# Finding the best subset of SNPs for distinguishing populations of Yukon River Chinook salmon

## William D. Templin and Anton B. Antonovich
### Division of Commercial Fisheries, Alaska Department of Fish and Game, Anchorage, Alaska, USA

## Introduction

Recent advances in laboratory methods have increased the number of genetic markers available for identifying stock components in mixtures. Although the cost and laboratory time to analyze each marker has decreased, overall costs have increased due to running ever larger numbers of markers. Prior to the availability of a large number of markers most researchers analyzed all available markers to maximize precision. Increasingly, selecting the most informative markers for specific applications will be critical to containing costs while maximizing the discrimination of key stock components. Here we explored three methods of choosing an optimal set of single nucleotide polymorphism (SNP) loci from the 24 SNPs available in Yukon River Chinook salmon populations. By ranking each locus using 1) mean interpopulation allelic frequency differences (delta), 2) mean interpopulation Fst, and 3) summed loadings for each locus from Principal Components Analysis (PCA), we developed sets of informative loci that were incrementally tested for precision and accuracy in simulated mixed stock analysis.
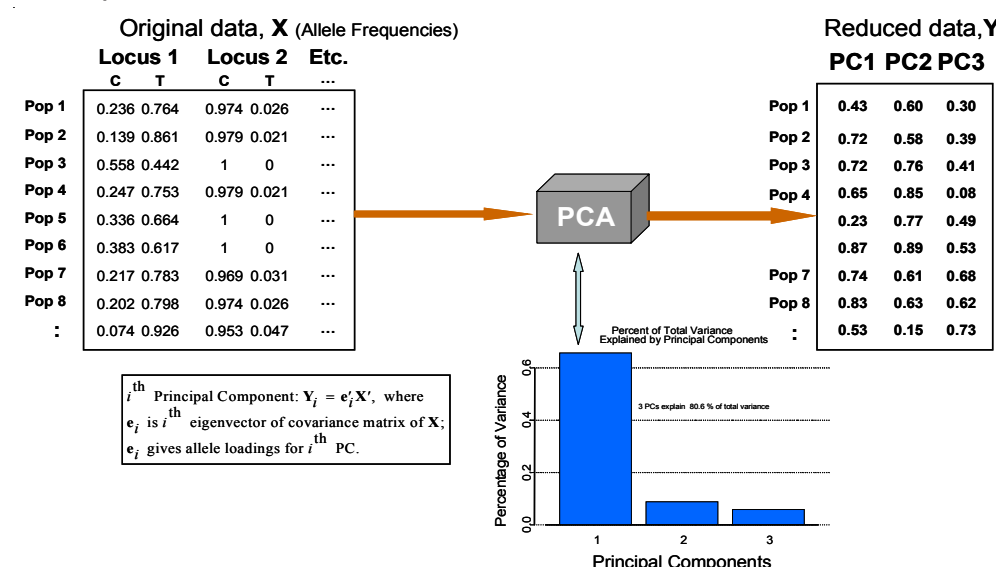
## Methods

### Ranking Markers -

*delta* - The *delta* statistic measures the genetic distance between population pairs as the sum of the absolute differences between allele frequencies. Markers were ranked by the mean of the interpopulation *delta* values calculated for that marker.

$$\delta_{AB} = \frac{1}{2}\sum_{i=1}^{L}\left|p_i^A - p_i^B\right|$$

*Fst* - The Fst statistic is a measure of genetic diversity based on partitioning the variance of allele frequencies within and among populations in a "weighted" ANOVA. Markers were ranked by the mean of the interpopulation Fst values calculated for that marker.
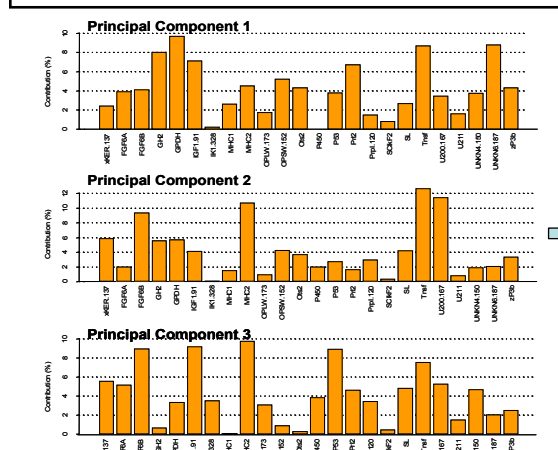
*PCA (Principal Components Analysis)* - PCA is used as a data reduction method that seeks to explain the variation in data (allele frequencies) with fewer parameters. We adapt this method to provide information about which markers are more closely associated with the variation in allele frequency within the data set. A brief description of the method follows:

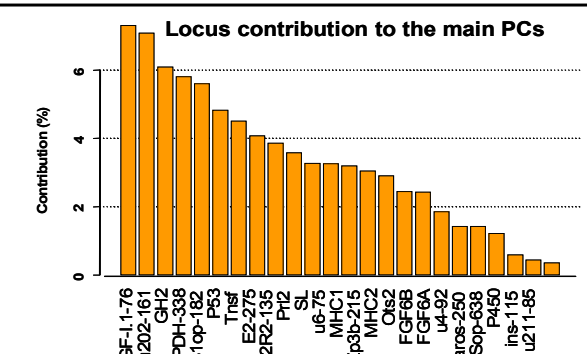Part 1: Find Principal Components that account for >80% of variation.



Part 2: Determine each marker contribution to each Principal Component and rank by average.



### Testing Marker Sets -

Sets of loci identified by the above methods were tested for usefulness for estimating relative contributions to mixed stock fish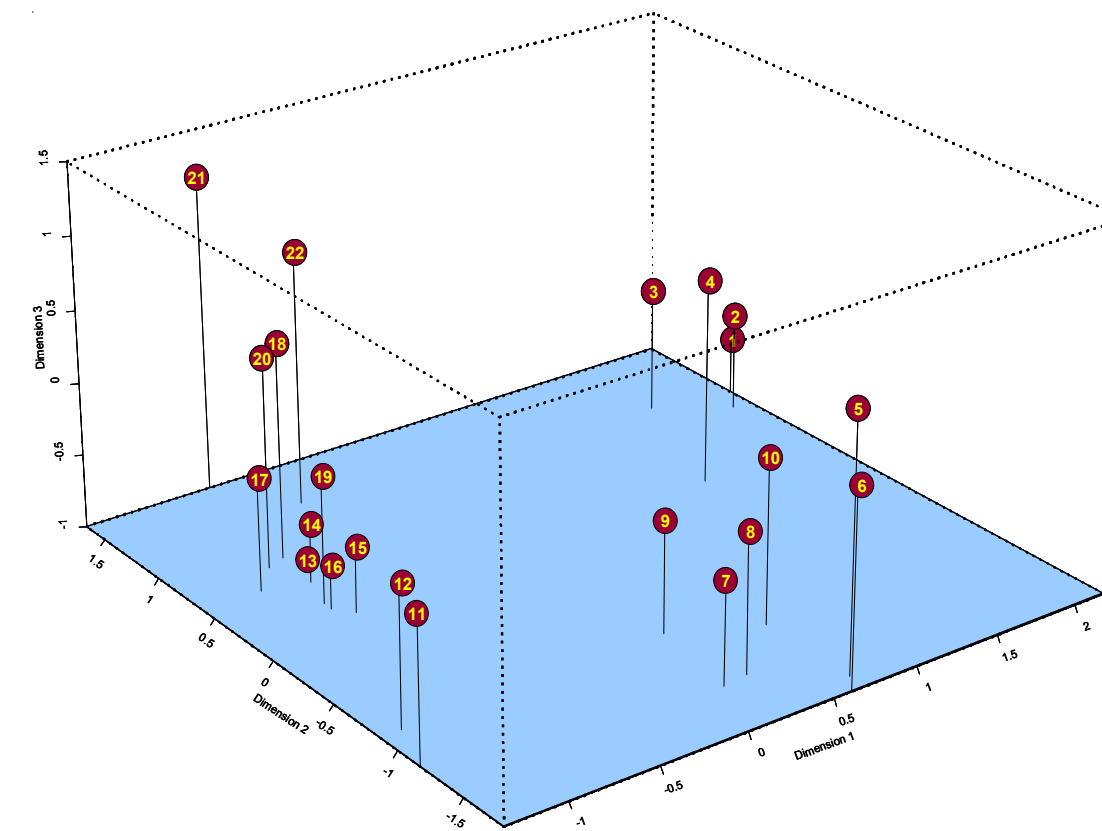eries in the Yukon River. Using simulations in which the relative contribution of all stocks in the baseline were estimated for simulated mixtures composed entirely of fish from a single population (1000 iterations), we measured the accuracy and precision of population composition estimates based on the reduced sets of markers. The performance of the selected sets of markers were compared with the perfomance of sets of randomly selected markers.

## Yukon River drainage and Chinook salmon populations



## Population Structure

Using the delta distances calculated between populations and plotting these distances in three dimensions (multidimensional scaling analysis) we can display relationships among populations. The dots and numbers match the populations on the map and in the legend. Populations can be segregated into groups based on genetic and geographic factors.



## Results: Locus Ranks

The three different methods of ranking the SNP markers did not sort the markers identically, but in general the "best" markers appeared near the top in each ranking. Twelve markers are common to all sets of the top fifteen ranked markers.

Spearman's rank correlation test indicated that the *delta* and the Fst ranks were more similar to each other (r=0.94; 90%CI [0.53-1.00]) than either was to the PCA ranking (r=0.87; 90%CI [0.46-1.00] and r=0.77; 90%CI [0.36-1.00], repectively).

| SNP | Ranks | | |
|---|---|---|---|
| Locus | delta | Fst | PCA |
| Ots_GH2 | 1 | 1 | 3 |
| Ots_IGF-i.1-76 | 2 | 2 | 1 |
| Ots_GPDH-338 | 3 | 4 | 4 |
| Ots_P53 | 4 | 8 | 6 |
| Ots_E2-275 | 5 | 5 | 8 |
| Ots_Tnsf | 6 | 3 | 7 |
| Ots_SClkF2R2-135 | 7 | 6 | 9 |
| Ots_Prl2 | 8 | 9 | 10 |
| Ots_FGF6A | 9 | 11 | 18 |
| Ots_MHC1 | 10 | 12 | 13 |
| Ots_u6-75 | 11 | 15 | 12 |
| Ots_FGF6B | 12 | 14 | 17 |
| Ots_SWS1op-182 | 13 | 7 | 5 |
| Ots_Zp3b-215 | 14 | 10 | 14 |
| Ots_SL | 15 | 16 | 11 |
| Ots_P450 | 16 | 20 | 22 |
| Ots_u202-161 | 17 | 13 | 2 |
| Ots_Ikaros-250 | 18 | 18 | 20 |
| Ots_Ots2 | 19 | 17 | 16 |
| Ots_MHC2 | 20 | 19 | 15 |
| Ots_LWSop-638 | 21 | 22 | 21 |
| Ots_u4-92 | 22 | 21 | 19 |
| Ots_ins-115 | 23 | 23 | 23 |
| Ots_u211-85 | 24 | 24 | 24 |

## Results: Stock Identification

The ranked sets of SNP markers show a rapid increase in mean population identification up to the top nine ranked markers, after which the addition of more markers demonstrate only small improvement (**Top Graph**). Perfect identification is 100% correctly assigned to the contributing group. All sets of "best" markers outperformed the randomly chosen sets. The coefficient of variation of the mean estimate dropped very rapidly from one to six marker sets (**Bottom Graph**). All three sets of ranked markers performed similarly and each was more accurate than randomly chosen sets of markers. The variation in the estimates with sets of three "best" markers was the same as with random sets, but precision improved more rapidly for chosen sets than for random sets as more markers were added.



## Conclusions

- The three methods of ranking the SNP markers by information content provided similar ranks.
- The ranked sets of SNP markers performed more accurately and with better precision than randomly chosen sets of markers for mixed stock analysis.
- Only a relatively small set of SNP markers (24) was available. This process may show greater differences between the ranking methods and improved mixed stock analysis performance when more SNP markers are available for analysis.

## Acknowledgements