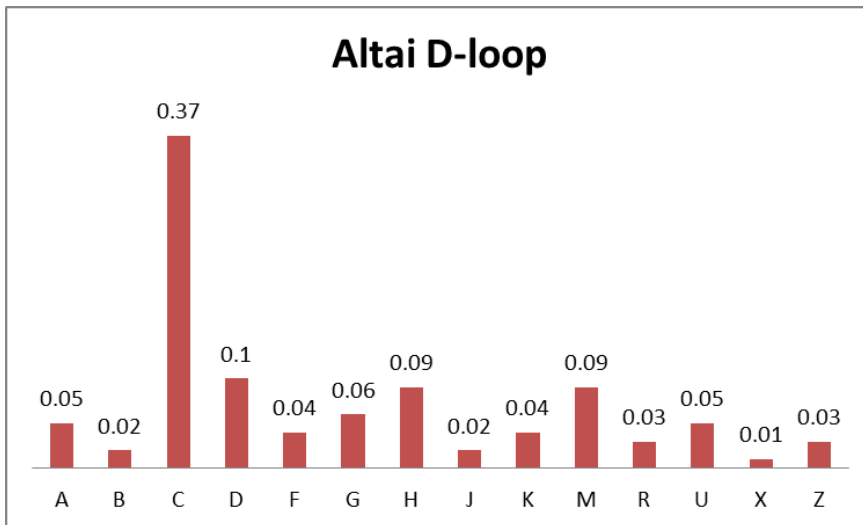
### 4.1 Haplogroup and Haplotype Results

Haplogroups were identified for samples using the Mitomaster haplogroup prediction tool published by *Mitomap*. (<http://www.mitomap.org/MITOMASTER/WebHome>). Samples were assigned to haplogroups based on sequence data from HVS 1 and compared to a second assignment using D-loop sequences to determine if haplogroup designation changed upon the addition of more SNPs. The findings are similar to previously reported haplogroups for these populations. All of the Native American haplogroups (A-D and X) were represented in the Altai, as well as various other Asian and European haplogroups. The highest reported haplogroups for the Altai HVS1 sequences were as follow: C – 23%, D – 10%, F – 16%, and H – 11 %. Using all three hypervariable segments the highest reported haplogroups were: C – 32 %, D – 10%, H – 9%, M – 9%, G – 6%). The Evenki also represented a mix of European, Siberian and East Asian haplogroups. The highest frequencies were of haplogroups C, D, H and A (28%, 17%, 13%, and 11% respectively) for HVS1 sequences, and of C, D, U, H, A (30%, 25%, 25%, 9%, and 8%) for all three HVS. The Udehe had the fewest number of haplogroups represented by the sample for both HVS1 and D-loop. Likewise, the frequencies of haplogroups did not change with the

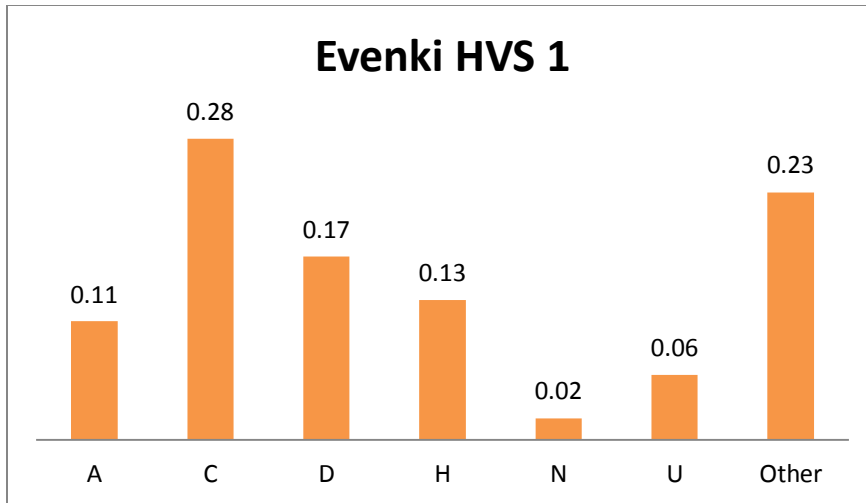
addition of HVS2 and 3 to the samples. In both cases, the reported haplogroups were as follows: C – 12%, M – 46%, N – 32%, and Y – 10%. In this sample, the only Siberian/Native American haplogroup present was C, at the second lowest frequency. The remaining haplogroups are of East Asian origin. The Yakut shared the highest number of haplogroups, representing all of the Native American/Siberian haplogroups (A-D) except for X. They also have various haplogroups of East Asian and Western Eurasian origin. The highest haplogroups for HVS1 data in the Yakut were: C – 46%, D – 24%, with many low frequency haplogroups. The D-loop sequence data is reported as strikingly similar with the haplogroups C and D sharing the highest frequencies (43% and 26%, respectively). See figure 4.1a-z for a breakdown of all haplogroups represented in the populations of study. It is interesting to note that in the cases in which HVS1 sequences left some undefined haplogroups, the addition of the second and third hypervariable segment resolved these undefined haplogroups, leaving none as “undefined.”



*Figure 4.1a displays mtDNA haplogroup predictions based on HVS 1 sequence data for the Altai*

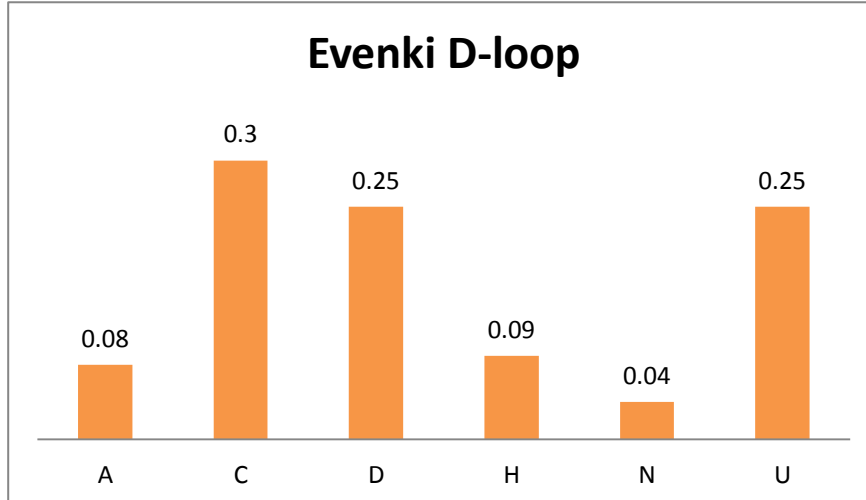


*Figure 4.1b displays mtDNA haplogroup predictions based on D-loop Sequence data for the Altai*



c.

Figure 4.1c displays mtDNA haplogroup predictions based on HVS 1 sequence data for the Evenki



d.

Figure 4.1d displays mtDNA haplogroup predictions based on D-loop sequence data for the Evenki

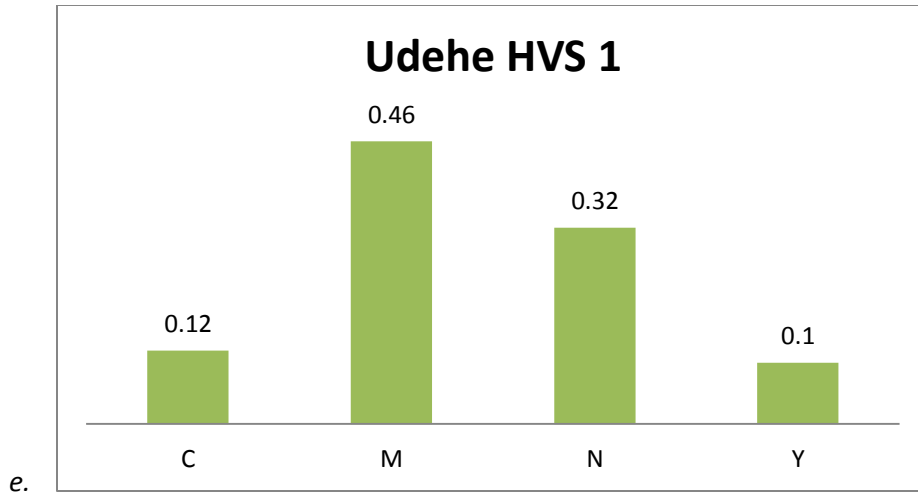


Figure 4.1e displays mtDNA haplogroup predictions based on HVS 1 sequence data for the Udehe

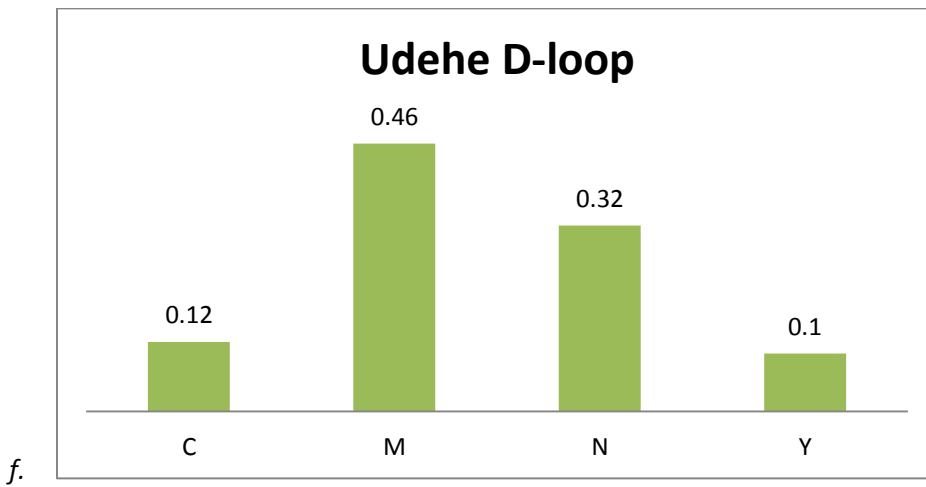
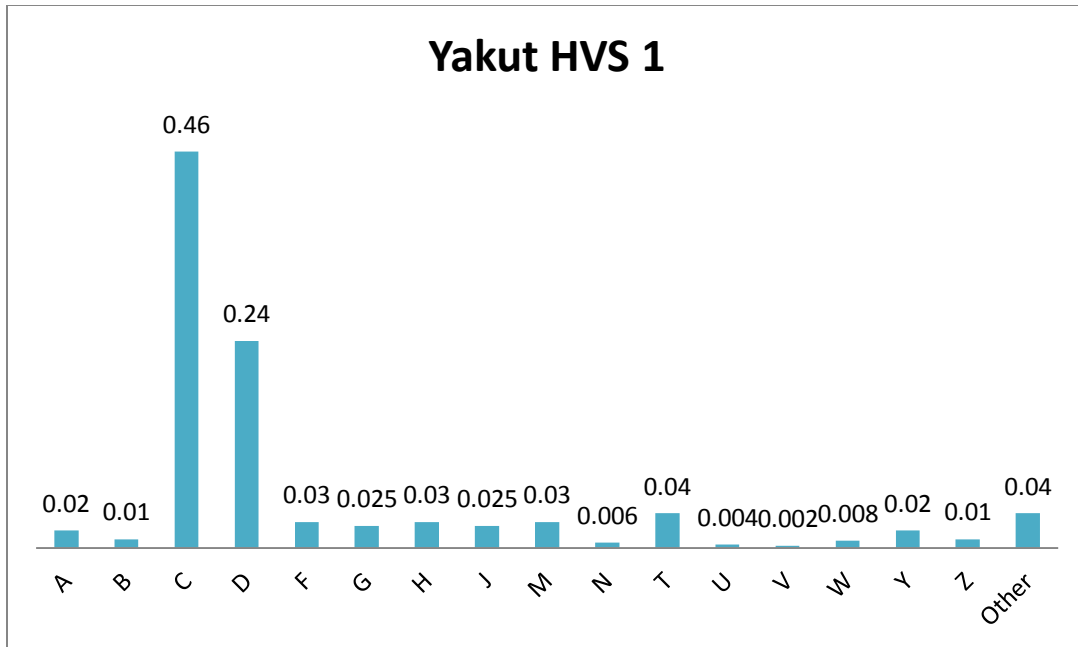
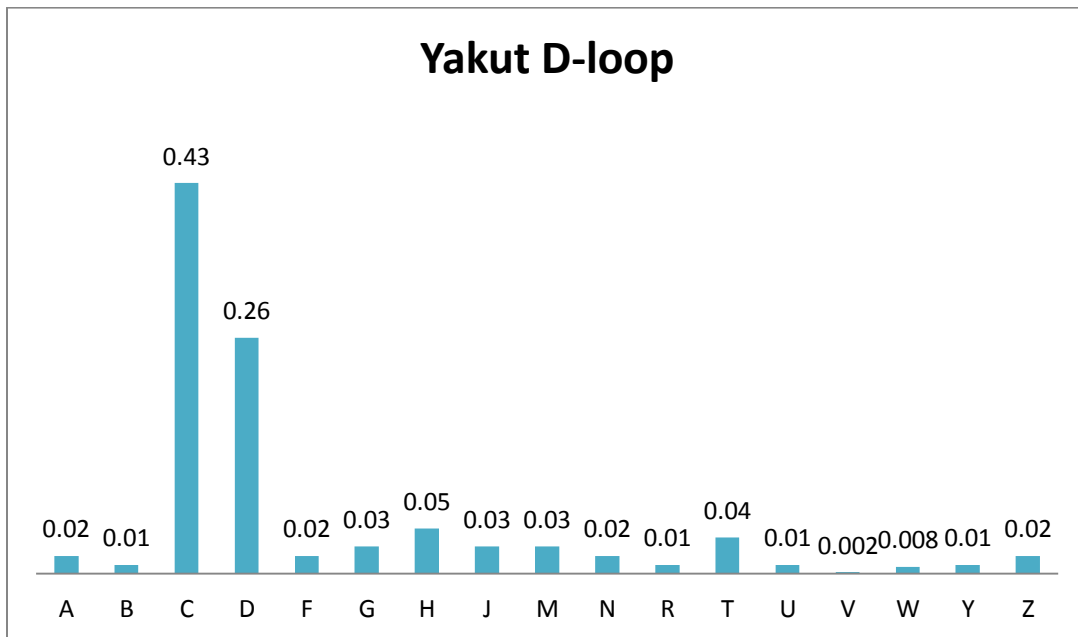


Figure 4.1f displays mtDNA haplogroup predictions based on D-loop sequence data for the Udehe



g.

Figure 4.1g displays mtDNA haplogroup predictions based on HVS 1 sequence data for the Yakut



h.

Figure 4.1h displays mtDNA haplogroup predictions based on D-loop sequence data for the Yakut

## 4.2 Hypervariable Segment Sequencing

The sequencing results of the mtDNA D-loop (nucleotide positions 1-600 for HVS 2 and 3, and 16,000-16,600 for HVS1) for all four populations are listed in table 4.2.1. Variable sites are shown relative to the Cambridge Reference Sequence (CRS), where an asterisk (\*) represents a deletion. For the Altai, Evenki, and Udehe, there were 44, 28, and 9 haplotypes, respectively. For the various Yakut communities (YakutCK, YakutDA, YakutDE, YakutDY, YakutYE, YakutYJ, YakutYK, YakutYM, and YakutYO), there were 27, 35, 30, 25, 21, 27, 50, 24, and 30 haplotypes respectively present. Haplogroups are listed along with each haplotype in the table. Molecular diversity indices for mtDNA sequences of the populations in this study were calculated with Arlequin ver. 3.5.1.2, (Excoffier et al., 2005; Excoffier and Lischer, 2010) and are presented in tables 4.2.2a and b.

	Altai	Evenki	Udehe	Yakut
Sample size (n)	101	53	41	471
Nucleotide Range	1-600; 16,000-16,600	1-600; 16,000-16,600	1-600; 16,000-16,600	1-600; 16,000-16,600
No. of haplotypes	44	28	9	269
No. of transitions	78	48	23	119
No. of transversions	3	6	2	17
No. of indels	13	12	4	13
No of polymorphic sites	94	66	28	143
Nucleotide composition	C: 33.07% T: 24.00% A: 35.27% G: 12.66%	C: 33.21% T: 23.92% A: 30.24% G: 12.64%	C:33.04% T: 24.01% A: 30.24% G: 12.71%	C: 33.08% T: 24.01% A: 30.22% G: 12.70

*Table 4.2.2.a Displays molecular diversity values of the four main populations used in this study based on D-loop sequence data.*

	YakutCK	YakutDA	YakutDE	YakutDY	YakutYE	YakutYJ	YakutYK	YakutYM	YakutYO
Sample size (n)	34	66	61	34	26	48	103	36	63
Nucleotide Range	1-600; 16,000-16,600	1-600; 16,000-16,600	1-600; 16,000-16,600	1-600; 16,000-16,600	1-600; 16,000-16,600	1-600; 16,000-16,600	1-600; 16,000-16,600	1-600; 16,000-16,600	1-600; 16,000-16,600
No. of haplotypes	27	35	30	25	21	27	50	24	30
No. of transitions	33	72	60	53	36	42	65	55	47
No. of transversions	3	3	3	4	2	3	8	2	11
No. of indels	6	9	9	9	10	9	9	7	6
No of polymorphic sites	41	83	72	66	48	54	80	63	64
Nucleotide composition	C: 33.11% T: 23.99% A: 30.20% G: 12.69%	C: 33.05% T: 24.03% A: 30.19% G: 12.73%	C: 33.10% T: 23.98% A: 30.24% G: 12.68%	C: 33.08% T: 23.99% A: 30.18% G: 12.74%	C: 33.09% T: 24.01% A: 30.19% G: 12.70%	C: 33.05% T: 24.02% A: 30.29% G: 12.64%	C: 33.08% T: 24.01% A: 30.19% G: 12.73%	C: 33.07% T: 24.01% A: 30.23% G: 12.90%	

Table 4.2.2.b Displays molecular diversity values of the various Yakut villages used in this study, based on D-loop sequence data.

Shared haplotypes among the populations of study are presented in Table 4.2.3. The greatest numbers of shared haplotypes are found between YakutDA and YakutYK, YakutYK and YakutYO, and YakutYM and YakutYO with 11 shared haplotypes. The Udehe shared the fewest number of haplotypes among all of the populations. Of the four haplotypes the Udehe shared with other populations, one was with YakutDA, YakutYK, YakutYM, and YakutYO each.

	Altai	Evenki	Udehe	YakutCK	YakutDA	YakutDE	YakutDY	YakutYE	YakutYJ	YakutYK	YakutYM
Evenki	2										
Udehe	0	1									
YakutCK	2	3	0								
YakutDA	1	4	1	6							
YakutDE	2	2	0	9	8						
YakutDY	1	2	0	5	5	4					
YakutYE	0	2	0	6	6	6	5				
YakutYJ	1	2	0	6	5	10	0	4			
YakutYK	2	3	1	9	11	7	9	9	4		
YakutYM	1	3	1	4	9	5	4	6	4	8	
YakutYO	3	3	1	10	10	9	10	9	6	11	11

Table 4.2.3 displays the shared Haplotypes between the populations used in this study, based on D-loop sequence data.

### 4.3 Within-Population Variation

#### 4.3.1 Genetic Diversity

Gene and nucleotide diversity ( $h$  and  $\pi$ ), were calculated for all the populations of study for both HVS1 and D-loop and are presented in Table 4.3.1.a and 4.3.1.b. The Udehe have the lowest gene diversity at 0.8256 (0.0347) for both hypervariable segments examined. They also have the lowest nucleotide diversity at 0.012829 (0.007050) for HVS1 and 0.007276 (0.003863) for the D-loop. For HVS1 sequences, the highest gene diversity is found among the Yakut (combined samples) at 0.9944 (0.0006) with YakutYE displaying the highest among the Yakut villages at 0.9754 (0.0152). The highest nucleotide

diversity for HVS 1 was also found among the combined Yakut samples at 0.016815 (0.008785), with the YakutYK displaying the highest diversity among Yakut villages at 0.018519 (0.009668). For D-loop sequences, gene diversity was highest among the Combined Yakut at 0.9960 (0.0005), with the YakutYE village measuring highest among all Yakut at 0.9846 (0.0144). The highest nucleotide diversity overall for D-loop was found in the YakutDA village at 0.12945 (0.006545), but when comparing the Yakut as one combined population to the other three (Altai, Evenki, and Udehe), the Evenki showed the highest nucleotide diversity at 0.011929 (0.006082).

Population	Gene diversity (S.D.)	Nucleotide Diversity (S.D.)
Altai	0.9632 (0.0082)	0.015210 (0.008084)
Evenki	0.9528 (0.0132)	0.015807 (0.008450)
Udehe	0.8256 (0.0347)	0.012829 (0.007050)
Yakut(combined)	0.9944 (0.0006)	0.016815 (0.008785)
YakutCK	0.9608 (0.0183)	0.015045 (0.008172)
YakutDA	0.9608 (0.0098)	0.017825 (0.009388)
YakutDE	0.9410 (0.0133)	0.016052 (0.008546)
YakutDY	0.9501 (0.0201)	0.015704 (0.008494)
YakutYE	0.9754 (0.0152)	0.016446 (0.008944)
YakutYJ	0.9371 (0.0144)	0.014253 (0.007714)
YakutYK	0.9659 (0.0069)	0.018519 (0.009668)
YakutYM	0.9651 (0.0144)	0.016437 (0.008837)
YakutYO	0.9524 (0.0094)	0.016414 (0.008716)

*Table 4.3.1.a. displays Gene and Nucleotide Diversity values with Standard Deviations (S.D.) for mtDNA HVS1 sequence data.*

Population	Gene diversity (S.D.)	Nucleotide Diversity (S.D.)
Altai	0.9669 (0.0080)	0.010522 (0.005352)
Evenki	0.9717 (0.0080)	0.011929 (0.006082)
Udehe	0.8256 (0.0347)	0.007276 (0.003863)
Yakut(combined)	0.9960 (0.0005)	0.011137 (0.005602)
YakutCK	0.9804 (0.0138)	0.009449 (0.004943)
YakutDA	0.9692 (0.0086)	0.012945 (0.006545)
YakutDE	0.9656 (0.0094)	0.010659 (0.005456)
YakutDY	0.9786 (0.0132)	0.010514 (0.005461)
YakutYE	0.9846 (0.0144)	0.009955 (0.005239)
YakutYJ	0.9619 (0.0129)	0.009679 (0.005008)
YakutYK	0.9688 (0.0072)	0.011859 (0.005991)
YakutYM	0.9762 (0.0114)	0.011141 (0.005757)
YakutYO	0.9672 (0.0079)	0.010277 (0.005270)

Table 4.3.1.b. displays Gene and Nucleotide Diversity values with Standard Deviations (S.D.) for mtDNA D-loop sequence data.

#### 4.4 Among-Population Variation

##### 4.4.1 AMOVA

Analysis of Molecular Variance (AMOVA) was performed in Arlequin ver. 3.5.1.2 (Excoffier et al., 2005; Excoffier and Lischer, 2010) in order to test whether there was population substructuring of the Siberian samples. Three models were tested by again using HVS1 as well as D-loop data. The first model that was tested grouped the twelve populations of study into their four respective cultural groups (Altai, Evenki, Udehe, and Yakut). Secondly, populations were grouped by language family, where the Evenki and Udehe were grouped as Tungusic speakers and the Altai and Yakut were grouped as Turkic speakers. The third model tested was geography, where the Altai were considered Southwest Siberia, the Evenki and Yakut were considered North-Central Siberia and the Udehe were considered Far East Siberia. In all cases, almost all of the variation was accounted for within communities (94.64-96.75%).

The results of the first AMOVA, grouping the populations by their respective cultures, revealed very little variation among the populations within-groups or among-groups, with a significant amount (95.11%) of the variation falling within the populations for HVS1 and 94.64% for the D-loop. Low fixation indices further illustrates this trend, where  $\Phi_{CT}$  = Among groups of populations,  $\Phi_{SC}$  = Among populations within groups, and  $\Phi_{ST}$  = within populations. Results of the first AMOVA are summarized in *Tables 4.4.1.a and 4.4.1.b.*

Source of Variation	d.f.	SSD	Variance	% of Variance	Fixation Indices	P-values
Among Groups	3	56.968	0.12973	3.82	FCT = 0.03821	<.01
Among Populations within Groups	8	40.711	0.03638	1.07	FSC = 0.01114	<.01
Within Populations	654	2111.884	3.22918	95.11	FST = 0.04893	<.00001
<b>Total</b>	665	2209.563	3.3953			

*Table 4.4.1.a displays results of AMOVA based on HVS1 sequence data, in which populations were grouped culture. (Significant variance values:  $P < .05$ )*

Source of Variation	d.f.	SSD	Variance	% of Variance	Fixation Indices	P-value
Among Groups	3	96.711	0.19103	3.38	FCT = 0.03383	<.05
Among Populations within Groups	8	88.32	0.11144	1.97	FSC = 0.02043	<.00001
Within Populations	654	3494.916	5.34391	94.64	FST = 0.05357	<.00002
<b>Total</b>	665	3679.946	5.64637			

*Table 4.4.1.b displays results of AMOVA based on D-loop sequence data, in which populations were grouped culture. (Significant variance values:  $P < .05$ )*

The second AMOVA grouped populations into two linguistic groups, Tungusic (Evenki and Udehe) and Turkic (Altai and Yakut). The results, presented in *Tables 4.4.1.c and 4.4.1.d*, also revealed that little of the variation was explained among groups or among populations within groups, with 95.01% of the variation found within individual populations for HVS1 and 94.47% for HVS1,2,and 3. Again, this is further supported by low Fixation indices.

Source of Variation	d.f.	SSD	Variance	% of Variance	Fixation Indices	P-values
Among Groups	1	21.766	0.08984	2.64	FCT = 0.02644	<.05
Among Populations within Groups	10	75.913	0.07942	2.34	FSC = 0.0240	<.00001
Within Populations	654	2111.884	3.22918	95.02	FST = 0.0498	<.00001
<b>Total</b>	665	2209.563	3.39844			

*Table 4.4.1.c displays results of AMOVA based on HVS1 sequence data, in which Populations were grouped by linguistic family. (Significant variance values:  $p < .05$ )*

Source of Variation	d.f.	SSD	Variance	% of Variance	Fixation Indices	P-value
Among Groups	1	36.688	0.13981	2.47	FCT = 0.02472	<.05
Among Populations within Groups	10	148.343	0.17278	3.05	FSC = 0.03132	<.00001
Within Populations	654	3494.916	5.34391	94.47	FST = 0.05526	<.00001
<b>Total</b>	665	3679.946	5.65649			

*Table 4.4.1.d displays results of AMOVA based on D-loop sequence data, in which Populations were grouped by linguistic family. (Significant variance values:  $p < .05$ )*

Results of the third AMOVA are presented in *Tables 4.4.1.e and 4.4.1.f*, where the populations were grouped by geography. In this case, the Udehe were considered Far-east Siberian, the Evenki and Yakut were considered North-central Siberian, and the Altai were considered Southwestern Siberian. As with the last two models, the AMOVA revealed very little variation explained among groups or among populations within groups, apportioning the vast majority of variation within populations. In fact, 94.9%

and 94.7% of the variation was within the populations of study respectively for HVS1 and D-loop, although in both cases, the Among group variation was not significant. Likewise, low fixation indices support these findings on all accounts, as is seen with the first two models.

<b>Source of Variation</b>	<b>d.f.</b>	<b>SSD</b>	<b>Variance</b>	<b>% of Variance</b>	<b>Fixation Indices</b>	<b>P-values</b>
Among Groups	2	41.669	0.1154	3.39	FCT = 0.03391	NS
Among Populations within Groups	9	56.011	0.05822	1.71	FSC = 0.01771	<.00001
Within Populations	654	2111.884	3.22918	94.9	FST = 0.05102	<.00001
<b>Total</b>	<b>665</b>	<b>2209.563</b>	<b>3.40281</b>			

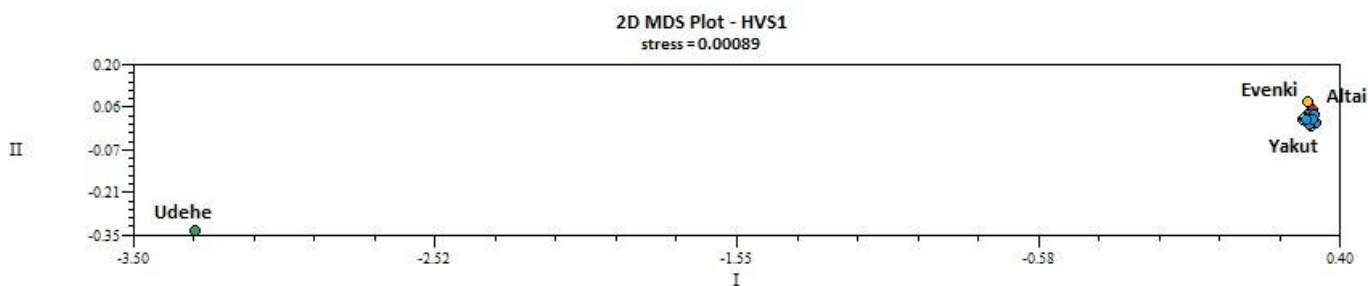
*Table 4.4.1.e displays results of AMOVA based on HVS1 sequence data, in which populations were grouped by geography. (Significant variance values:  $p < .05$ )*

<b>Source of Variation</b>	<b>d.f.</b>	<b>SSD</b>	<b>Variance</b>	<b>% of Variance</b>	<b>Fixation Indices</b>	<b>P-value</b>
Among Groups	2	66.174	0.1463	2.59	FCT = 0.02593	NS
Among Populations within Groups	9	118.856	0.15289	2.71	FSC = 0.02781	<.00001
Within Populations	654	3494.916	5.34391	94.7	FST = 0.05302	<.00001
<b>Total</b>	<b>665</b>	<b>3679.946</b>	<b>5.6431</b>			

*Table 4.4.1.f displays results of AMOVA based on D-loop sequence data, in which populations were grouped by geography. (Significant variance values:  $p < .05$ )*

#### 4.4.2 Multi-dimensional Scaling

Multidimensional scaling plots (MDS) were constructed in NTSYSpc ver 2.2 (Rohlf, 2008) for both HVS1 data as well as D-loop data based on Kumura 2P distance matrices (pairwise distances) that were calculated in Arlequin ver. 3.5.1.2 (Excoffier et al., 2005; Excoffier and Lischer, 2010) and are presented in *figures 4.4.2a-c*. The stress value for the 2D MDS plot for HVS1 data is 0.00089, which is well below the upper bound of 0.199 (Kruskal, 1964). This suggests that the plot is an excellent goodness of fit between the original distance matrix and the population distances displayed in the MDS plot (Sturrock and Rocha, 2000). In the plot, 11 of the 12 populations cluster tightly in the upper right corner, while the Udehe stand alone in the lower left. This would indicate that the Udehe are markedly distinct from the remainder of the populations of this study.



*Figure 4.4.2a displays the 2-dimensional MDS plot for HVS 1 sequence data for the populations used in this study*

To get a better resolution of the relationships between the Altai, Evenki, and various Yakut populations, the Udehe were then removed from the HVS 1 distance matrix and a new MDS plot was generated (*Figure 4.4.2b*). The stress value for this plot is 0.11304, which is an intermediate level, but still considered an acceptable goodness of fit with the original distance matrix. With the Udehe distances no longer being considered, the relationship between the remainder of the Siberian

populations is much clearer. The various Yakut populations all fall on the left half of the plot, whereas the Altai are separated along the top of the Y-axis (2<sup>nd</sup> Dimension), and in the right half of the X-axis (1<sup>st</sup> Dimension). Likewise, the Evenki appear to be more dissimilar from the Yakut populations and Altai, falling along the extreme boundary of the X-axis and within the lower half of the Y-axis.

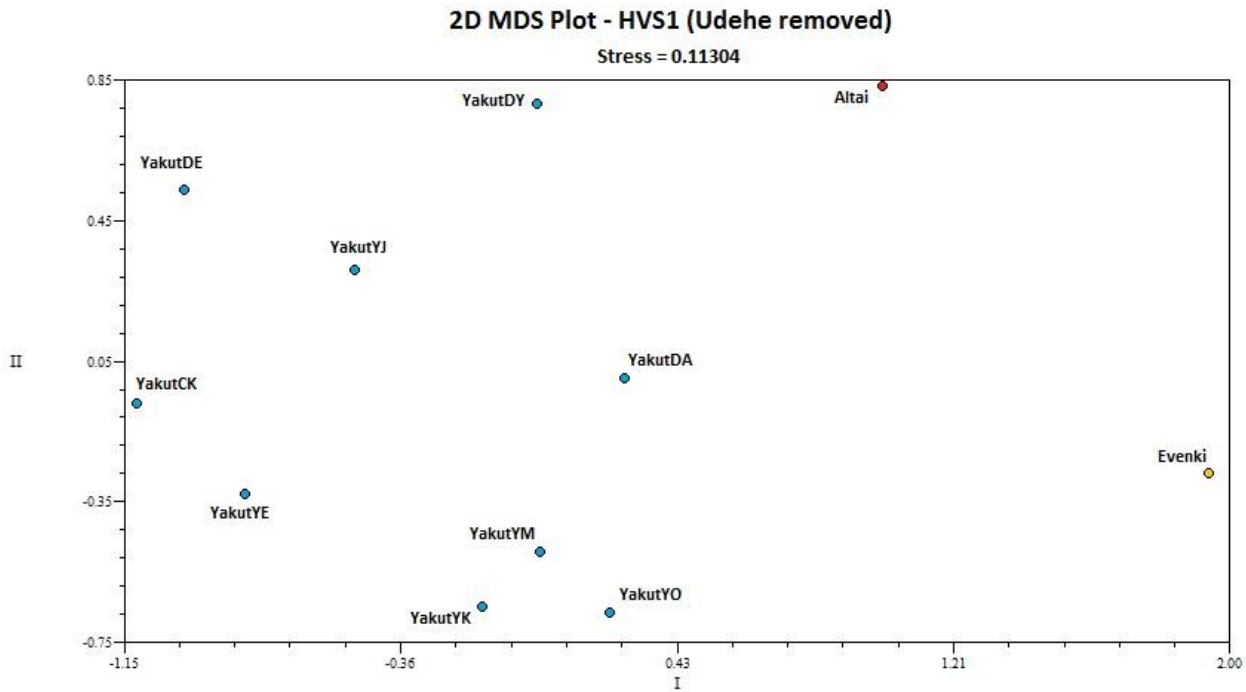


Figure 4.4.2b displays the 2-dimensional MDS plot for HVS1 sequence data for the populations used in this study. The Udehe have been removed.

The 2D MDS plot for D-loop distance matrix reveals 8 of the 9 Yakut populations clustering together toward the bottom and center of the plot, with the Altai close by. The Evenki are still removed, toward the top of the Y-axis, and the center of the X-axis, suggesting more dissimilarity with the other populations. Though the Udehe are still dissimilar from everyone else, it appears that the addition of HVS2 and 3 loci have reduced their dissimilarity significantly. The stress for this plot is 0.07996, which is

well below the upper bounds of an excellent goodness of fit with the original distance matrix (Kruskal 1964).

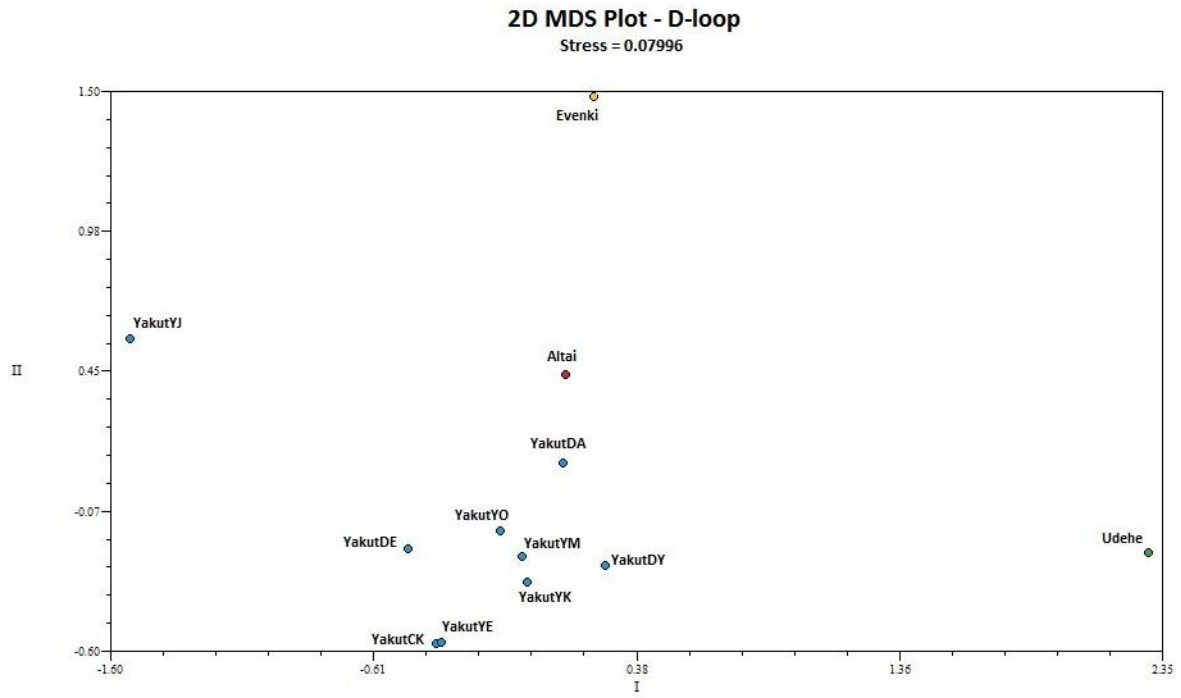


Figure 4.4.2c displays the 2-dimensional MDS plot from D-loop sequence data for the populations used in this study.

### 4.4.3 Neighbor-Joining Trees

Neighbor-Joining Trees (NJ) were constructed using Kimura 2p distances for both HVS1 and D-loop data and are shown in *Figures 4.4.3a and 4.4.3b*. Notice that all of the Yakut, in both NJTs cluster on one branch. A cophenetic distance matrix was generated from the NJTs and compared to the original distance matrix with a Mantel randomization test. The Mantel test revealed an extremely high correlation between the HVS1 NJT and the original distance matrix that was also statistically significant ( $r=0.92213$ ,  $p<0.001$ ). The NJT based on the whole D-loop revealed a moderately high correlation that was also statistically significant ( $r=0.7533$ ,  $p<0.001$ ). Therefore, the NJTs are accepted as an accurate representation of the relationship among populations.

A second pair of Neighbor-Joining trees was constructed based on sequence data for all of the four populations, considering each individual sequence as a taxon, or operational taxonomic unit (OTU). This set of NJTs was constructed with 1000 bootstrap replicates. The resulting trees were incredibly large, as the combined sample size is well over 500 individuals, and were therefore not depicted in this study. Only bootstrap values of  $\geq 50\%$  are reported. A simple tally of the branches with bootstrap values revealed interesting differences between the two trees. The HVS1 NJT contained seventy-one branches with bootstrap values at  $\geq 50\%$ , twelve branches at  $\geq 80\%$ , and only two branches at  $\geq 95\%$ . The D-loop NJT showed an increased amount of highly probable bootstrap values for all instances, with ninety-five branches at  $\geq 50\%$ , thirty-seven branches at  $\geq 82\%$ , and eight branches at  $\geq 95\%$ .

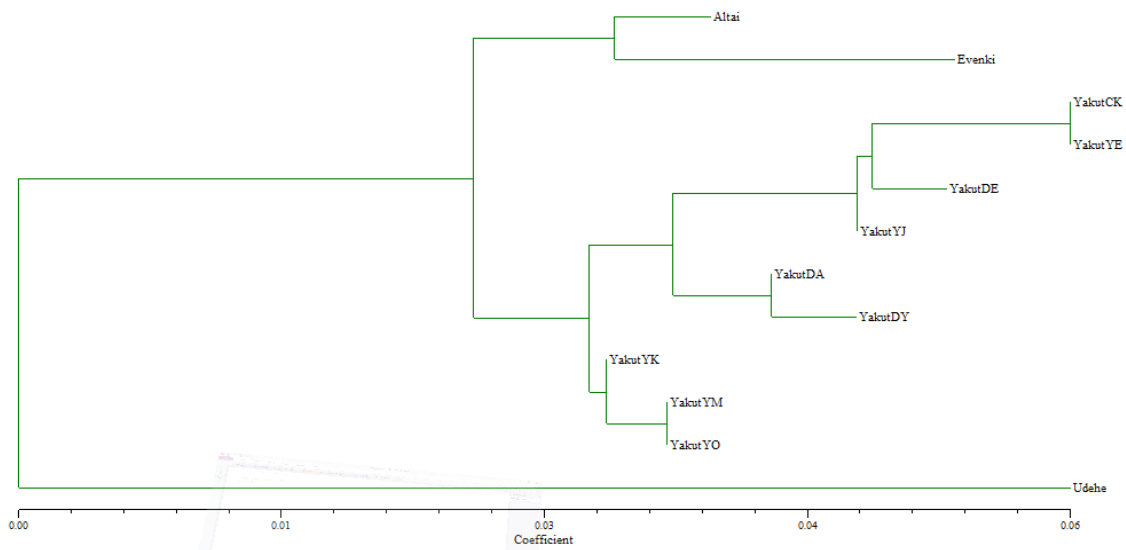


Figure 4.4.3a displays the Neighbor-joining tree based on HVS 1 sequence data of all the populations examined in this study.

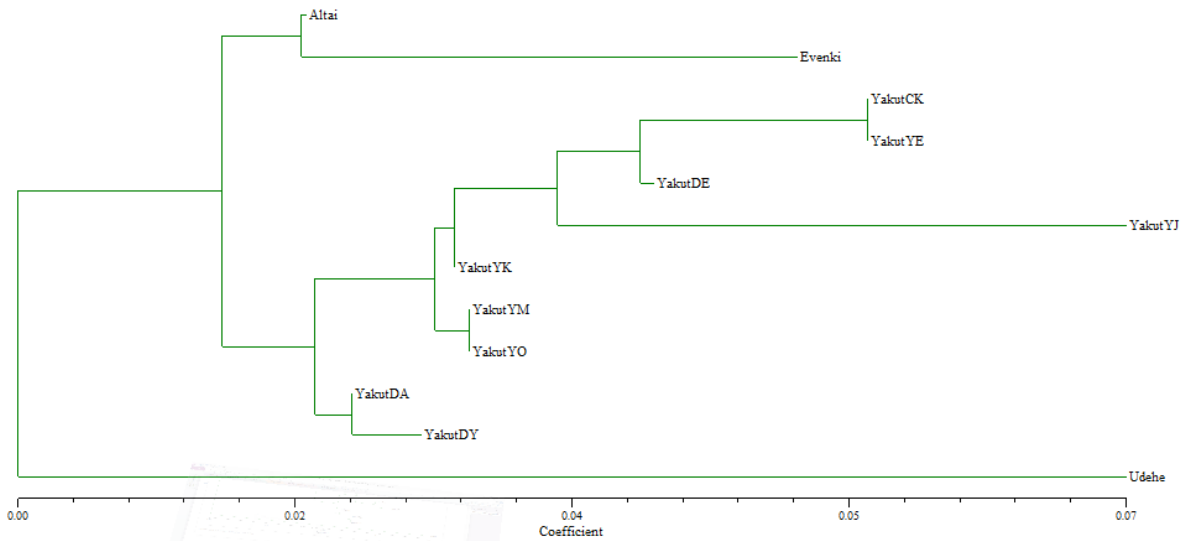


Figure 4.4.3b displays the Neighbor-Joining tree based on D-loop sequence data of all the populations examined in this study.

## 4.5 Forces of Evolution

### 4.5.1 Neutrality Test Statistics

Neutrality tests such as Tajima's  $D$  and Fu's  $F_s$  are used to quantify the probability that a population has undergone demographic events such as genetic drift, or expansion. Results of these tests for the populations of study are presented in *Table 4.5.1a* and *4.5.1b*. For both HVS1 and D-loop sequence data, all populations have negative values for both  $D$  and  $F_s$ , except for the Udehe. Though the Udehe have a small, negative  $D$  value for HVS 1 sequences, the value for  $F_s$  is positive. Likewise, both statistics for the D-loop data of the Udehe stand alone as positive. The combined Yakut population has the highest  $D$  value among all of the populations for both HVS1 and D-loop. In HVS1, it is the only  $D$  value that is significant, whereas for the D-loop, the Altai have a  $D$  value that is similar in magnitude and also is significant. In both cases, these two populations'  $D$  values are considered to be largely negative, compared to the rest of the populations sampled. For both HVS1 and D-loop data, all populations expressed largely negative  $F_s$  values except the Udehe, which was positive. The Altai, combined Yakut, YakutCK, YakutDA, and YakutYK resulted in the largest significant  $F_s$  values for both HVS1 as well as the D-loop.

<b>Population</b>	<b>Tajima's D Statistic</b>	<b>Fu's Fs Statistic</b>
Altai	-1.3912	-18.58077**
Evenki	-0.80126	-5.3752
Udehe	-0.05764	2.18427
Yakut (combined)	-1.50105*	-24.33068**
YakutCK	-0.48471	-9.19640**
YakutDA	-1.16028	-9.08029*
YakutDE	-1.10208	-5.0962
YakutDY	-1.33287	-6.13438*
YakutYE	-0.77522	-6.08865*
YakutYJ	-0.84141	-3.41203
YakutYK	-0.94891	-16.63439**
YakutYM	-1.27731	-7.58133**
YakutYO	-1.20745	-4.67467

Table 4.5.1a Displays Tajima's D and Fu's Fs Neutrality Test Statistics for HVS 1 sequence data for Altai, Evenki, Udehe, and various Yakut populations.

\*  $P < 0.05$ ; \*\* $p < 0.01$

<b>Population</b>	<b>Tajima's D Statistic</b>	<b>Fu's Fs Statistic</b>
Altai	-1.47143*	-11.87125**
Evenki	-0.73136	-3.98721
Udehe	0.10322	4.18028
Yakut (combined)	-1.57814*	-23.75223**
YakutCK	-0.51087	-12.15991**
YakutDA	-1.19001	-6.36063
YakutDE	-1.23662	-5.17932
YakutDY	-1.32018	-7.77624**
YakutYE	-0.90287	-7.18593**
YakutYJ	-0.81958	-6.03300*
YakutYK	-0.98124	-15.73679**
YakutYM	-1.15021	-5.22724*
YakutYO	-1.0828	-5.24439

Table 4.5.1b Displays Tajima's D and Fu's Fs Neutrality Test Statistics for D-loop sequence data for Altai, Evenki, Udehe, and various Yakut populations.

\*  $P < 0.05$ ; \*\* $p < 0.01$

#### 4.5.2 Mismatch Analysis

Mismatch distributions of pairwise differences, utilizing mtDNA sequences, were created for all of the populations of study considering first HVS1 data and again using the D-loop. Multimodal distributions tend to be indicative of populations at demographic equilibrium, and unimodal distributions reflect population expansion. Recently admixed populations tend to display ragged distributions, whereas un-admixed populations are likely to be more smooth (Rogers and Harpending, 1992; Hudson and Slatkin, 1991). This analysis was conducted for the Yakut as a singular population as well as for each of the distinct villages. These distributions are depicted in *Figures 4.5.2a-4.5.2l*, and raggedness indices are presented in *Table 4.5.2a*. *Table 4.5.2b* illustrates the Mean Pairwise Differences, and 95% Confidence Intervals for each population. The HVS1 distributions for the Altai, Evenki, Udehe and Yakut (combined) are all relatively unimodal with peaks at 5-6 pairwise differences and a minor mode at zero difference. The Udehe reveal a raggedness index of approximately 0.421, which is the highest of all of the populations. The distributions for these populations have also been depicted in a combined distribution for ease of comparison in *Figure 4.5.2b*.

Population	HVS1 Raggedness Index	D-Loop Raggedness Index
Altai	0.01295842	0.00821072
Evenki	0.01732544	0.01166949
Udehe	0.42075699	0.13977543
Yakut (Combined)	0.00566828	0.00247487
YakutCK	0.01182317	0.01247772
YakutDA	0.00831510	0.00947615
YakutDE	0.02227418	0.00800054
YakutDY	0.02520328	0.01460341
YakutYE	0.02040237	0.01123787
YakutYJ	0.03036819	0.02395818
YakutYK	0.00943150	0.00608382
YakutYM	0.02513983	0.00688587
YakutYO	0.03726828	0.01002934

*Table 4.5.2a Displays Raggedness Indices of Mismatch Distributions for HVS 1 and D-loop sequence data for Altai, Evenki, Udehe, and the various Yakut populations.*

Population	HVS1		HVS1		D-Loop		D-Loop	
	Mean Pairwise Difference (SD)	Confidence Interval	Mean Pairwise Difference (SD)	Confidence Interval	Mean Pairwise Difference (SD)	Confidence Interval	Mean Pairwise Difference (SD)	Confidence Interval
Altai	6.099406 (2.926868)	4.812-8.015	10.616436 (4.875858)	7.337-14.191				
Evenki	6.322932 (3.046393)	3.619-10.523	11.964441 (5.497946)	7.383-16.800				
Udehe	5.131707 (2.538708)	3.652-8.834	7.268293 (3.474271)	5.727-11.254				
Yakut (Combined)	6.725934 (3.177193)	4.338-11.699	11.237638 (5.110423)	7.531-19.081				
YakutCK	6.014825 (2.939783)	3.018-16.316	9.477718 (4.458530)	5.216-25.752				
YakutDA	7.130070 (3.387164)	4.475-10.428	13.009324 (5.932579)	9.313-17.456				
YakutDE	6.420765 (3.082557)	3.995-10.442	10.712568 (4.944199)	6.647-17.983				
YakutDY	6.281640 (3.055802)	3.487-13.939	10.577540 (4.940514)	6.087-16.795				
YakutYE	6.578462 (3.211413)	3.745-11.849	10.024615 (4.735524)	5.889-20.194				
YakutYJ	5.701241 (2.780031)	3.175-8.060	9.727837 (4.534591)	6.348-13.499				
YakutYK	7.407577 (3.491561)	5.354-10.281	11.906910 (5.430867)	8.790-16.482				
YakutYM	6.574603 (3.180104)	4.006-10.589	11.174603 (5.194719)	7.516-18.024				
YakutYO	6.565796 (3.144210)	3.244-11.842	10.308244 (4.766930)	6.137-20.037				

Table 4.5.2b Displays Mean Pairwise Differences and associated 95% Confidence Intervals for HVS 1 and D-loop sequence data for the Altai, Evenki, Udehe, and various Yakut Populations

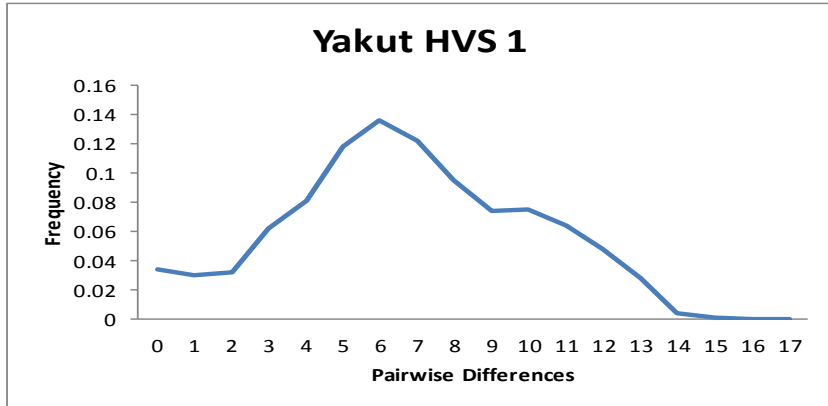
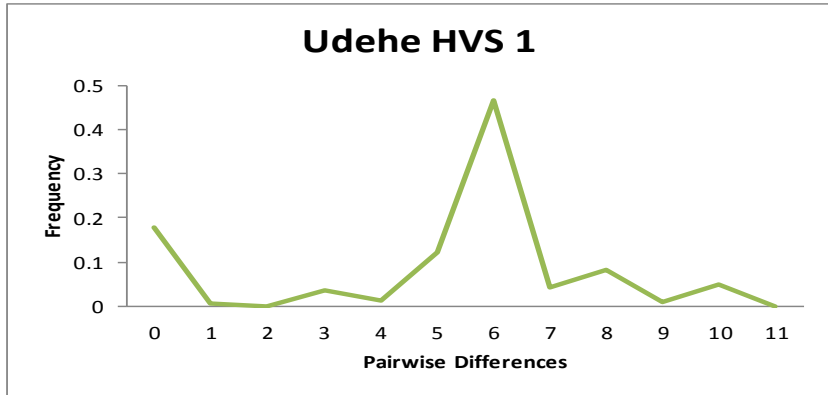
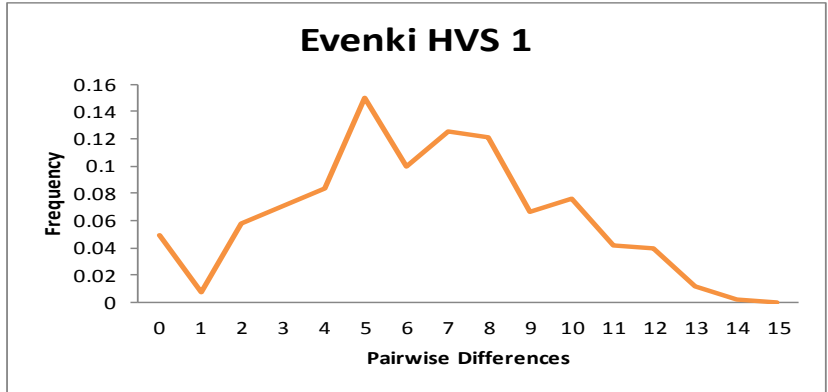
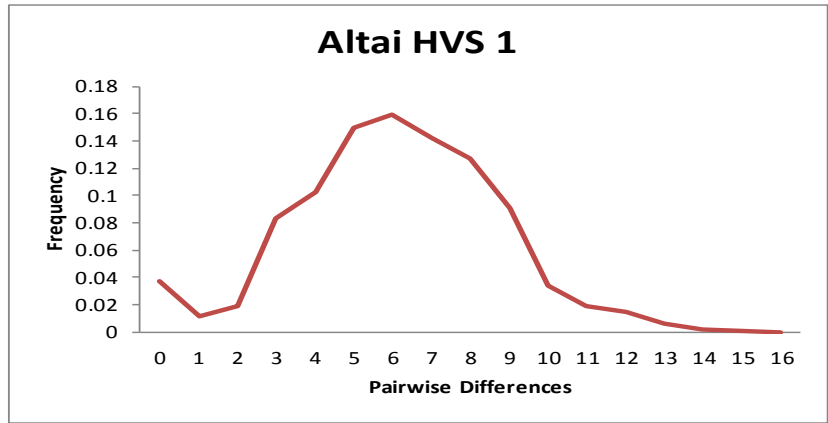


Figure 4.5.2a Displays Mismatch Distributions of Pairwise Differences from HVS 1 Sequence Data for the Altai, Evenki, Udehe and Yakut Populations.

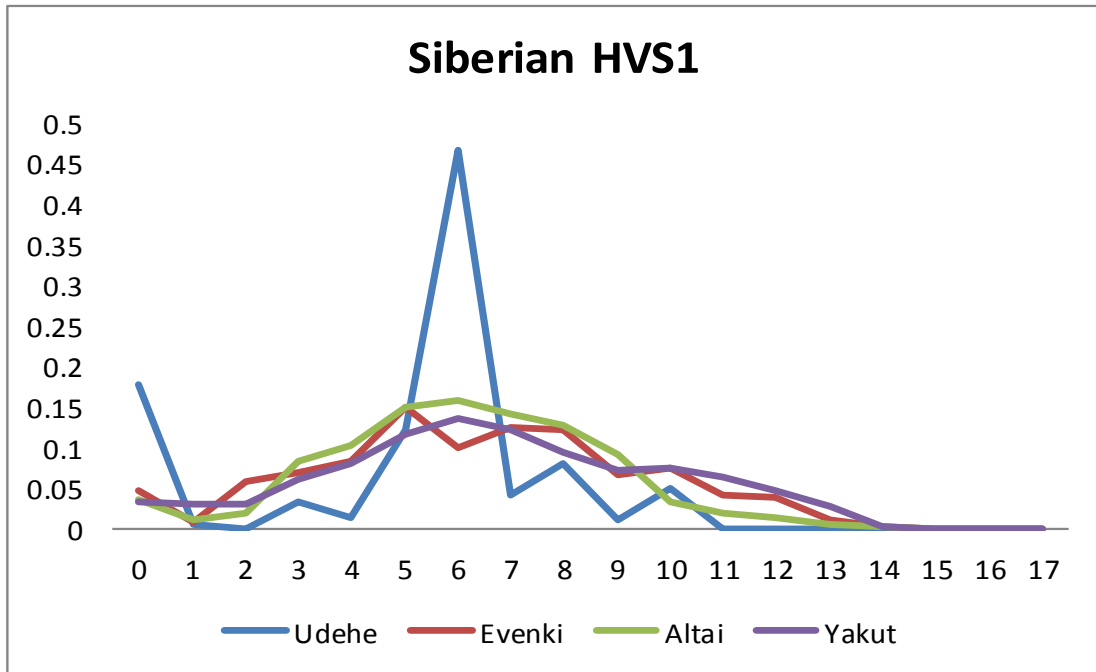


Figure 4.5.2b Displays Mismatch Distributions of Pairwise Differences for HVS 1 Sequence data for the Udehe, Evenki, Altai, and Yakut populations.

Since the Yakut are sampled from nine different villages and have been treated, whenever possible, as nine distinct populations for comparative purposes, HVS1 mismatch distributions were generated for the respective Yakut communities and are represented in Figure 4.5.2c, 4.5.2d, and 4.5.2e. The majority of the Yakut populations (Yakut-DA, DE, YJ, YK, YM, and YO) revealed unimodal distributions with major peaks between 5-8 pairwise differences. In all Yakut populations but the YakutCK and YakutYE, there was a minor mode at 0 pairwise differences. YakutCK, YakutDY, and YakutYE are considered multimodal with highest peaks at 5, 7, and 7 pairwise differences, respectively. YakutCK and YakutYE showed secondary peaks at 11 pairwise differences. The YakutDY also revealed a tertiary peak at 13 differences. All of the Yakut villages showed low raggedness indices ( $r < 0.04$ ). Figure 4.5.2f depicts a combined mismatch distribution of all of the nine Yakut villages.

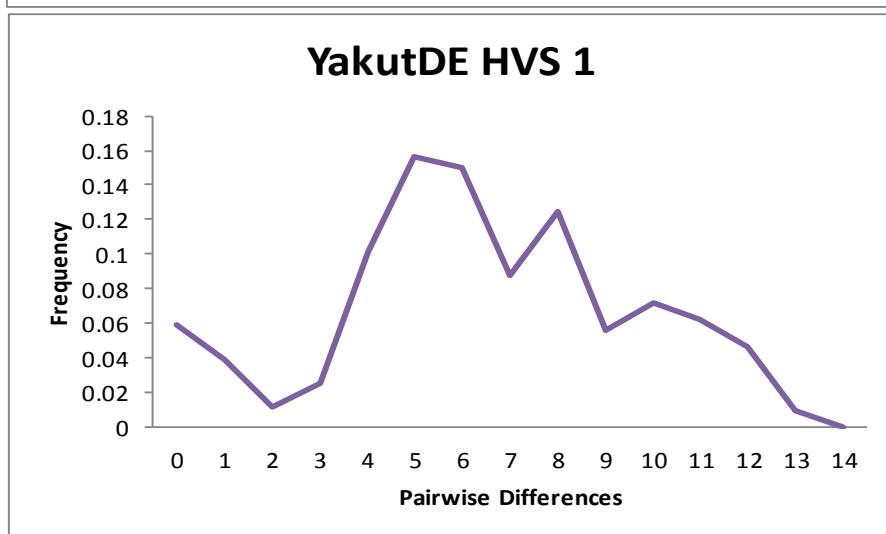
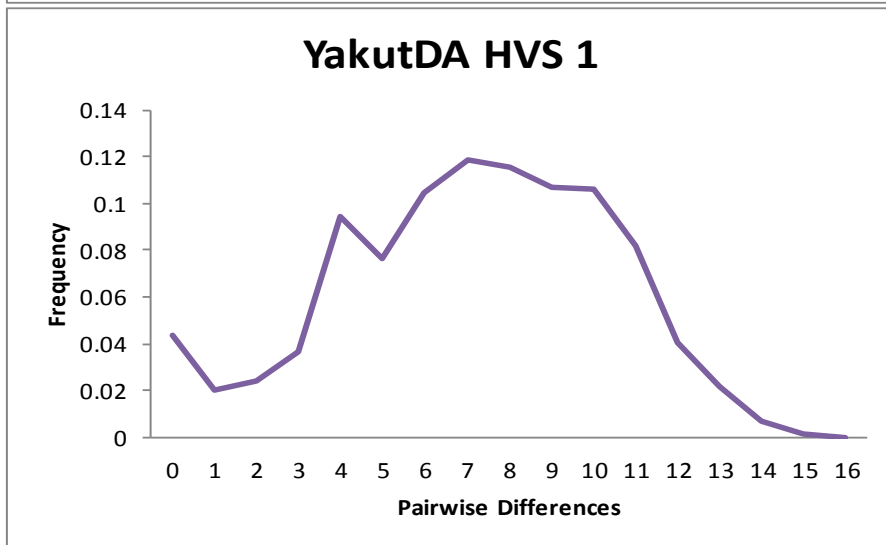
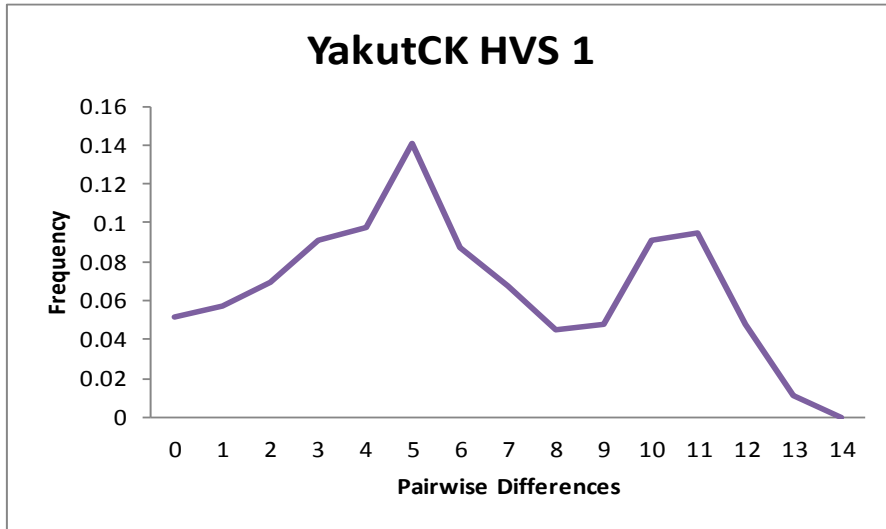


Figure 4.5.2c Displays Mismatch Distributions of Pairwise Differences from HVS 1 Sequence Data for the YakutCK, YakutDA, and YakutDE.

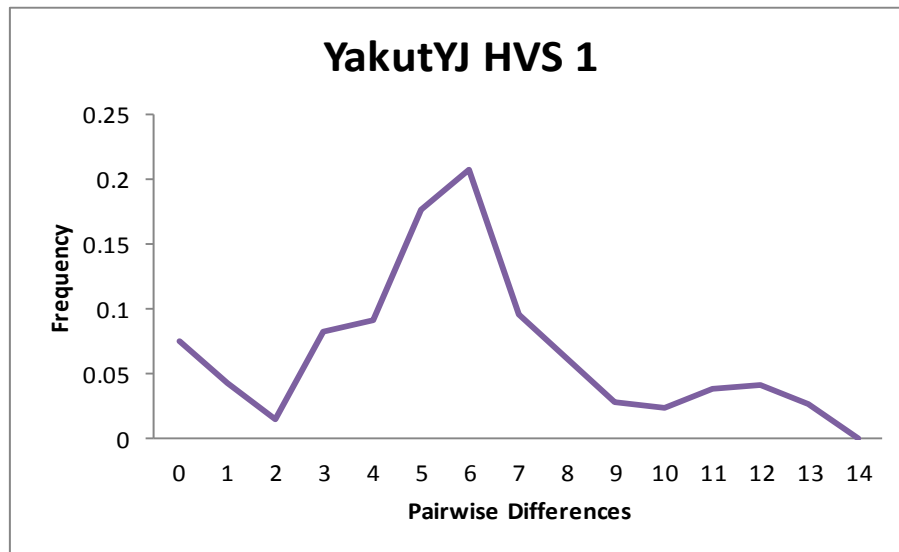
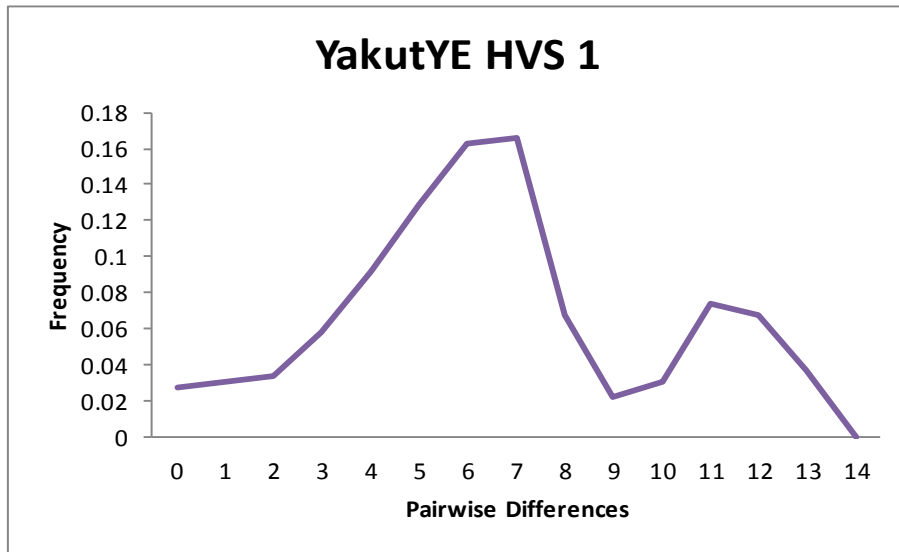
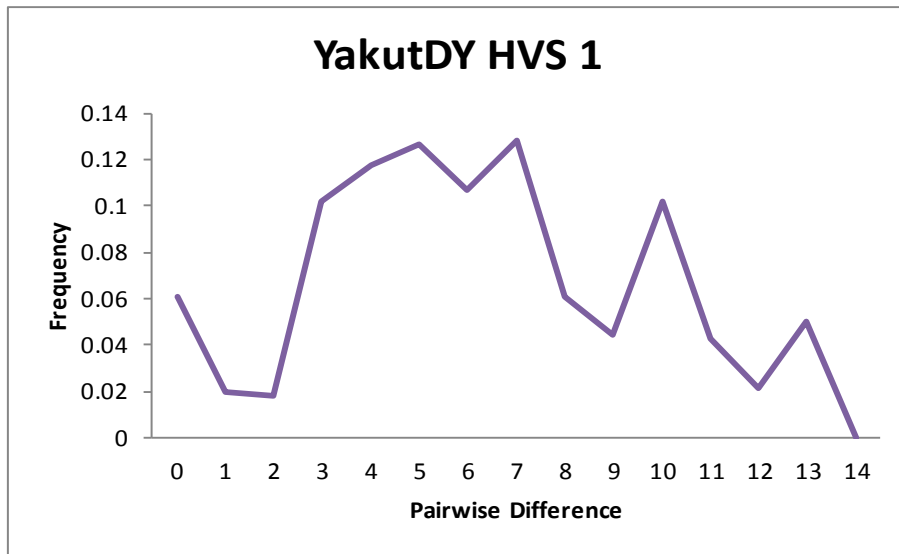


Figure 4.5.2d Mismatch Distributions of Pairwise Differences from HVS 1 Sequence Data for the YakutDY, YakutYE, and YakutYJ.

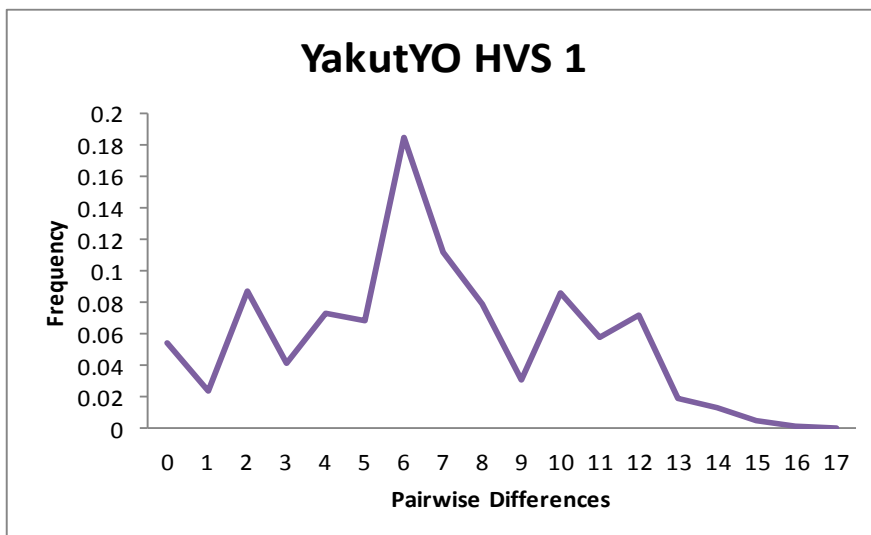
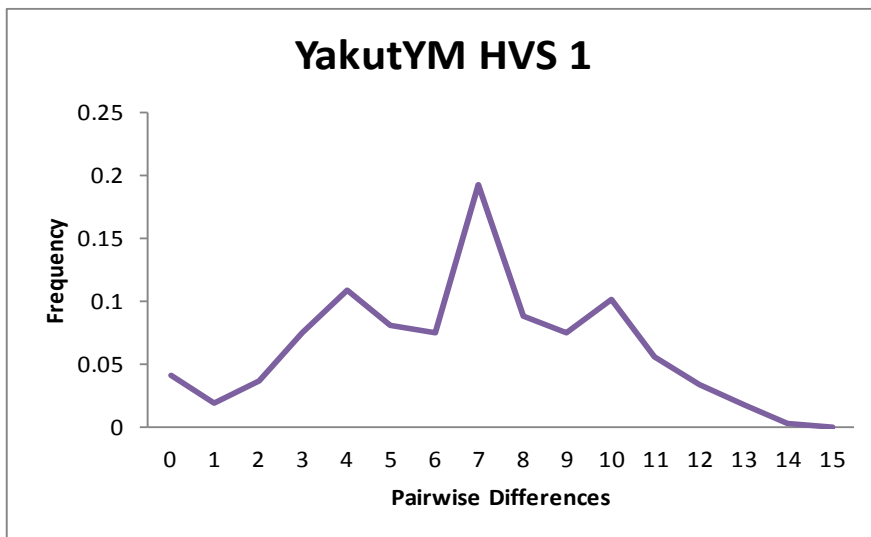
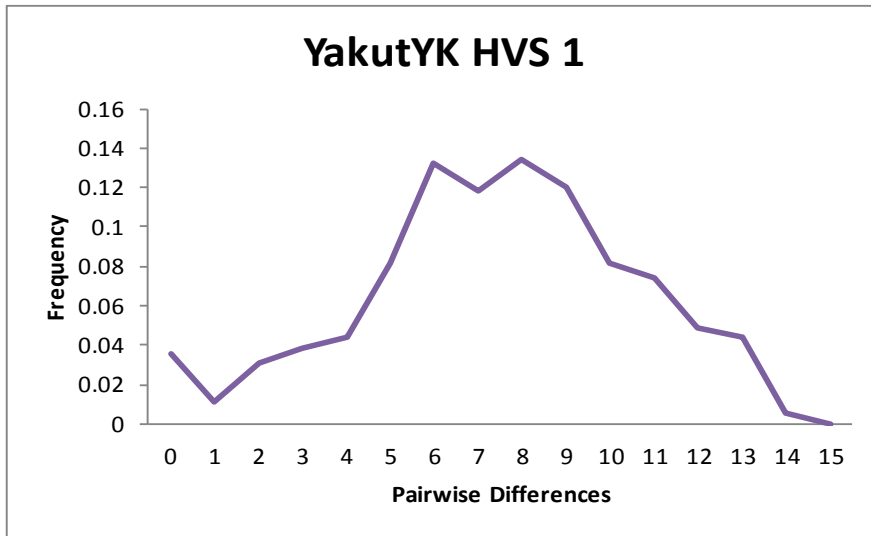


Figure 4.5.2e Displays Mismatch Distributions of Pairwise Differences from HVS 1 Sequence Data for the YakutYK, YakutYM, and YakutYO.

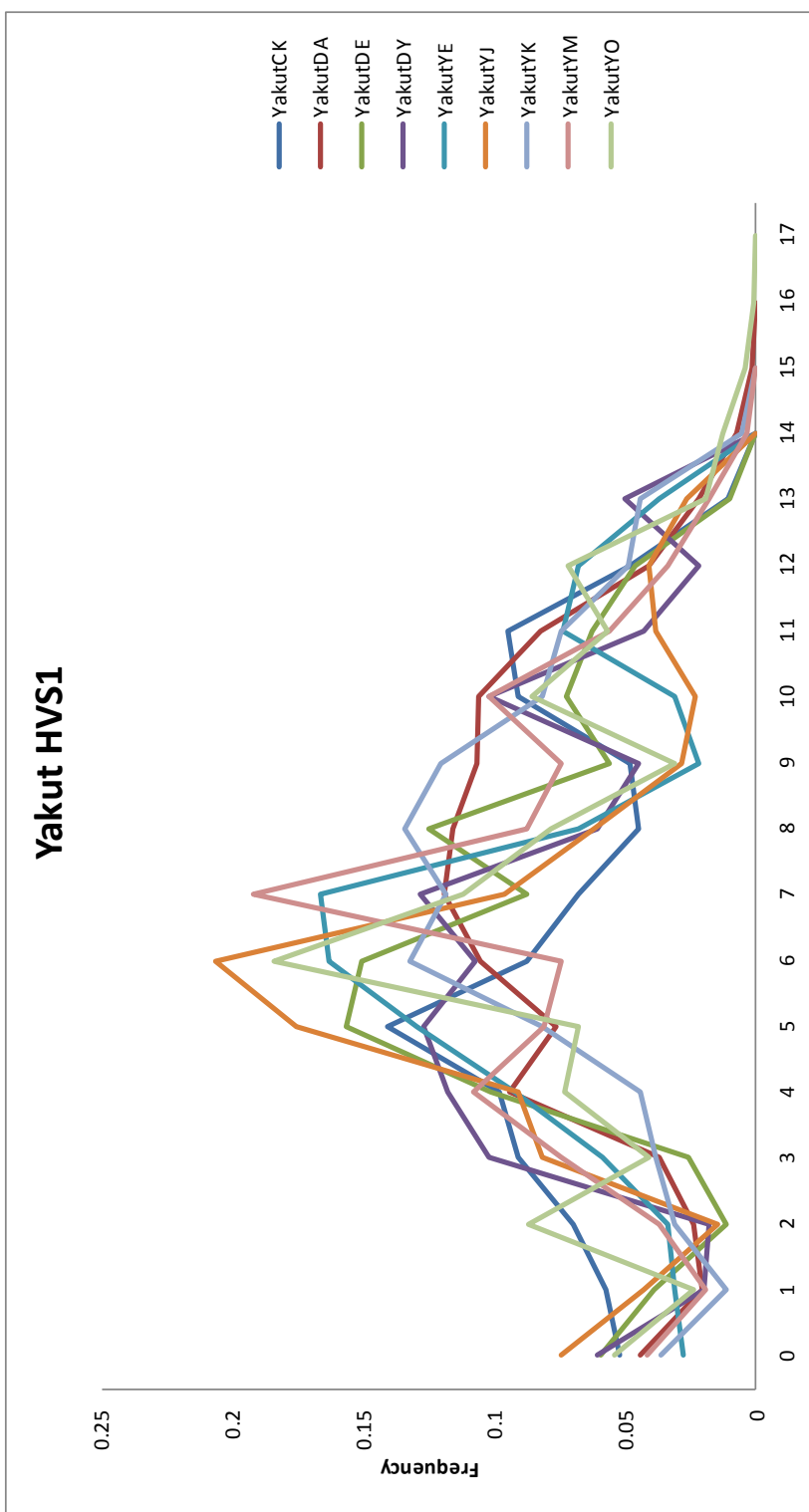


Figure 4.5.2f Displays Mismatch Distributions of Pairwise Differences from HVS 1 Sequence Data for each of Yakut Villages.

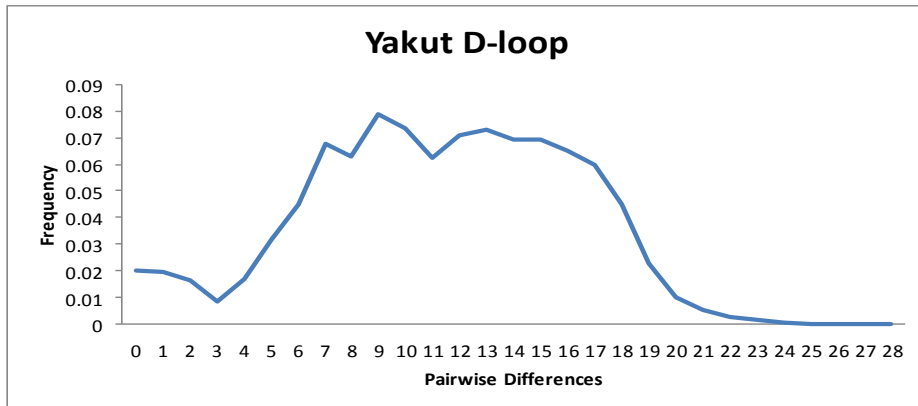
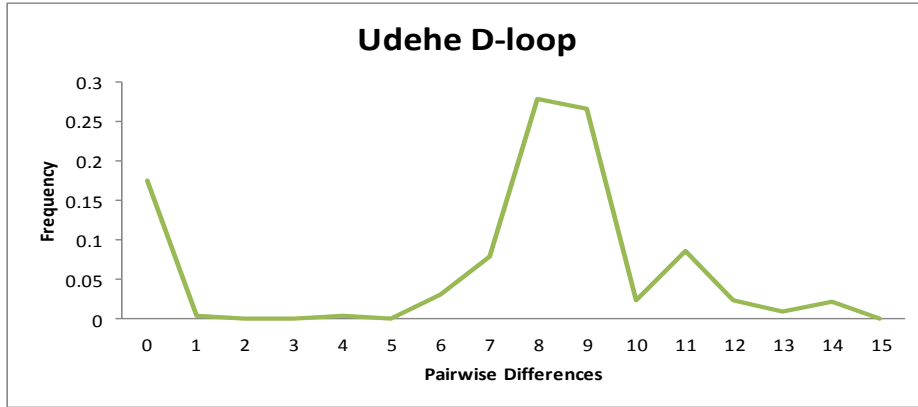
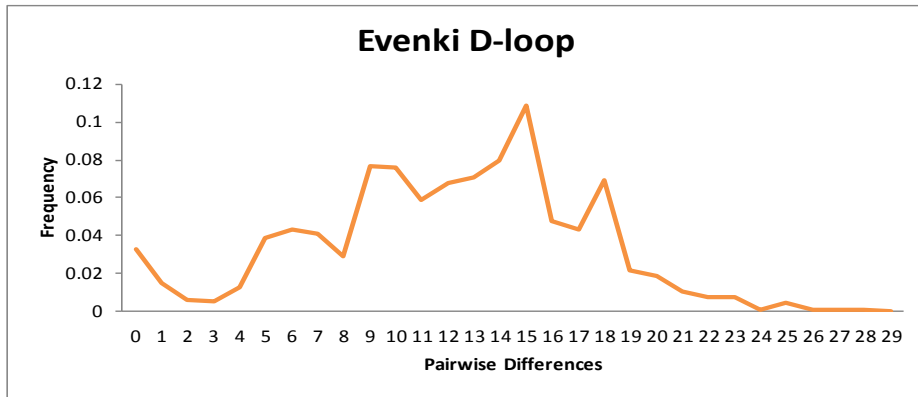
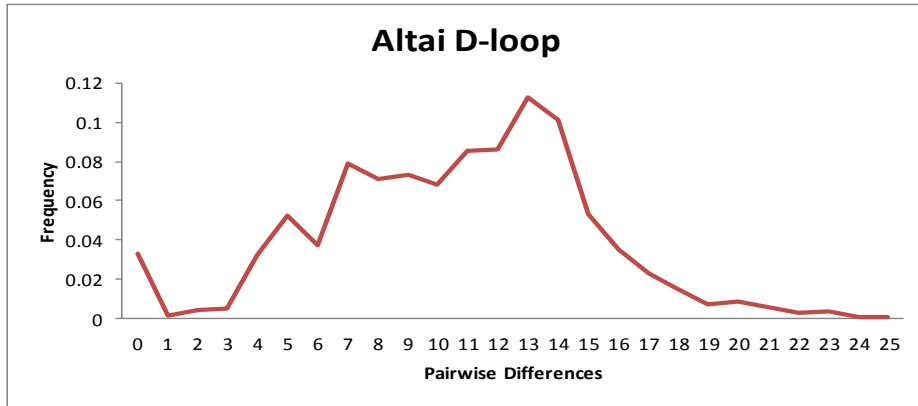


Figure 4.5.2g Displays Mismatch Distributions of Pairwise Differences of D-loop Sequence data for the Altai, Evenki, Udehe and Yakut.

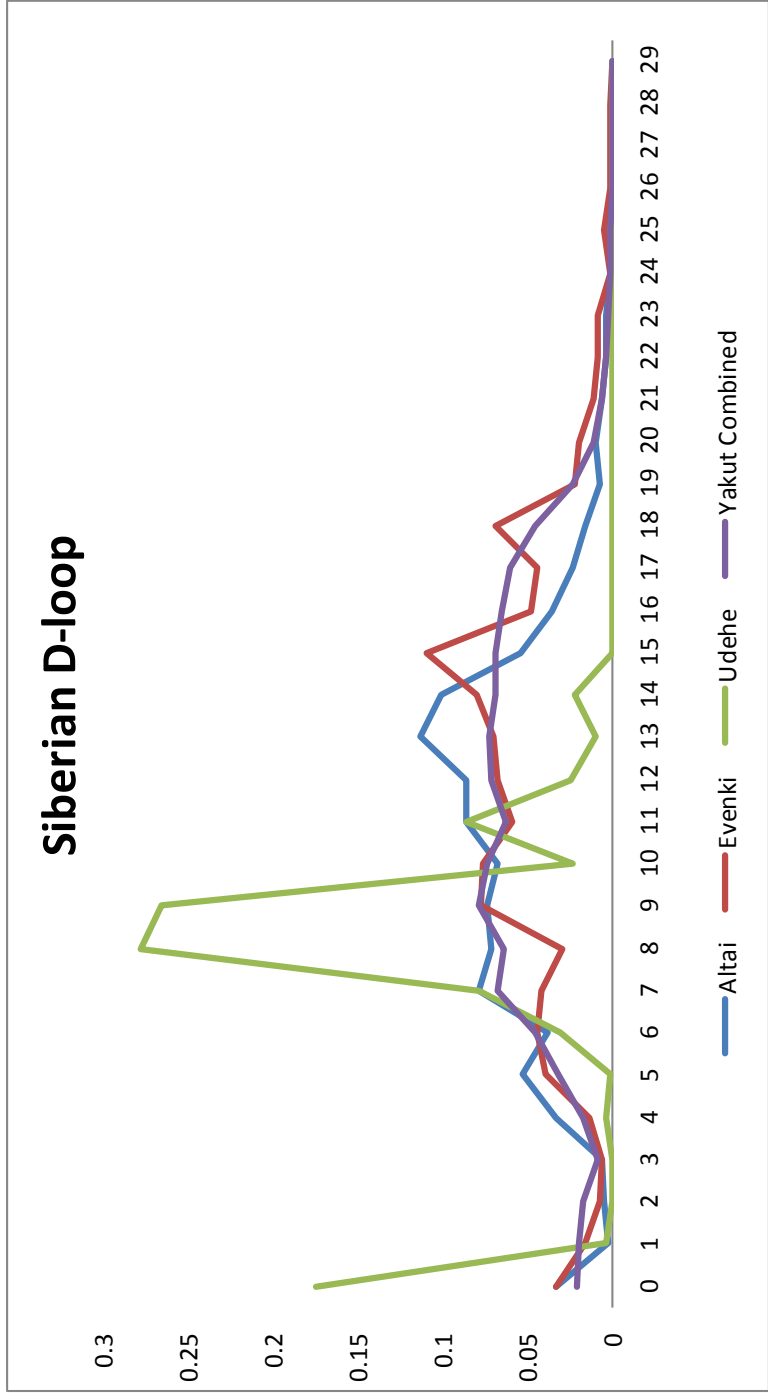


Figure 4.5.2h Displays Mismatch Distribution of Pairwise Differences of D-loop Sequence Data for the Altai, Evenki, Udehe, and Yakut.

Similar to the HVS1 data, the Altai, Evenki, Udehe and combined Yakut all displayed relatively unimodal distributions with highest peaks between 8-13 pairwise distances, and all showed minor modes at zero mismatches. The highest raggedness index was again seen among the Udehe ( $r = 0.139$ ). Overall, the D-loop sequence data reveals distributions that have lower raggedness indices for all populations except for YakutCK and YakutDA compared to those for HVS1. The mismatch distributions for these four populations are depicted together in *Figure 4.5.2h*. Comparing these results, it appears that not only does the addition of SNPs reveals more minor modes, but the number of pairwise differences observed also appears to grows. In all cases, the Mean Pairwise Differences among each population increased, as did the associated confidence intervals.

Mismatch distributions of Pairwise differences for the D-loop of the individual Yakut villages all appear to be generally unimodal, except for the YakutCK and YakutYE, which are multimodal; these distributions are displayed in *Figure 4.5.2i*, *4.5.2j*, and *4.5.2k*. Likewise, all Yakut villages save the YakutCK had minor modes at zero pairwise differences. The YakutCK had the largest peak at 16 differences with a secondary peak at 5. The YakutYE demonstrate a multimodal distribution with a primary peak at 8 pairwise differences and a large secondary peak at 17 differences. *Figure 4.5.2l* depicts all of the D-loop mismatch distributions for the various Yakut communities together.

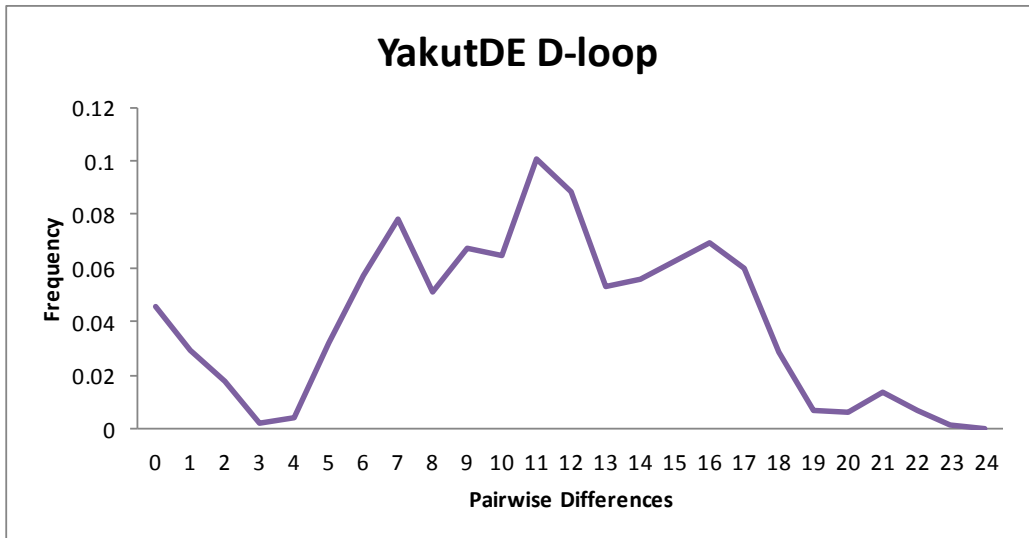
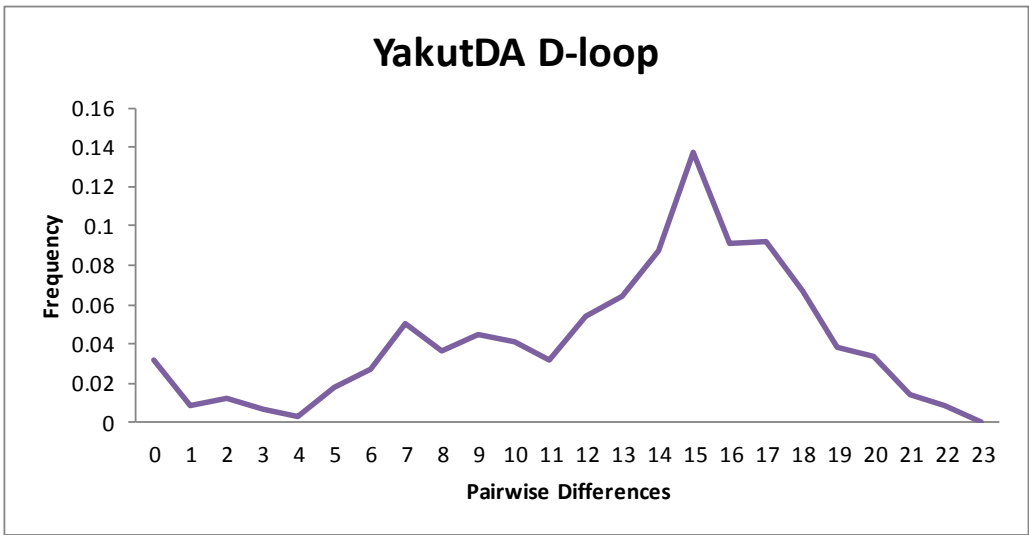
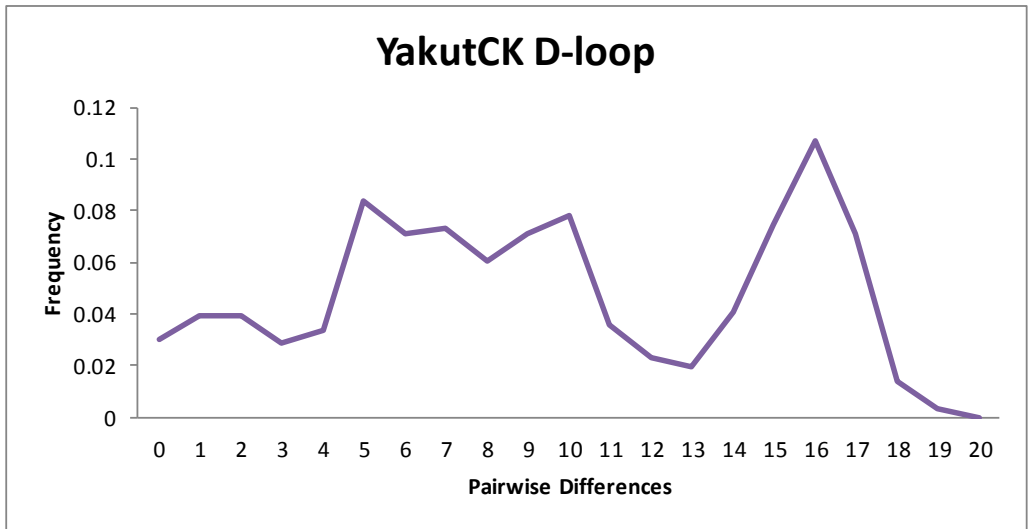


Figure 4.5.2i Displays Mismatch Distributions of Pairwise Differences from D-loop Sequence Data for the YakutCK, YakutDA, and YakutDE.

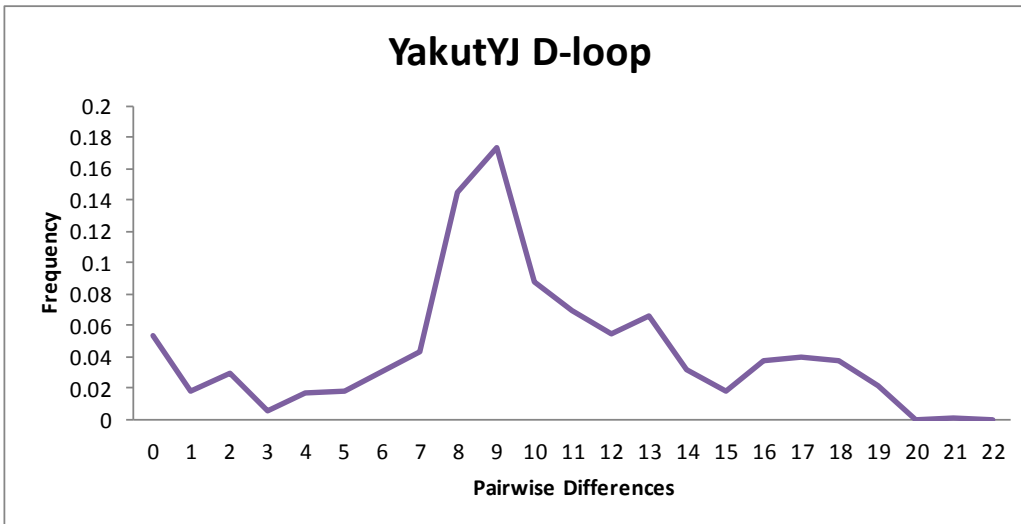
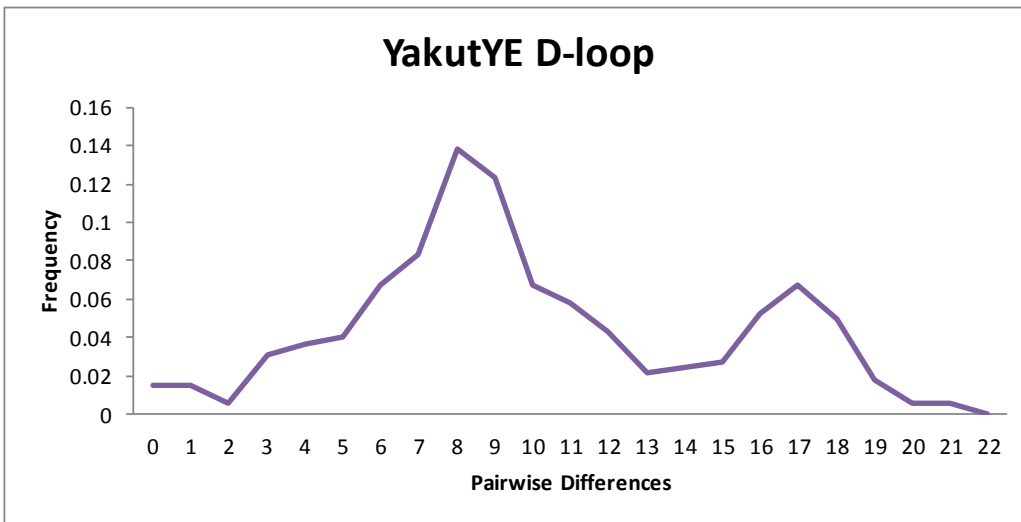
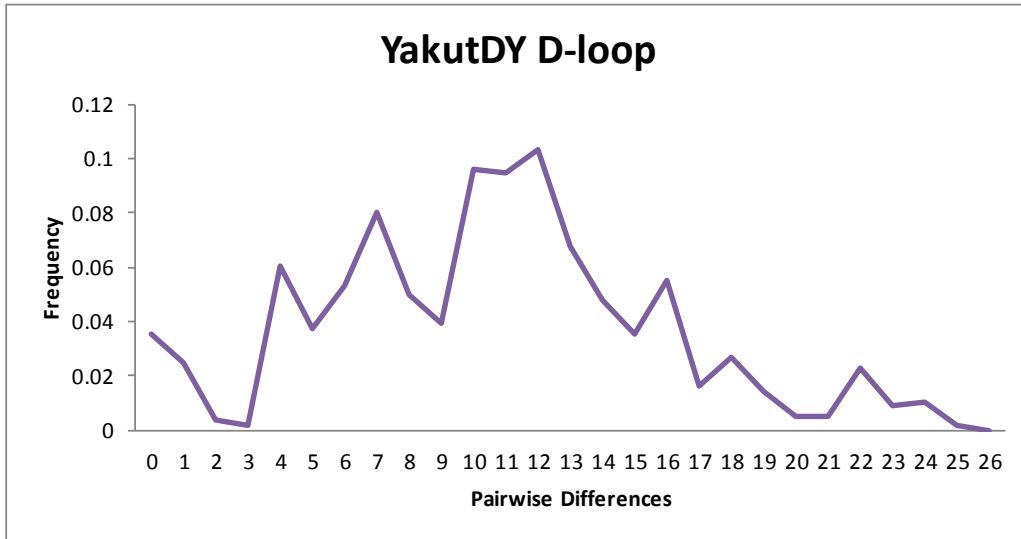


Figure 4.5.2j Displays Mismatch Distributions of Pairwise Differences from D-loop Sequence Data for the YakutDY, YakutYE, and YakutYJ.

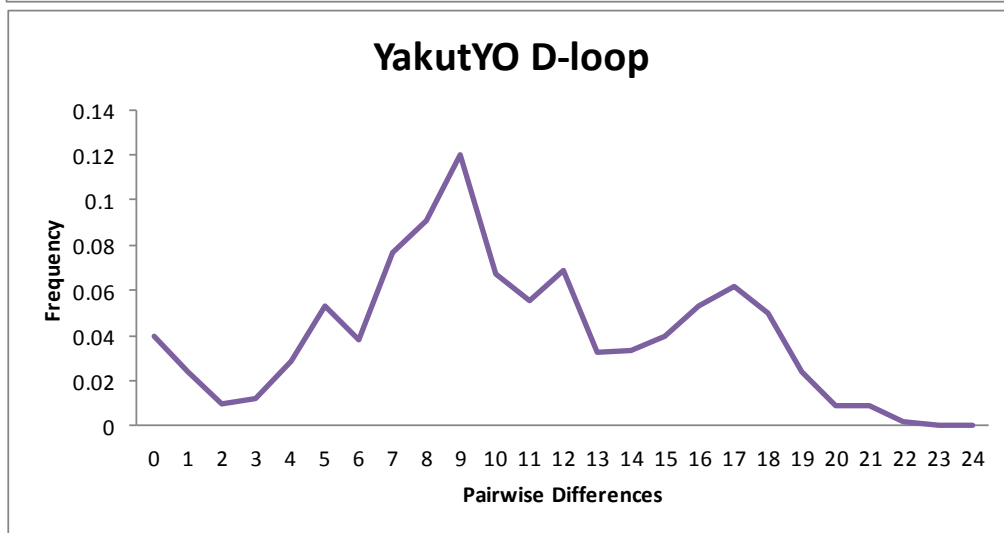
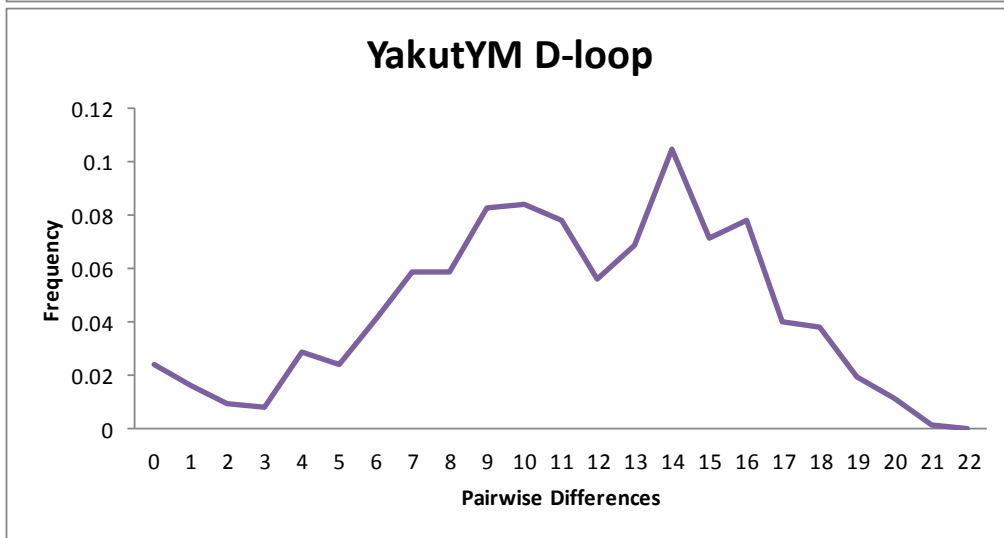
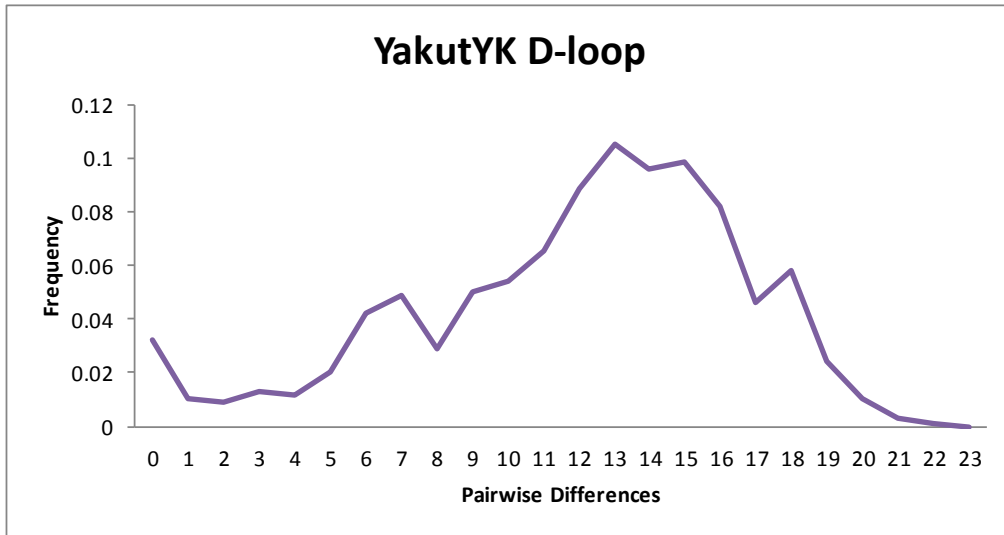


Figure 4.5.2k Displays Mismatch Distributions of Pairwise Differences from D-loop Sequence Data for the YakutYK, YakutYM, and YakutYO.

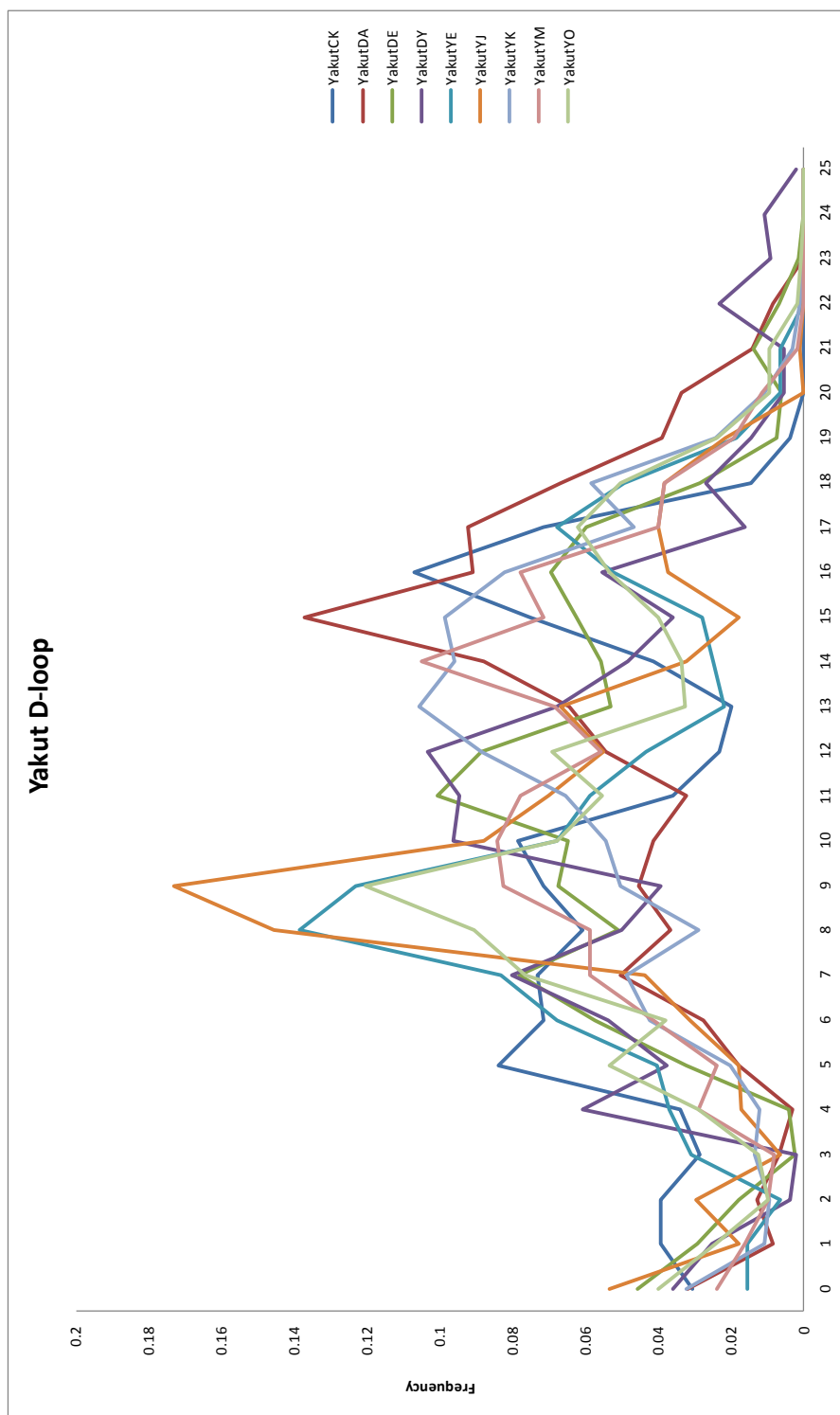


Figure 4.5.2| Displays Mismatch Distributions of Pairwise Differences from D-loop Sequence Data for each of the Yakut Villages.

## 4.6 Phylogeography

### 4.6.1 Mantel Randomization

Mantel Randomization tests were conducted in order to test the relationship between genetics and geography, which seeks to determine the significance of the correlation between a pairwise distance matrix of the mtDNA sequence data and a distance matrix created for geographical distances of the populations of study. Tests were performed using bootstrap replication of 1000 iterations and the results of these tests are presented in *Table 4.6.1a*. Geography was held as a constant (Matrix X) in all cases, and tested against HVS1 and D-loop genetic distance matrices (Matrix Y) for a.) the four main populations of study, b.) the same populations with the Yakut treated as nine distinct populations, and c.) a random sample of all of the Yakut, which allowed for a sample size much closer to that of the rest of the populations of study. Out of all of these tests, the only correlation coefficient that was significant was that comparing geographical distances to the whole D-loop genetic distances with the Yakut separated into the nine villages ( $r = 0.42$ ,  $p < .05$ ). Interestingly, the same mantel test grouping comparing geographical distances to HVS1 genetic distances of all of the Siberian populations with the Yakut separated showed a stronger correlation, but the result was less significant. ( $r=0.58$ ,  $p=0.077$ ).

*Table 4.6.1a Mantel Randomization*

<b>Matrix X</b>	<b>Matrix Y</b>	<b>Correlation Coefficient <i>r</i></b>	<b>P-Value</b>
Geography	HVS1 (Combined Yakut)	-0.4097	0.8261
Geography	HVS1 (Separated Yakut)	0.58	0.077
Geography	HVS1 (Random Yakut Sample)	0.0435	0.4783
Geography	D-loop (Combined Yakut)	-0.4122	0.8261
Geography	D-loop (Separated Yakut)	0.42*	<b>0.039</b>
Geography	D-loop (Random Yakut Sample)	0.1691	0.5217

*This figure displays results from Mantel Randomization Tests comparing geographical to HVS 1 and D-loop distance matrices.*

*\* Statistically significant*

#### 4.6.2. Mantel Randomization and HVS1 vs. D-loop

To determine if there was a significant correlation between the HVS1 data and whole D-loop data, a Mantel Randomization Test was performed with 1000 permutations so as to compare the two respective pairwise distance matrices (FSTs). The results for these tests are presented in *Table 4.6.1b*. This analysis was performed considering: a.) the Yakut as a single population, b.) the Yakut separated into populations based on village from which they were collected, and c.) with a random sample of Yakut to make the sample sizes more similar among the four populations. In all cases, the correlation coefficient was extremely high (>0.9), with high significance values ( $p < .01$ ).

*Table 4.6.1b Mantel Randomization*

<b>Matrix X</b>	<b>Matrix Y</b>	<b>Correlation Coefficient <i>r</i></b>	<b>P-Value</b>
HVS1 (Combined Yakut)	D-loop (Combined Yakut)	0.9877*	<.0001
HVS1 (Separated Yakut)	D-loop (Separated Yakut)	0.9429*	0.004
HVS1 (Random Yakut Sample)	D-loop (Random Yakut Sample)	0.9735*	<.0001

*This figure displays results from Mantel Randomization Tests comparing HVS1 to D-loop genetic distance matrices.*

*\*Statistically Significant*

## CHAPTER 5: DISCUSSION

The following chapter seeks to interpret the results from the previous chapter in order to make observations about the questions under investigation. The discussion will divide the results of the all analyses conducted for both HVS1 and the D-loop data and will examine comparisons between the two. These outcomes will be integrated into a discussion of the proposed research questions from chapter 1, which include: 1.) if the increase in the number of SNPs sequenced revealed different phylogenetic relationships between Siberian populations; 2.) if additional genetic variation can be revealed by the addition of more genomic regions; and 3.) whether additional SNPs reveal stronger relationships between genetics, linguistics, and geography than using the HVS1 alone.

### 5.1 mtDNA Lineages and Within Group Variation

Mitochondrial DNA lineages of the four sample populations are consistent with previous findings within Siberian groups. For the Majority of Siberian populations, representative haplogroups include: A, B, C, D, F, G, and X, with haplogroup C and D being the most prevalent. Of these, haplogroups, B and X tend to be the least frequent (Torroni et al., 1993; Ballinger et al., 1992; Starikovskaya et al., 1998; Puzyrev et al., 2003). Unlike the Altai, Evenki and Yakut, the Udehe are outliers, characterized solely by haplogroups C, M, N, and Y, where 88% of the sample was not considered to be one of the typical Siberian haplogroups. This is similar to the findings of Torroni et al. (1993), indicating that the Udehe are composed of approximately 19% haplogroup C and 20% “other.” This is likely due to gene flow and heavy East Asian influence resulting from geographical proximity. Between 4 and 23% of individuals sampled from the remaining three populations studied (Altai, Evenki, and Yakut), could not be characterized for mtDNA haplogroups based on HVS1 data alone. However, when the rest of the D-loop was added to each group, all remaining ambiguous haplogroups were

resolved. This shows that adding more loci will provide finer resolution when defining haplogroups based on SNPs.

Mitochondrial DNA results based on HVS1 data illustrated that the Altai, Evenki, and Yakut are represented by high gene diversity ( $h = 0.9944$ ) with the YakutYE displaying the highest levels among the Yakut villages at  $h = 0.9754$ . Results from D-loop data revealed the similar results, verifying that the Yakut had the highest gene diversity at  $h = 0.9960$ , with the YakutYE village measuring highest among them at  $h = 0.9846$ . Nucleotide diversity was the highest for combined Yakut samples based on HVS1 data (0.016815), with the YakutYK displaying the highest among Yakut villages (0.018519). The highest nucleotide diversity overall for D-loop data was found in the YakutDA village at 0.12945. However, when comparing the Yakut as one population to the Altai, Evenki, and Udehe, the Evenki showed the higher nucleotide diversity at 0.011929. These values are concurrent with previous Siberian studies and are likely to be indicative of populations that have undergone demographic and spatial expansion (Phillips-Krawczak et al., 2006; Zlojutro et al., 2008). The lowest gene and nucleotide diversity was present in the Udehe for both HVS1 and D-loop sequences, where gene diversity was 0.8256 in both instances and nucleotide diversity was 0.012829 and 0.007276 respectively. This should be expected due to the low levels of haplogroup diversity present in the sample. The mtDNA gene diversity seen in the Udehe is low compared to other Siberian populations, which is suggestive of a small, isolated population. Low gene diversity could also be indicative of genetic drift, which can have dramatic effects on small populations; the Udehe number about 1000 today (Bermisheva et al., 2005). This could also be due to an overrepresentation of related individuals in the sample. When considering results comparing both HVS1 and the entire D-loop, it appears that the addition of HVS 2 and 3 does not significantly alter the genetic diversity measures.

## 5.2 Variation among Populations

### 5.2.1 AMOVA

According to AMOVA results, little substructure exists among populations; when the groups were organized by linguistics, geography, or ethnicity. Over 94% of the variation was explained within populations and little is explained among populations or among groups within populations. This is true for both HVS1 as well as the D-loop sequence data. Greater genetic diversity occurring within a population as opposed to among populations is typical of human population studies. Barbujani et al. (1997) illustrated this in a study involving 16 populations from around the globe, demonstrating that over 84% of the genetic variation was explained within populations. Further, Stone and Stoneking (1998) found similar results using AMOVA, suggesting that over 74% of the variation in Native American populations was within groups as well. The results of this study are concordant with previous findings of variance apportionment in the region, possibly illustrating the homogeneity of these populations (Derenko et al., 2003; Jin et al., 2003). The addition of the rest of the D-loop (HVS 2 and 3) did not affect the partitioning of explained variation within samples when using AMOVA. Likewise, the addition of HVS2 and 3 revealed little changes to the statistical significance of the results (see *Tables 4.4.1a-f*).

### 5.2.2 Multi-dimensional Scaling

All multidimensional scaling (MDS) plots were characterized by low stress values (below 0.19), indicating an acceptable goodness of fit when compared to the original distance matrices (Kruskal, 1964). Initial results for HVS1 data were unexpected, depicting the Udehe as extremely divergent in both dimensions of the plot from the other Siberian groups, resulting in the remaining populations to form a tight cluster in the opposite of both dimensions of the plot. This spatial representation of genetic distance was misleading, as the Udehe were not expected to be as distinct compared to the other

Siberian populations of study, due to their similar geographical proximities. Furthermore, it has been suggested that the Udehe may be closely related to the Evenki (Crawford et al. 2002; Jin et al., 2009) as they share a variety of mtDNA haplotypes, and speak a similar Tungusic language. Because of this highly unexpected MDS plot, an AMOVA was also run without including the Udehe to see if this was altering the apportionment of variation among the populations. The result was that even without the Udehe (HVS1) data, there was still over 94% of the variation explained within groups ( $p < .01$ ). Removing the Udehe from the MDS plot, then allowed for a better resolution of the relationship of the Altai, Evenki, and Yakut populations based on HVS 1 data. The remaining 11 populations fell into three main clusters, categorized by ethnicity. Though the Yakut were broadly scattered within their respective cluster, they fell completely along the left side of the plot, while the Altai and Evenki were further removed on the right. The Altai appear to cluster more closely to the Yakut group than to the Evenki, which may be indicative of shared ancestral ties between the two groups; since both speak Turkic languages and are postulated to have originated in Southern Siberia. The addition of the sequence data from HVS 2 and 3 significantly impacted MDS plots based on genetic distances. The Udehe, most notably, were not extremely divergent when considering all three hypervariable regions, though they seemed almost completely unrelated with HVS 1 data alone. With the exception of the YakutYJ, all Yakut villages clustered with Altai in close proximity. This could possibly be due to a higher proportion of European haplotypes within the village. Compared to the rest of the Yakut villages, the YakutYJ have the highest representation of haplogroup H (17%), based on D-loop sequences data. One explanation for YakutYJ being non-divergent in MDS based on HVS 1 data is possibly due to an inability to characterize any of these samples as haplogroup H based solely on this genomic region (HVS 1). The addition of the remaining hypervariable regions (2 and 3) significantly affect haplogroup prediction and spatial characterization of populations based on genetic distances in multidimensional scaling.

### 5.2.3 Neighbor-Joining Trees

Two different types of neighbor-joining trees (NJT) were constructed based on Kimura 2P genetic distances for both HVS1 and D-loop data. Based on the HVS 1, the first set involved using a population based distance matrix in order to examine the relationship between populations that resulted in a tree composed of two main branches, where the Udehe occupy one branch as a separate operational taxonomic unit (OTU). This branch is also the longest, (i.e., the longest genetic distance), suggesting the Udehe have deviated the most from the remaining populations. These inferences were not surprising, based on the results of the MDS plots. The second branch bifurcates into tertiary branches, with the first leading to the Altai and Evenki as separate OTUs, and the other depicting all of the Yakut villages as sub-branches. Considering the MDS plots, as well as the similar composition of mtDNA Haplogroups, it is surprising, however, that the large Yakut limb is not closer to that of the Altai.

Various notable differences emerged upon the addition of HVS 2 and 3 sequence data. The clusters of Yakut OTUs have been reorganized, with YakutDA and YakutDY no longer considered closest “neighbors” to the Yakut-YJ, -DE, -CK, and -YE. According to the D-loop data, these four Yakut populations are more closely related to YakutYK. Additionally, this second tree shows variation in branch length for many populations. Specifically, the Evenki appear more genetically distant from the Altai, though they have remained in the same position on the NJT. The YakutYJ also show an increased branch length, second in length only to the Udehe, which suggests a “pseudo-outlier” when compared with the rest of the Yakuts. This was expected following the MDS plots and diversity indices, and again may signify a sampling of a higher frequency of European genes.

Though it appears that the additional loci available when adding the 2<sup>nd</sup> and 3<sup>rd</sup> hypervariable region of the D-loop yields different relationships between populations, the goal of this project is to determine whether or not this increases the resolution of mtDNA findings. To determine whether the

tree was a good representation of the original matrix, cophenetic distance matrices were generated from the NJTs and then compared to the original distance matrices using a Mantel randomization test. Mantel test revealed a high correlation between the two HVS1 trees ( $r=0.92213$ ,  $p<0.001$ ). Results for the D-loop data indicated a moderately high correlation value ( $r=0.7533$ ,  $p<0.001$ ). The reduction in correlation coefficient value suggests that for creating NJTs based on distance matrices, a dendrogram based on HVS 1 data reveals population relationships that more closely represent the observed genetic distances. Results did not significantly improve with the addition of the second and third hypervariable regions.

The second set of NJTs was constructed with individual sequences instead of population distance matrices as the input data. This is a method often used by biologists examining the evolutionary relationship among various same species individuals within a sample collection, or among various species, using one individual from each species to test the evolutionary relationship between species in general. For the purposes of this study, every individual was considered its own taxon, translating as its own OTU within the dendrogram. This method allowed for bootstrap replication of the trees in order to better compare the statistical accuracy of each dendrogram (HVS1 and D-loop). Increased branching was observed in dendrograms using D-loop, likely due to the presence of SNPs, increasing chances for additional deviations between individuals. Furthermore, increases in sequence haplotypes were observed within each population, therein increasing observable genetic variation. In both phylogenies, the Altai, Evenki and Yakut dispersed throughout the entire tree, often occupying the same branch as individuals from the other groups. The Udehe on the other hand typically cluster alone on sub-branches without any of the other three populations present. This is not surprising based on other representations of inter-population genetic variation described in this study. It also may be indicative of genetic drift affecting the Udehe, or due to a small sample size. It is important to note that

because this method considered every individual and not population data, the relationships inferred in the dendrograms are likely a reflection of ancestral lineages or haplotypes among haplogroups. (e.g., D2a), and not inter-population relationships.

Similar to the prior set of NJTs, when adding the rest of the D-loop sequence data, the branch length and orientation of various branches and sub-branches are altered slightly. A more noteworthy observation, however, is the fact that bootstrap values change as well. As reported in the results section of this work, the magnitude of bootstrap values above 50%, 80%, and 95% all increase in the D-loop tree. Again, bootstrapping is a resampling technique that will provide a percentage of how many times a particular branch was created out of a certain number of iterations. For the purposes of this study, 1000 replicates were used. This increase in high bootstrap values suggests that adding the additional sequence data associated with the HVS 2 and 3 sequence data provides better resolution of the evolutionary relationship among individuals.

### **5.3 Forces of Evolution**

#### **5.3.1 Neutrality Tests**

Neutrality test statistics (Tajima's  $D$  and Fu's  $F_s$ ) and mismatch analysis were used to infer demographic change and evolutionary forces at work on the genetic composition of study. Results of these analyses were compared for HVS1 and the entire D-loop, considering the Yakut both as one population and as individual populations with respect to individual villages. Tajima's  $D$  results for HVS 1 data resulted in negative values for all populations, suggesting these populations are undergoing expansion. This statistic for the Udehe, however, was close to zero (-0.05764), which is not as strong of an indicator of expansion as large negative values. Although all  $D$  statistics were negative, these values were only statistically significant ( $p < 0.05$ ) for the Altai and the Combined Yakut group. Moreover, these

two populations had the two largest negative D values for HVS1 data. Not surprisingly, the Altai and combined Yakut are comprised of the largest sample size (n) as well as the largest number of mtDNA lineages or haplogroups among the selected populations. This suggests that further analysis with larger sample sizes may reveal different results for the rest of the populations. In comparison, results from the entire D-loop yielded similar patterns with one notable difference. Tajima's D for the Udehe resulted in a positive value, suggesting genetic drift has affected this group. As was observed with HVS 1 data, the Altai and combined Yakut were the only two populations with a statistically significant result. This overall consistency in Tajima's D values suggests that little evolutionary information can be gained from the addition of more SNPs, or with the addition of the rest of the D-loop. However, the switch from a small negative value to a small positive number in the Udehe that was only observed by the addition of HVS 2 and 3, illustrates that it may be useful if the degree to which a resulting D value is negative or positive. The larger the negative value is with HVS 1, the less likely the addition of HVS 2 and 3 will diverge from that trend.

The neutrality test statistic Fu's  $F_s$  has been reported to be more sensitive to population growth and genetic drift compared to the Tajima's D in part because this assessment uses differences between haplotypic distributions, instead of the more conservative use of genetic pairwise differences (Fu, 1997; Zloturo et al., 2006). Results for HVS1 illustrated again an overall theme of expansion for all populations except the Udehe, demonstrated by negative values. Combined with a positive Tajima's D value, these  $F_s$  values indicate that the Udehe have experienced increased homozygosity due to genetic drift, possibly as a consequence of being geographically and culturally isolated from the three other Siberian groups. The Udehe appear to be more influenced ethnically by East Asian populations, as is seen in the high frequency of mtDNA haplogroups M, N, and Y. Concurrent with HVS1 results, the Altai and combined Yakut had the largest negative  $F_s$  values (-18.58077 and -24.33068, respectively). Overall,

$F_s$  values tended to be more statistically significant compared to  $D$  values, where three populations (YakutDA, YakutDY, and YakutYE) exhibited  $p$ -values less than 0.05 and five (Altai, Combined Yakut, YakutCK, YakutYK, and YakutYM) exhibited  $p$ -values less than 0.01. It should be noted, however, that according to Fu (1997), only results at  $p < 0.02$  should be considered statistically significant, leaving  $p < 0.05$  as marginally significant. Results from D-loop data indicate an overall trend of expansion for all populations except for the Udehe, which had a positive value, suggesting genetic drift.

Results for these two neutrality tests (Tajima's  $D$  and Fu's  $F_s$ ) were consistent with previously reported findings for these populations (Derenko et al., 2000, 2003, 2007; Phillips-Krawczak, et al. 2006; Shields, et al. 1993; Zlojutro et al., 2006), suggesting these Siberian populations have experienced population expansion. Although the results of Tajima's  $D$  were not greatly affected by the addition of HVS 2 and 3, several Fu's  $F_s$  values were altered, by strengthened statistical significance, suggesting that the addition of the SNPs of the HVS 2 and 3 allow for better statistical inference in some cases. However, since the increase in significant results from HVS1 to D-loop is only noted in a third of the populations, further research is needed to reach a more substantial conclusion.

### **5.3.2 Mismatch Analysis**

Mismatch analyses is another methodology that can be useful in describing forces of evolution that may be affecting populations. HVS 1 data revealed mismatch distributions that were predominately unimodal with a secondary minor peak at zero mismatches for all groups except for the YakutCK, YakutDY, and YakutYE. Of these three, the YakutYE also had a minor peak at zero mismatches. Unimodal distributions are typical of populations undergoing demographic expansion, whereas bi- or multimodal distributions are indicative of populations at demographic equilibrium (Hudson and Slakitin, 1991; Rogers and Harpending, 1992; Joblong et al., 2004). A small mode at zero pairwise differences may result from an overrepresentation of related individuals due to sampling bias and overall small

population sizes. These results are consistent with neutrality test statistics of these groups. Contrary to the neutrality tests that suggested all of the Yakut populations were expanding, however, the three Yakut groups (YakutCK, YakutDY, and YakutYE) that show multimodal distributions may be considered stable populations or have experienced genetic bottlenecks in the past. It could also be the result of not getting an appropriate sample size, therein overlooking much of the variation within the group. In all cases, the raggedness index is quite low, near zero for all populations except for the Udehe. Raggedness indices near zero are interpreted as populations that are expanding, while large raggedness indices suggest demographic equilibrium and possible genetic drift. Raggedness values closely mimic conclusions drawn from Tajima's D results, suggesting that all populations are expanding except for the Udehe. Additionally, Yakut villages with multimodal mismatch distributions, suggesting demographic stability, had conflicting low raggedness values, which may mean that the demographic inferences made about these populations may not be accurate and more research should be conducted to better elucidate the evolutionary forces acting on these populations. Furthermore, increasing the sample size of these populations may allow for a better representation of the population and provide a better picture of what evolutionary forces are acting on the population.

Although the D-loop distributions exhibited curves that appeared much less (visually) smooth than the HVS1 results, which could be indicative of populations experiencing either: genetic drift or genetic bottlenecks in both the distant or recent past, the addition of HVS 2 and 3 sequence data caused lower raggedness indices for all populations--except for YakutCK and YakutDA. However, these villages still maintained relatively low raggedness values, which suggest demographic expansion. Consistent with HVS1 distribution results, the Udehe were the only population that had high raggedness index, demographic equilibrium, genetic drift, and increased homogeneity. Likewise, all of the distributions are unimodal except for the YakutCK and YakutYE, providing additional support for the idea that such

populations are numerically expanding. Furthermore, as was observed with the HVS1 distributions, all populations had auxiliary modes at zero mismatches, which could be the result of sampling error due to the overrepresentation of related individuals. Though larger sample sizes may help to give a better representation of what evolutionary forces are impacting these populations (Sherry et al., 1994), it appears that the addition of HVS 2 and 3 data does not significantly alter the information gained from mismatch distributions.

## **5.4 Phylogeography**

### **5.4.1 Mantel Randomization**

Mantel Randomization tests were used to evaluate the relationship between molecular variation and geographical distribution for all populations using HVS 1 and D-loop sequence data. A Mantel test compares two matrices, in this case a geographical distance and genetic distance matrix, resulting in a correlation coefficient between them. Genetic distance matrices were constructed for the four main populations, for all populations with the Yakut separated into home villages, and a randomly generated sample of all of the Yakut. The results of this analysis revealed moderate correlation between genetics and geography for both HVS 1 and D-loop; however, none of these results were statistically significant. The genetic distance matrix with the highest  $r$  value (0.58) included all populations with the Yakut separated into home village for HVS 1 data. This value was marginally non-significant at  $P=0.077$ . Although these findings are non-significant, they are consistent with the correlation found by Crawford et al. (2007), suggesting the existence of a strong relationship between genetics and geography ( $r=0.55$ ,  $p<0.000$ ) for Native Siberian populations. Mantel test results for D-loop data mirrored those of HVS 1. The highest correlation between molecular and geographic distances was found with the Yakut separated into respective villages ( $r = 0.42$ ,  $p=0.039$ ). Interestingly, though correlation decreased after the addition of the second and third hypervariable segment, a moderately strong correlation remained

between the two matrices, and the results became statistically significant. This suggests that for phylogeographic analysis utilizing mantel randomization tests, the addition of sequence data from the rest of the D-loop may not show any stronger correlated relationships between geography and genetics, but it can improve the statistical significance and therein the power of the test.

Mantel randomization tests were also used to test whether there was a significant correlation between HVS1 sequence data and D-loop sequence data. A significantly low correlation between the data sets would mean that either one of the two sets is markedly different and therein better able to characterize the genetic structure of a population than the other, or the difference in number of SNPs would heavily influence the outcome of all analyses. The correlation between HVS 1 and D-loop data was high and significant, all above  $r = 0.94$  ( $p < .005$ ). The intimate correlation between matrices for HVS 1 and D-loop data provides further weight to the analyses conducted in this study, suggesting that the additional SNPs attained from adding HVS 2 and 3 should not drastically alter results due to overt differences in matrix composition. Rather, one matrix is not inherently better than the other. However, even with this marked similarity, it appears that the addition of HVS 2 and 3 does increase the phylogenetic resolution of findings for some of the analyses conducted, significantly altering some of these tests.

## CHAPTER 6: CONCLUSION

This project analyzed mtDNA sequence data from the D-loop (HVS 1,2, and 3) of four Siberian populations: Altai, Evenki, Yakut, and Udehe, in order to gain insight on the utility and efficacy of sequencing the entire D-loop instead of only using the HVS 1 region when characterizing population data. By comparing these two data sets using multivariate statistical methods commonly employed in anthropological genetics, this project investigated: 1.) whether the increase in the number of SNPs sequenced revealed different phylogenetic relationships between Siberian populations; 2.) if additional genetic variation can be revealed by the addition of more genomic regions; and 3.) whether additional SNPs reveal stronger relationships between genetics, linguistics, and geography than using the HVS1 alone.

The results of the analyses conducted in this study were largely consistent with previously reported findings of other Siberian populations, maintaining similar molecular and genetic diversity indices, within and between population variation, effects of forces of evolution and relationships between geography, language and genes. When considering the mtDNA lineages present in the four populations of study, each group showed a high frequency of Native Siberian haplogroups with European and East Asian lineages at varying frequencies. It appears that the addition of the second and third hypervariable segment does in fact aid in further characterizing the populations. In three of the four (Altai, Evenki, and Yakut), all of the haplotypes that were not characterized by HVS1 alone were resolved by adding HVS 2 and 3. However, no tangible differences were reported between the two sets of data for gene and nucleotide diversity.

As often noted in Native Siberian populations, a high proportion of variation is explained within the various groups. In fact, in each grouping of populations (language, culture, or geography), the

amount of within group variation was always above 94.0%. This result exhibited no change by adding the HVS 2 and 3, and the statistical significance remained the same as well. This suggests that using AMOVA, no additional genetic variation is explained by analyzing other genomic regions of the control region of the mitochondrial genome.

One of the most interesting and unexpected results occurred between the two different data sets for multidimensional scaling. HVS1 data showed the Udehe being extremely divergent from the other of the populations. However, adding the additional loci from the HVS 2 and 3, it revealed that they were much more similar than originally thought. This was perhaps the clearest example of the utility of additional genomic information in the entire study.

The two different types of neighbor-joining trees utilized for this project revealed different conclusions on whether or not the entire D-loop gives better resolution of phylogenetic relationships between populations. When considering NJTs based on FSTs, no significant changes were noted between HVS1 and D-loop. However, visualizing the phylogenetic relationship between individuals from all populations, the addition of the extra genomic regions allow for a better resolution and higher statistical power, as the bootstrap values of many branches increased dramatically.

Various analyses were conducted to determine the effect of evolutionary forces that may be acting on populations in this study, such as Tajima's  $D$ , Fu's  $F_s$ , and Mismatch Analysis. The overall conclusion is that the majority of populations are undergoing expansion, whereas the Udehe are consistently shown to be in demographic equilibrium or genetic drift. This is not surprising as the Udehe are the smallest population present in this study, numbering just one-thousand individuals, and are the most geographically and culturally isolated from the rest of the Siberian populations, along with being heavily influenced by East Asian cultures. The Udehe sample for this study was small ( $n = 34$ ), and so although this sample was not characterized by any European mtDNA lineages, it is possible that they

do contain influence from Eurasians, as the Udehe territory is one of the final stops on the Trans-Siberian Railroad. Since its inception, there has been a major influx of Russians and Ukrainians in the region. Although few moderate differences were noted with statistical significance, neutrality tests and mismatch analysis indicated that adding the second and third hypervariable segments do not substantially alter results obtained from HVS1 alone.

Mantel randomization tests were conducted to help determine whether there was a relationship between geography and genetics among these populations. These results indicate that adding HVS 2 and 3 can significantly alter correlations. One of the most telling results of this study illustrated this effect when a non-significant correlation became significant after this addition, providing evidence that the entire D-loop provides considerable improvement in the resolution of results in population studies.

Overall conclusions drawn from this study suggest that there is an element of analytical “relativism.” Rather, it depends on the analyses employed in the study. For half of the analyses conducted: Haplogroup characterization, MDS, Neighbor Joining Tree construction with bootstrap replication based on haplotypic data, and Mantel randomizations, the additional SNPs from HVS 2 and 3 sequence data significantly impacted the results of the study, suggesting it is preferable to characterize populations based on the entire D-loop instead of HVS 1 alone. However, the addition of these two regions did not appear to significantly alter results obtained from diversity measures, AMOVA results, NJTs based on distance matrices, and examining the effects evolutionary forces on populations. However, in some cases raggedness indices for mismatch distributions did improve, this did not change the overall outcome. This indicates that sequencing the entire D-loop does not significantly improve these types of tests and it may be more beneficial to use resources to increase the sample size utilized to better represent the population of interest.

The use of mtDNA markers is an important tool for anthropological studies that can be used to characterize the distribution of genes within a population and gain insight into the biological consequences of historical and demographic contexts of specific populations. Though the financial burden of conducting molecular research is decreasing, it is still an expensive discipline. That said, the debate over whether or not one should analyze the entire mitochondrial D-Loop or simply the seemingly standard HVS 1 remains a valid query. Results from this study indicate that the answer is not as simple as one might expect. This work has shown that it is not necessarily as dichotomous as whether the D-loop is inherently better than HVS 1 at characterizing populations, as nearly half of the methodologies conducted in this study revealed that the use of the entire D-loop does not provide a higher degree of phylogenetic resolution compared to that of HVS 1 data. The decision over the quantity of SNPs to be analyzed is therefore dependent upon the research questions, the types of analyses conducted, and the sample sizes of the populations of study.

## Bibliography

- Alekseev AN (1996). Ancient Yakutia: The Iron Age and the Medieval Epoch (in Russian). Izdatel'stvo Instituta Arkheologii I Etnografii, Novosibirsk, Russia.
- Alexeev VP (1989). *Historical Anthropology and Ethnogenesis*. Nauka, Moscow.
- Alexeev VP, Gohman II (1984). *Anthropology of the Asiatic part of USSR*. Nauka, Moscow.
- Anderson S, Bankier AT, Barrell BG, De Bruijn MHL, Coulson AR, Drouin J, Eperon IC, Nierlich DP, Roe BA, Sanger F, Schreier PH, Smith AJ, Staden R, and Young IG (1981). Sequence and Organization of the Human Mitochondrial Genome. *Nature* 290: 457-467.
- Andrews RM, Kubacka I, Chinnery PF, Lightowlers RN, Turnbull DM and Howell N (1999). Reanalysis and Revision of the Cambridge Reference Sequence for Human Mitochondrial DNA. *Nature Genetics* 23: 147.
- Aris-Brosou S, and Excoffier L (1996). The Impact of Population Expansion and Mutation Rate Heterogeneity on DNA Sequence Polymorphism. *Molecular Biology and Evolution*. 13(3): 494-504.
- Armstrong T (1968). Farming on the Permafrost. *The Geographical Magazine*. March: 961-967.
- Avise JC (2000). *Phylogeography: The History and Formation of Species*. Cambridge, MA; Harvard University Press, London.
- Avise JC, Arnold J, Ball RM, Bermingham E, Lamb T, Neigel JE, Reeb CA, and Saunders NC (1987). Intraspecific Phylogeography: The Molecular Bridge between Population Genetics and Systematics. *Annual Review of Ecology and Systematics*. 18: 489-522.
- Ballinger SW, Schurr TG, Torroni A, Gan YY, Hodge JA, Hassan K, Chen K-H, and Wallace DC (1992). Southeast Asian Mitochondrial DNA Analysis Reveals Genetic Continuity of Ancient Mongoloid Migrations. *Genetics* 130: 139-152.
- Balzer MM (1994). Yakut. In: *Encyclopedia of World Cultures, eds Friedrich P, Diamond N, volume VI – Russia and Europe/China*. G.K. Hall & Co. New York, 404-407.
- Bandelt H-J, Herrnstadt C, Yao YG, Kong QP, Kivisild T, Rengo C, Scozzari R, Richards R, Richards M, Villems R, Macaulay V, Howell N, Torroni A, and Zhang YP (2003). Identification of Native American Founder mtDNAs through the Analysis of Complete mtDNA Sequences: Some Caveats.
- Barbujani G, Magagni A, Minch E, and Cavalli-Sforza L (1997). An Apportionment of Human DNA Diversity. *Proceedings of the National Academy of Sciences, USA* 90: 4670-4673.
- Bowles GT (1977). *The Peoples of Asia*. Weidenfeld & Nicholson, London.

- Brown M, Hosseini S, Torroni A, Bandelt H, Allen J, Schurr T, Scozzari R, Cruciani F, Wallace D (1998). mtDNA Haplogroup X: An Ancient Link between Europe/West Asian and North America? *American Journal of Human Genetics* 63: 1852-1861.
- Crawford MH (1973). The Use of Genetic Markers of the Blood in the Study of the Evolution of Human Populations. In *Methods and Theories of Anthropological Genetics*. Eds Crawford MH and Workman GL. University of New Mexico Press, Albuquerque. 19-38.
- Crawford MH (2007). Genetic Structure of Circumpolar Populations: A synthesis. *American Journal of Human Biology* 19:203-217.
- Crawford MH (2007b). *Anthropological Genetics; Theory, Methods and Applications* ed. Crawford MH. Cambridge University Press, New York.
- Crawford MH and Leonard WR (2002). The Biological Diversity of Herding Populations: An Introduction. In: *Human Biology of Pastoral Populations*, eds. Leonard WR and Crawford MH. Cambridge University Press, Cambridge. 1-9.
- Crawford MH and Workman PL eds. (1973) *Methods and Theories of Anthropological Genetics*. University of New Mexico Press, Albuquerque, NM.
- Crick FHC and Watson JD (1954). The Complementary Structure of Deoxyribonucleic Acid. *Proceedings of the Royal Society A: Mathematical Physical and Engineering Sciences*. 223(1152): 80-96.
- Derbeneva OA, Starikovskaya EB, Volodko NV, Wallace DC, and Sukernik RI (2002). Mitochondrial DNA Variation in Kets and Nganasans and the Early Peopling of Northern Eurasia. *Russian Journal of Genetics*. 38: 1554-1560.
- Derbeneva OA, Starkovskaya EB, Wallace DC and Sukernik RI (2002b). Traces of Early Eurasians in the Mansi of Northwest Siberia Revealed by Mitochondrial DNA Analysis. *American Journal of Human Genetics*. 70: 1009-1014.
- Derenko MV, Denisova GA, Malyarchuk BA, Dambueva IK, Luzina FA, Lotosh EA, Dorzhu CM, Karamchakova ON, Solovenchuk LL, and Zakharov IA (2001) The Structure of the Gene Pools of the Ethnic Populations of Altai-Sayan Region Based on Mitochondrial DNA Polymorphism Data. *Russian Journal of Genetics*. 37: 1177-1184.
- Derenko MV, Gryzbowski T, Malyarchuk BA, Dambueva IK, Denisova GA, Czarny J, Dorzhu CM, Kakpakov VT, Miścicka- Śliwka D, Woźniak M, and Zakharov IA (2003). Diversity of Mitochondrial DNA Lineages of South Siberia. *Annals of Human Genetics*. 67: 391-411.
- Derenko MV, Malyarchuk BA, Grzybowski T, Denisova G, Dambueva I, Perkova M, Dorzhu C, Luzina F, Lee HK, Vanecek T, Villems R, and Zakharov I (2007). Phylogenetic Analysis of Mitochondrial DNA in Northern Asian Populations. *American Journal of Human Genetics* 81: 1025-1041.
- Derenko MV and Shields GF (1997). Mitochondrial DNA Sequence Diversity in Three Northern Asian Aboriginal Population Groups. *Molecular Biology*. 31: 665-669.

- Derenko MV and Shields GF (1998). Polymorphism in Region V of Mitochondrial DNA in Indigenous Populations of Northern Asia. *Russian Journal of Genetics* 34: 553-557.
- Derenko MV and Shields GF (1998b). Variation of Mitochondrial DNA in Three Groups of Indigenous Populations of Northern Asia. *Russian Journal of Genetics*. 34: 321-324.
- Derevianko AP (1998). Human Occupation of Nearby Regions and the Role of Population Movements in the Paleolithic of Siberia. In: *The Paleolithic Siberia: New Discoveries and Interpretations*. Eds Derev'yanko AP. 336-351. University of Illinois Press. Urbana, IL.
- Efron B (1982). *The Jackknife, the Bootstrap and Other Resampling Plans*. Vol. 38. Society for Industrial and Applied Mathematics. Philadelphia, PA.
- Excoffier L, Laval G, and Schneider S (2005). Arlequin (version 3.0): An Integrated Software Package for Population Genetics Data Analysis. *Evolutionary Bioinformatics*. 1: 47-50.
- Excoffier L and Lischer H (2010). Arlequin Suite ver 3.5: A New Series of Programs to Perform Population Genetics Analyses Under Linux and Windows. *Molecular Ecology Resources*. 10;3: 564-567.
- Excoffier L, Smouse PE, and Quattro JM (1992). Analysis of Molecular Variance Inferred from Metric Distances among DNA Haplotypes: Application to Human Mitochondrial DNA Restriction Data. *Genetics* 131: 479-491.
- Fagundes NJR, Kanitz R, Bonatto SL (2008). A Reevaluation of the Native American mtDNA Genome Diversity and Its Bearing on the Models of Early Colonization of Beringia. *PLoS ONE* 3(9): e3157.
- Federova SA, Bermisheva MA, VILLEMS R, Maksimova NR, and Khusnutdinova EK (2003). Analysis of Mitochondrial DNA Lineages in Yakuts. *Molecular Biology*. 37: 544-554.
- Felsenstein J (1985). Confidence Limits on Phylogenies: An Approach Using the Bootstrap. *Evolution*. 783-791.
- Forster P, Harding R, Torroni A, Bandelt H-J (1996). Origin and Evolution of Native American mtDNA Variation: A Reappraisal. *American Journal of Human Genetics*. 59: 935-945.
- Forsyth J (1992). *A History of the Peoples of Siberia: Russia's North Asian Colony 1581-1990*. Cambridge University Press, Cambridge, England.
- Francalacci P, Montiel R, and malgosa A (1999). A Mitochondrial DNA Database. In *Genomic Diversity: Applications in Human Population Genetics*, eds. Papiha S, Deka R, and Chakraborty R. Kluwer Academic/Plenum Publishers, New York. 103-119.
- Fu XY (1997). Statistical Tests of Neutrality of Mutations against Population Growth, Hitchhiking, and Background Selection. *Genetics*. 147: 915-925.

- Han KC. (2005). Parhae's Succession to Koguryo Based on Inhabitatns. In: ed Kim CB, *Pharae's Sucession to Koguryo, Vol. 7 Seoul: The Koguryo Research Foundation Press.* 103-191.
- Helgason A, Nicholson G, Stefánsson K, and Donnelly P (2003). A Reassessment of Genetic Diversity in Icelanders: Strong Evidence from Multiple Loci for Relative Homogeneity Caused by Genetic Drift. *Annals of Human Genetics.* 67: 281-297.
- Herrnstadt C, Elson JL, Fahy E, Preston G, Turnbull DM, Anderson C, Ghosh SS, Olefsky JM, Beal MF, Davis RE, and Howell N (2002). Reduced-Median-Network analysis of Complete Mitochondrial DNA Coding-Region Sequences for the Major African, Asian, and European Haplogroups. *American Journal of Human Genetics.* 70: 1152-1171.
- Jin HJ, Kim KC, and Kim Wook (2010). Genetic Diversity of Two Haploid Markers in the Udegey Population from Southeastern Siberia. *American Journal of Physical Anthropology* 142: 303-313.
- Jin HJ, Kwak KD, Hammer MF, Nakahori Y, Shinka T, Lee JW, Jin F, Jia X, Tyler-Smith C, and Kim W (2003). Y-Chromoosomal DNA Haplogroups and Their Implications for the Dual Origins of the Koreans. *Human Genetics* 114: 27-35.
- Jin HJ, Tyler-Smith C, and Kim W (2009). The Peopling of Korea Revealed by Analyses of Mitochondrial DNA and Y-Chromosomal Markers. *PLoS ONE* 4: e4210.
- Jobling MA, Hurles ME, and Tyler-Smith C (2004). *Human Evolutionary Genetics: Origins, Peoples and Disease.* Garland Science, New York.
- Jordan BB, Jordan-Bychkov TG and Holz RK (2001). *Siberian Village: Land and Life in the Sakha Republic.*
- Kimura M (1980). A Simple method for Estimating Evolutionary Rate of Base Substitutions through Comparative Studies of Nucleotide Sequences. *Journal of Molecular Evolution.* 16: 111-120.
- Kolman CJ, Sambuughin N, and Bermingham E (1996). Mitochondrial DNA Analysis of Mongolian Popualtions and Implications for the Origin of New World Founders. *Genetics* 142: 1321-1334.
- Kong Q-P, Yao Y-G, Sun C, Bandelt H-J, Zhu C-L and Zhang Y-P (2003). Phylogeny of East Asian Mitochondrial DNA Lineages Inferred from Complete Sequences. *American Journal of Human Genetics* 73: 671-676.
- Kostantinov IV (1975). The Origins of the Yakut People and Their Culture (in Russian). *Publications of the Prilenskaya Archaeological Expedition.* Yakutskiy Filial SO AN SSSR, Yakutsk. 106-173.
- Krauss ME (1988) Many Tongues: Ancient Tales. In: *Crossroads of Continents.* Eds Crowell WW. Smithsonian Institution, Washington D.C. 145-150.
- Kruskal JB (1964). Multidimensional Scaling by Optimizing Goodness of Fit to a Nonmetric Hypothesis. *Psychometrika* 29: 1-27.
- Kruskal JB (1964b). Nonmetric Multidimensional Scaling: A Numeric Method. *Psychometrika* 29: 28-42.

- Levin MG and Potapov LP (1964). *The Peoples of Siberia*. Translated from Russian by Stephen Dunn. University of Chicago Press, Chicago, IL
- Luick J (1978). *Reindeer, Horse, and Yak Production in Yakutia USSR*. Department of the Interior Bureau of Indian Affairs, Juneau, AK.
- Meyer S, Weiss G, Von Haeseler A (1999). Pattern of Nucleotide Substitution and Rate Heterogeneity in the Hypervariable Regions I and II of Human mtDNA. *Genetics* 152: 1103-1110.
- Miller RA (1991). Genetic Connections among the Altaic Languages. In *Sprung from Some Common Source: Investigations into the Prehistory of Languages*, eds Lamb SM, and Mitchell ED, 39-327 Stanford University Press, Stanford.
- Mullis KB, Faloona FA, Scharf SJ, Saiki R, Horn G, and Erlich H (1985). Specific Enzymatic Amplification of DNA In Vitro: The Polymerase Chain Reaction. *Cold Spring Harbor Symposia on Quantitative Biology, Volume LI*.
- Nei M (1987). *Molecular Evolutionary Genetics*. Columbia University Press, New York, USA.
- Nei M and Kumar S (2000). *Molecular Evolution and Phylogenetics*. Oxford University Press, USA.
- Nei M, and Li WH (1979). Mathematical Model for Studying Genetic Variation in Terms of Restriction Endonucleases. *Proceedings of the National Academy of Sciences, USA*. 76: 5269-5373.
- Nicholson G, Smith AV, Jónsson F, Gústafsson Ó, Stefánsson K, and Donnelly P (2002). Assessing Population Differentiation and Isolation from Single-Nucleotide Polymorphism Data. *Journal of the Royal Statistical Society: Series B* 64:695-715.
- Non AL, Kitchen A, and Mulligan CJ (2006). Identification of the Most Informative Regions of the Mitochondrial Genome for the Phylogenetic and Coalescent Analysis. *Molecular Phylogenetics and Evolution* 44: 1164-1171.
- Oaks J and Riewe R (1998) *Spirit of Siberia: Traditional Native Life, Clothing, and Footwear*. Smithsonian Institution Press, Washington D.C.
- Okladnikov AP (1955). The History of the Yakut ASSR. Volume 1: Yakutia Before Its Incorporation into the Russian State. Izdatel'stvo Akademii Nauk SSSR, Moscow, Russia.
- Okladnikov AP (1964). Ancient Population of Siberia and Its Culture. In: *The Peoples of Siberia*, eds Levin MG and Potapov LP 13-98. University of Chicago Press, Chicago, IL.
- Okladnikov AP (1970). Yakutia: Before Its Incorporation into the Russian State. McGill-Queen's University Press, Montreal, Canada.
- Pakendorf B, Spitsyn VA, and Rodewald A (1999). Genetic Structure of the Sakha Population from Siberia and Ethnic Affinities. *Human Biology*. 71: 231-244.

- Pakendorf B, Wiebe V, Tarskaia LA, Spitsyn VA, Soodyall H, Rodewald A, and Stoneking M (2003). Mitochondrial DNA Evidence for Admixed Origins of Central Siberian Populations. *American Journal of Physical Anthropology*. 120: 211-224.
- Petrishchev VN, Jutueva AB, and Rychkov IG (1993). Deletion-Insertion Polymorphism in the V-Region of Mitochondrial DNA in Ten Mongoloid Populations of Siberia, Frequency of Deletion Correlates with Geographic Coordinates of the Locality. *Genetica*. 29: 1196-1204.
- Petrova T (1967). *Language of Oroks (Ul'ta)*. Nauka, Moscow.
- Phillips-Krawczak C, Zlojutro M, Devor E, Crawford MH, and Wilson K (2006). mtDNA Variation in the Kizhi Population of Gorno Altai: A Comparative Study. *Human Biology*
- Puzyrev VP, Stepanov VA, Golubenko MV, Puzyrev KV, Maximova NR, Kharkov NV, Spiridonova MG, and Gogvitsina AN (2003). mtDNA and Y-Chromosomal Lineages in the Yakut Population. *Russian Journal of Genetics*. 39: 819-822.
- Reidla M, Kivisild T, Metspalu E, Kaldma K, Tambets K, Tolk H-V, Parik J, Loogv äli E-L, Derenko M, Malyarchuk B, Bermisheva M, Zhadanov S, Pennarun E, Gubina M, Golubenko M, Damba L, Fedorova S, Gusar V, Grechanina E, Mikerezi I, Moisan J-P, Chaventré A, Khusnutdinova E, Osipova L, Stepanov V, Voevoda M, Achilli A, Rengo C, Rickards O, De Stefano GF, Papiha S, Beckman L, Janicijevic B, Rudan P, Anagnou N, Michalodimitrakis E, Koziel S, Usanga E, Geberhiwot T, Herrnstadt C, Howell N, Torroni A, and VILLEMS R (2003). Origin and Diffusion of mtDNA Haplogroup X. *American Journal of Human Genetics*. 73: 1178-1190.
- Relethford, JH (2003). *Reflections of Our Past; How Human History is Revealed in Our Genes*. Westview Press.
- Relethford, JH (2012). *Human Population Genetics*. Wiley-Blackwell. Hoboken, NJ.
- Rolf FJ (2000). NTSYSpc, version 2.1 Exeter Publishing, Ltd. Setauket, NY.
- Rolf FJ (2008). NTSYSpc: Numerical Taxonomy System, version 2.2 Exeter Publishing, Ltd. Setauket, NY.
- Rolf FJ and Sokal RR (1981). Comparing Numerical Taxonomic Studies. *Systematic Biology* 30(4): 459-490.
- Rogers AR and Harpending HC (1992). Population Growth Makes Waves in the Distribution of Pairwise Genetic Differences. *Molecular Biology and Evolution* 9: 552-569.hudson
- Rubicz RC (2001). *Origins of the Aleuts: Molecular Perspectives*. M.A. Thesis, Department of Anthropology, University of Kansas.
- Rubicz RC, Schurr TG, Babb PL, and Crawford MH (2003). Mitochondrial DNA Variation and the Origins of the Aleuts. *Human Biology*. 75: 809-835.

- Saillard J, Forster P, Lynnerup N, Bandelt HJ, and Norby S (2000). mtDNA Variation among Greenland Eskimos: The Edge of the Beringina Expansion. *American Journal of Human Genetics*. 67: 718-726.
- Saitou N and Nei M (1987). The Neighbor-Joining Method: A New Method for Reconstructing Phylogenetic Trees. *Molecular Biology and Evolution* 10: 471-483.
- Sambuughin N, Petrishchev VN, and Rychkov YG (1991). DNA Polymorphisms in a Mongolian Population: Analysis of Restriction Endonuclease Polymorphism of Mitochondrial DNA. *Genetika*. 27: 2143-2151.
- Schönig C (1990). Classification Problems of Yakut in Dor R (ed) *L'Asie Centrale et ses Voisins*. INALCO, Paris, 91-102.
- Schurr T and Wallace D (1999). Mitochondrial DNA Variation in Native Americans and Siberians and Its Implications for the Peopling of the New World. In: *Who Were the First Americans: Proceedings of the 58<sup>th</sup> Annual Biological Colloquium, Oregon State University*. Editor: Bonnicksen R. centre for the Study of the First Americans Corvallis, OR 41-77.
- Sherry ST, Rogers AR, Harpending H, Soodyall H, Jenkins T, and Stoneking M (1994). Mismatch Distributions of mtDNA Reveal Recent Human Population Expansions. *Human Biology*. 66: 761-775.
- Shields GF, Hecker K, Voevoda MI, and Reed JK (1992). Absence of the Asian-Specific Region Mitochondrial marker in Native Beringians. *American Journal of Human Genetics*. 50: 758-765.
- Shields GF, Schmiechen AM, Frazier BL, Redd A, Voevoda MI, Reed JK and Ward RH (1993). mtDNA Sequences Suggest a Recent Evolutionary Divergence for Beringian and Northern North American Populations. *American Journal of Human Genetics*. 53: 549-562.
- Slatkin M and Hudson RR (1991). Pairwise Comparisons of Mitochondrial DNA Sequences in Stable and Exponentially Growing Populations. *Genetics* 129(2): 555-562.
- Smolyak AK, (1975). Ethnic Processes in the Populations of Lower Amur Region and Sakhalin. Nauka, Moscow.
- Smouse PE, Long JC, and Sokal RR (1986). Multiple Regression and Correlation Extensions of the Mantel Test of matrix Correspondence. *Systematic Zoology*. 35: 627-632.
- Starikovskaya YB, Sukernik RI, Schurr TG, Kogelnik AM, and Wallace DC (1998). mtDNA Diversity in Chukchi and Siberian Eskimos: Implications for the Genetic History of Ancient Beringia and the Peopling of the New World. *American Journal of Human Genetics*. 63: 1473-1491.
- Starikovskaya YB, Sukernik RI, Derbeneva OA, Volodko NV, Ruiz-Pesini E, Torroni A, Brown MD, Lott MT, Hosseini SH, Huoponen K, and Wallace DC (2005). Mitochondrial DNA Diversity in Indigenous Populations of the Southern Extent of Siberia and Origins of Native American Haplogroups. *Annals of Human Genetics* 69: 67-89.

- Stone AC and Stoneking M (1999). Analysis of Ancient DNA from a Prehistoric Amerindian Cemetery. *Philosophical Transactions of the Royal Society of London. Series B: Biological Sciences.* 354: 153-159.
- Stoneking M (2000). Hypervariable Sites in the mtDNA Control Region Are Mutational Hotspots. *American Journal of Human Genetics.* 67: 1029-1032.
- Stoneking M, and Soodyall H (1996). Human Evolution and the Mitochondrial Genome. *Current Opinion in Genetics and Development.* 731-736.
- Sturrock K and Rocha J (2000). A Multidimensional Scaling Stress Evaluation Table. *Field Methods* 12(1): 49-60.
- Tajima F (1983). Evolutionary Relationship of DNA Sequences in Finite Populations. *Genetics.* 123: 585-595.
- Tajima F (1989). Statistical Method for Testing the Neutral Mutation Hypothesis DNA Polymorphism. *Genetics.* 123: 585-595.
- Tamura K, Dudley J, Nei M, and Kumar S (2007). MEGA4: Molecular Evolutionary Genetics Analysis (MEGA) software version 4.0 *Molecular Biology and Evolution* 24: 1596-1599.
- Tamura K, Peterson D, Peterson N, Stecher G, Masatoshi N, and Kumar S (2011). MEGA5: Analysis Using Maximum Likelihood, Evolutionary Distance and Maximum Parsimony. *Molecular Biology and evolution.* 28(10): 2731-2739.
- Tarskaia LA, Bychkovskaia LS, Pai GV, Makarov SB, Pakendorf B, Spitsyn VA (2002). Distribution of the ABO Blood Groups and the HP, TF, GC, PI, and C3 Serum Proteins in Yakuts. *Russian Journal of Genetics* 38: 548-553.
- Tarskaia LA, Makarov SB, Bychkovskaia LS, Pai GV, Pakendorf B, Elchinova GI, Deriabin BE, and Spitsyn VA (2002b). Ethnogenetics of Yakuts from the Three Regions of the Republic of Sakha (Yakutia) Inferred from Frequencies of Biochemical Gene Markers. *Russian Journal of Genetics* 38: 1088-1097.
- Tarskaia LA, Bychkovskaya LS, Pai GV, Makarov SV, Pakendorf B, and Spitsyn VA (2002c). Genetic Polymorphism of Erythrocytic Enzymes in Yakut Populations. *Russian Journal of Genetics* 38: 426-429.
- Tarskaia L, Gray RR, Burkley B, and Mulligan CJ (2006). Genetic Variation at the Mitochondrial DNA 9-bp Repeat Locus in the Sakha of Siberia. *Human Biology.* 78(2): 179-198.
- Tokarev SA and Gurvich IS (1956). The Yakuts (in Russian) in Levin MG and Potapov LP (eds) *Narody Sibiri.* Russian Academy of Science, Moscow. 267-328.
- Torrioni A, Schurr TG, Cabell MF, Brown MD, Neel JV, Larsen M, Smith DG, Vullo CM, and Wallace DC (1993a). Asian Affinities and Continental Radiation of the Four Founding Native American mtDNAs. *American Journal of Human Genetics.* 53: 563-590.

- Torrioni A, Sukernik RI, Schurr TG, Starikovskaya YB, Cabel MF, Crawford MH, Commuzie AG, and Wallace DC (1993b). mtDNA Variation of Aborigininal Siberians Reveals Distinct Genetic Affinities with Native Americans. *American Journal of Human Genetics*. 53: 591-608.
- Vasiliev SA (1993). The Upper Palaeolithic of Northern Asia. *Current Anthropology*. 34: 82-92.
- Vasilevich GM, and Smolyak AV (1964). The Evenks. In: *The Peoples of Siberia*, eds Levin MG and Ptotapov LP. University of Chicago Press, Chicago, IL.
- Wallace DR and Lott MT (2004). MITOMAP: A Human Mitochondrial Database. <http://www.mitomap.org>. Accessed March 2013.
- Ward RH, Frazier BL, Dew-Jager K, and Pääbo S (1991). Extensive Mitochondrial Diversity within a Single Amerindian Tribe. *Proceedings of the National Academy of Science*. 88: 8720-8724.
- Young KL (2009). *The Basques in the Genetic Landscape of Europe*. PhD Dissertation. University of Kansas.
- Zlojutro M, Rubicz R, Devor EJ, Spitsyn VA, Makarov SV, Wilson K, and Crawford MH (2006). Genetic Structure of the Aleuts and Circumpolar Populations Based on Mitochondrial DNA Sequences: A Synthesis. *American Journal of Physical Anthropology*. 129: 446-464.
- Zlojutro M, Tarskaia LA, Sorensen M, Snodgrass JJ, Leonard WR, and Crawford MH (2009). Coalescent Simulations of Yakut mtDNA Variation Suggest Small Founding Population. *American Journal of Physical Anthropology* 139: 474-482.









































