

**Title:** Thermal Mark Recovery Data Quality Assurance and Quality Control Procedures by the ADF&G Mark, Tag and Age Laboratory  
**Authors:** Agler, B., L. Wilson, and M. Lovejoy  
**Date:** April 29, 2016

**Version:** 1.0

1 **Abstract**

2 Origin of Pacific salmon (*Oncorhynchus* spp.) sampled for the Alaska Hatchery Research  
3 Program can be determined by examining otoliths (ear stones) for thermal marks. Thermal mark  
4 presence indicates that a fish originated from a hatchery; whereas, thermal mark absence  
5 indicates wild origin. Identification of such marks provides information about a fish's age,  
6 hatchery of origin, and release location. The Mark, Tag and Age Lab, Alaska Department of Fish  
7 and Game is responsible for conducting mark recovery operations for a variety of statewide  
8 management and research projects. Thermal-marked fish typically are not given a secondary  
9 mark, so multiple readings among readers and across geographic areas are used to estimate  
10 reader ability to detect a thermal mark and to calculate agreement of thermal mark  
11 identifications. Thus, we compare first and second reads with an agreement matrix to determine  
12 whether there are any significant problems in reader training or challenging marks that might be  
13 re-examined. We then use the *kappa* statistic to examine overall agreement between readers as  
14 well as agreement by specific thermal mark. At the end of each project, we estimate the error  
15 rates of each reader using latent class models, because although useful, *kappa* statistics are  
16 influenced by the true proportion of marked fish. Analyzing the thermal mark read results in this  
17 manner provides a method to ensure quality control among projects and a measure of accuracy of  
18 thermal mark recoveries of fish sampled for the Alaska Hatchery Research Program.

19 **Background of AHRP**

20 Extensive ocean-ranching salmon aquaculture is practiced in Alaska by private non-profit  
21 corporations (PNP) to enhance common property fisheries. Most of the approximately 1.7B  
22 juvenile salmon that PNP hatcheries release annually are pink salmon in Prince William Sound  
23 (PWS) and chum salmon in Southeast Alaska (SEAK; Vercessi 2014). The large scale of these  
24 hatchery programs has raised concerns among some that hatchery fish may have a detrimental  
25 impact on the productivity and sustainability of natural stocks. Others maintain that the potential  
26 for positive effects exists. To address these concerns ADF&G convened a Science Panel for the

---

<sup>1</sup> This document serves as a record of communication between the Alaska Department of Fish and Game Commercial Fisheries Division and other members of the Science Panel of the Alaska Hatchery Research Program. As such, these documents serve diverse ad hoc information purposes and may contain basic, uninterpreted data. The contents of this document have not been subjected to review and should not be cited or distributed without the permission of the authors or the Commercial Fisheries Division

27 Alaska Hatchery Research Program (AHRP) whose members have broad experience in salmon  
28 enhancement, management, and natural and hatchery fish interactions. The AHRP was tasked  
29 with answering three priority questions:

- 30 I. *What is the genetic stock structure of pink and chum salmon in each region (PWS and*  
31 *SEAK)?;*
- 32 II. *What is the extent and annual variability in straying of hatchery pink salmon in PWS and*  
33 *chum salmon in PWS and SEAK?;* and
- 34 III. *What is the impact on fitness (productivity) of natural pink and chum salmon stocks due*  
35 *to straying of hatchery pink and chum salmon?*

## 36 **Introduction**

37 An important consideration in fisheries management is the ability to identify the origins of  
38 captured and harvested fish. The development of mass-marking techniques, such as thermal  
39 manipulation of water temperature to mark otoliths, permits millions of hatchery-incubated  
40 juvenile salmon to be marked simultaneously. These techniques have successfully applied  
41 species-specific, thermal mark patterns to otoliths (ear stones) of hatchery-reared salmon  
42 throughout Alaska and the Pacific Rim over the past 26 years (Hagen et al. 1995; Volk et al.  
43 1990). For the AHRP, accurate mark interpretation is vital to the assessment of stray rates  
44 associated with hatchery-reared salmon and provides validation of genetic stock identifications.

45 There are many potential sources of error in any research project, and the extent that these errors  
46 can be minimized increases confidence in a study's findings and conclusions. There are two  
47 categories of reliability with respect to data collectors: reliability across multiple data collectors,  
48 or *inter-rater* reliability, and reliability of a single data collector, or *intra-rater* reliability.  
49 Presented with the same situation and phenomenon every time, the assumption is that a  
50 laboratory staff would react the same way every time; however, Gwet (2014) provided examples  
51 of where this was false and affected intra-rater reliability. Reader reliability is affected by the  
52 fineness of discriminations required by the samples. If a variable only has two possible states,  
53 and the states are sharply differentiated, reliability is likely to be high. For example, if the  
54 outcome variable is that a fish either survived or did not, or the otolith is marked or not marked,  
55 there is likely to be high reliability in data comparisons between readers. On the other hand, if  
56 readers are required to make judgements or determinations regarding the width and amount of  
57 thermal mark rings, both inter- and intra-rater reliability declines. Careful training of laboratory  
58 staff is critical to reader reliability.

59 To determine the presence or absence of a thermal mark in an otolith, laboratory staff use pattern  
60 recognition and image matching (Blick and Hagen 1998). Although hatcheries follow strict  
61 rearing protocols to produce consistent thermal marks, natural variation in otolith development  
62 and growth patterns can obscure these patterns and interfere with the ability to detect a mark,  
63 reducing mark identification. In addition, stress on fish at the hatchery caused by temperature

64 fluctuations, water quality, rearing density, noise and light fluctuations, lot size, maintenance  
65 procedures, and handling protocols can affect mark consistency and clarity (Hagen et al. 1995).

66 Examination of accuracy rates of thermal mark identifications, including correct assignment of  
67 age, hatchery, and release site, provide useful knowledge regarding reliability of mark recoveries  
68 to assess stray rates and validate genetic stock identifications. However, otolith thermal marks  
69 are typically applied without a secondary mark, such as a coded-wire tag or a passive integrated  
70 transponder (PIT) tag, thus there is no reliable method to assess the true accuracy of thermal  
71 mark presence and identification. Consequently, the Mark, Tag, and Age (MTA) Lab uses latent  
72 class models (LCMs) to estimate a reader's ability to distinguish between hatchery and wild fish.  
73 In addition, *kappa* statistics (Cohen 1960) are used to assess reader agreement among individual  
74 mark patterns. Agreement matrices combined with the *kappa* statistic assist in identifying  
75 problematic mark patterns.

#### 76 *Goal*

77 Our goal is to describe the methods used by the MTA Lab in Juneau, Alaska to find errors in  
78 thermal mark classification and correct them. We describe the methods used to assess the  
79 accuracy of reader's ability to correctly ascertain presence and absence of a thermal mark and to  
80 identify specific mark patterns.

#### 81 **Methods**

82 Hatcheries apply thermal marks to incubating salmon eggs and fry by raising and lowering the  
83 water temperature at set intervals. Cycling temperature, or thermal marking, leaves patterns of  
84 optically-dense rings in the otolith (Volk et al. 1990). A thermal mark consists of rings, which  
85 are optically dark circles visible in the otolith and bands, which consist of one or more rings  
86 separated by a space from other rings (Figure. 1). We describe the thermal mark with a  
87 specialized notation termed the "hatch code" (Josephson et al. 2006). For example, a hatch code  
88 of 5,1,2H describes a set of three bands: the first band is composed of five rings, the second band  
89 includes one ring, and the last band contains two rings. The capital "H" indicates the mark was  
90 applied before hatching. In this example, all three bands occur prior to the hatch mark (Figure  
91 1). Varying the number and spacing of the induced rings produces unique patterns used to  
92 distinguish among similarly treated hatchery fish and wild stock (Hagen et al. 1995).

#### 93 *Thermal mark reference collection*

94 Initially, laboratory personnel are trained to dissect, prepare, and process otoliths from reference  
95 specimens or representative samples of salmon eggs, fry, and smolt obtained from the hatchery  
96 and preserved in alcohol before release. Upon receipt, five otoliths from each sample are  
97 dissected, mounted to glass slides (see "Otolith mounting" section of AHRP MTA Processing  
98 Tech Doc 7), and examined with a compound microscope. These reference specimens become  
99 the standard or authoritative mark pattern for that thermal mark after laboratory staff compare the  
100 observed marked with the assigned mark. Laboratory staff then measures the specimens,

101 because mark locations and ring spacing can vary among individuals of the same thermal mark  
102 group due to variability of fry developmental stage during marking. When a thermal mark is  
103 applied during different developmental stages, the distance from the core of the otolith to the  
104 initial band varies among fish making the mark challenging to identify. Successful thermal mark  
105 application at hatcheries is the first step to correct determination of fish origin and is  
106 fundamental to the success of this project.

107 Thermal marks are described in the Mark Characteristic Report (available online – see below).  
108 This report includes: brood year, release year, thermal mark identification, species, brood stock,  
109 release site(s), assigned (target) mark, actual mark observed, mark quality assessments, number  
110 of samples received, measurements of the otolith, information about temperature profile (if  
111 available), various comments, and an authoritative image of the mark as well as images of any  
112 variants of the mark. Measurements include minimum, maximum, and average distance ( $\mu\text{m}$ )  
113 from the core of the otolith to the first band; minimum, maximum and average width of each  
114 band; and distance among bands (Figure 2). The authoritative image, which represents the mark  
115 pattern observed in the majority of voucher samples, is annotated with measurements and a  
116 comment about the thermal mark (Figure 1). Occasionally, thermal marking procedures can  
117 produce errant mark patterns or multiple pattern variants of the planned mark (Figure 3). When  
118 this occurs, images of these mark variants are included in the reference collection. The Mark  
119 Characteristic Report and the thermal mark reference collection are both available online:

120 <http://www.taglab.org/OTO/reports/VoucherSummary.asp>

### 121 *Reader Training*

122 Prior to each field season, laboratory staff (or “readers”) gain familiarity with the thermal mark  
123 patterns likely to appear in AHRP samples by studying the physical and online reference  
124 collection of marked otoliths maintained at the MTA Lab. Familiarization with thermal mark  
125 patterns is important because growth rings in otoliths of wild salmon can occasionally appear to  
126 be similar to marks create during the thermal marking process. This review of known marks  
127 helps to minimize the chance of labeling an otolith as marked when it is actually wild as well as  
128 helps to increase reader accuracy and precision with regards to mark identification.

129 Laboratory personnel are trained to process adult otoliths using surplus otoliths to practice  
130 grinding to visually enhance the core or the “primordia” of the otolith. Staff learns to reduce  
131 processing time by controlling the pressure exerted during grinding and by becoming familiar  
132 with variations in otolith patterns and shapes. After approximately two to four weeks of training,  
133 laboratory staff begins to examine samples containing a mixture of marked and unmarked  
134 otoliths. Experienced personnel work with new staff members until their reader agreement is at  
135 least 95%.

### 136 *First and Second Reads*

137 All chum salmon (*O. keta*) otoliths are examined twice. In other words, these samples are read  
138 independently by a first reader and then read a second time by a different reader. The second  
139 reader typically knows who read the first sample but has no knowledge of the previous read  
140 results. Thus, we consider these to be a blind second read. The AHRP stream and pedigree  
141 samples are stratified into four areas (Figure 4). Disagreements between first and second readers  
142 are resolved by a third reader examining the otolith. The third read is not independent. The third  
143 reader knows who conducted both first and second reads and is cognizant of the results of each  
144 read. Second reads are performed as first reads are completed, and readers review the results. If  
145 disagreements occur, these are discussed, increasing familiarity with challenging patterns.

### 146 *Study Design*

147 Samples are assigned to readers by sample location (area) and over time. The MTA Lab  
148 currently uses four readers, thus there are six reader-pair combinations, which is critical for data  
149 analysis using a latent class model (see below). For the AHRP, the stream strata include four  
150 geographic areas in Southeast Alaska (Figure 4). Four streams were chosen for the pedigree  
151 sites, and each pedigree stream is treated as one stratum.

### 152 *Read Assessment Methods*

153 The MTA Lab uses three methods to assess a reader's ability to determine the presence or  
154 absence of a thermal mark. These methods include two agreement measures (agreement matrix  
155 and *Kappa*) and a latent class model, part of a family of models that allow estimation of reader  
156 classification error through the use of spatial data and multiple independent readings.

#### 157 *1) Agreement Matrices*

158 As otoliths are examined, a preliminary review of results is conducted by cross-tabulating the  
159 first read and second read results (Table 1). Common in reliability studies (Blick and Hagen  
160 1998), this matrix highlights results to review in detail. The matrix also highlights thermal  
161 marks that are mistakenly termed wild fish, as well as thermal mark identifications with a high  
162 percentage of disagreement. The first reader's results are listed on the rows, while the second  
163 reader's results are listed in the columns. Table 1 shows the number of thermal marked fish as  
164 well as the number not marked (e.g. wild) and unreadable. The numbers on the diagonal  
165 between the rows and columns indicates the number of thermal marks upon which the two  
166 readers agreed. Numbers off the diagonal highlight the disagreements (Table 1). For example,  
167 reader one and two agreed that 34 otoliths were TM3, but reader one called two otoliths TM4  
168 and reader two labeled them TM3. Discrepancies in whether the otoliths are marked or  
169 unmarked are located on the edge of the matrix, and differences in readability may also be found  
170 by examining the matrix. For example, six otoliths were labeled TM4 by reader one but were  
171 called "wild" by reader two, and three otolith were called wild by reader one but labeled TM3 by  
172 reader two. Examination of the matrix provides a preliminary analysis during a project and  
173 allows biologists to target areas for review. Deviations from the diagonal are reviewed, and

174 sometimes otoliths are read a third time to ensure consistency. This matrix has been a useful tool  
 175 for highlighting when a reader missed a mark. Often such errors are caused by incorrect sample  
 176 preparations. If an otolith is not ground enough, the thermal mark will not be visible. In such  
 177 cases, the sample is simply ground some more until the core is visible. Conversely, if an otolith  
 178 is ground too much, the mark will be removed. In this instance, the other otolith can be prepared  
 179 for mark recovery since both left and right otoliths will exhibit a thermal mark.

180 *2) Latent Class Model*

181 Latent class models (LCMs) provide an alternative approach to estimating agreement (Hui and  
 182 Walter 1980). LCMs incorporate an estimate of reader classification error, so that the variability  
 183 of reader agreement may be estimated. These models hypothesize the existence of unobservable  
 184 (i.e. “latent”) variables about which information can only be obtained through measurements on  
 185 observable (i.e. “manifest”) variables (Blick and Hagen 1998). LCMs use categorical variables  
 186 for the latent and manifest variables. For the AHRP, the latent variable is whether an otolith is  
 187 hatchery or wild; whereas, the manifest variables are a reader’s classifications. Because the true  
 188 error rate for each reader is unknown, latent class models provide a method to assess the  
 189 accuracy of thermal mark results. Blick and Hagen (1998) demonstrated that LCMs could be  
 190 successfully applied to thermal mark results by setting additional constraints or collecting  
 191 additional information.

192 The most economical LCM method is to separate the study area into strata and use two readers.  
 193 Use of three or more readers would give more degrees of freedom (*df*) and improve model  
 194 results, but the cost of the project would increase. Maximum likelihood models are the preferred  
 195 method for estimating LCMs. Assuming readings are independent among readers and among  
 196 otoliths, the likelihood function is as follows:

197

$$\prod_{i=H,W} \prod_{j=H,W} \prod_{k=H,W} \left\{ p\pi_{i|H}^{(1)}\pi_{j|H}^{(2)}\pi_{k|H}^{(3)} + (1-p)\pi_{i|W}^{(1)}\pi_{j|W}^{(2)}\pi_{k|W}^{(3)} \right\}^{n_{ijk}}$$

198 *where*

199 H = hatchery (thermal marked)

200 W = wild (unmarked)

201 *n* = sample size

202  $\pi_{i|j}^{(k)}$  = probability that reader *k* classifies an otolith as *i* when its true state is *j*

203 *p* = proportion of hatchery fish

204

205 The likelihood functions used to estimate the above parameters are maximized using Solver in  
 206 Microsoft Excel. Standard errors are estimated using the jackknife method (Haddon 2001).

207 When there are only two readers, neither is a standard, and there are five parameters to estimate  
 208  $\pi_{H|H}^{(1)}$ ,  $\pi_{H|H}^{(2)}$ ,  $\pi_{W|W}^{(1)}$ ,  $\pi_{W|W}^{(2)}$  and  $p$ , which gives only three  $df$  (four data points – one due to fixed  
 209 sample size,  $n$ ). To prevent overparameterization, constraints on the parameters or more data are  
 210 needed. Possible constraints include: 1) considering two parameters as known (e.g.;  $\pi_{W|W}^{(1)} =$   
 211  $\pi_{W|W}^{(2)} = 1$ , both readers will call a wild stock correctly); or 2) considering two sets of parameters  
 212 equal (e.g.;  $\pi_{H|H}^{(1)} = \pi_{H|H}^{(2)} = \pi_{W|W}^{(1)} = \pi_{W|W}^{(2)}$ , the accuracy rates are the same for both readers).  
 213 These constraints are likely unrealistic, thus more data are necessary. One way to generate more  
 214 information is to have a third independent reader (Walter 1984). Three readers provide seven  
 215 parameters:  $\pi_{H|H}^{(1)(2)(3)}$ ,  $\pi_{W|W}^{(1)(2)(3)}$ , and  $p$ , thus there are  $2^3 - 1 = 7$   $df$ , so all parameters may be  
 216 estimated. On the other hand, adding a third reader is usually logistically unfeasible given the  
 217 financial constraints of a project.

218 Hui and Walter (1980) proposed an alternative method to generate information. They suggested  
 219 that if there are two or more strata with different hatchery proportions in each strata (Blick and  
 220 Hagen 1998), then reader results could be stratified temporally or spatially. We can then assume  
 221 that  $\pi_{H|H}^{(k)}$  and  $\pi_{W|W}^{(k)}$  remains constant across strata (Blick and Hagen 1998), reducing model  
 222 parameters to eight with 12  $df$ . Thus, a two reader – four strata model would have 4  $df$  extra for  
 223 goodness-of-fit, preventing overparameterization of the model.

224 The following is the likelihood function for the two independent reads with  $S$  strata (Hui and  
 225 Walter 1980):

$$\prod_{g=1}^S \prod_{i=H,W} \prod_{j=H,W} \{p_g \pi_{i|H}^{(1)} \pi_{j|H}^{(2)} + (1-p) \pi_{i|W}^{(1)} \pi_{j|W}^{(2)}\}^{n_{gij}}$$

226 To estimate the latent variable for each reader, the stream samples collected during the AHRP  
 227 project were separated into four spatial strata (Figure 4). These spatial strata included: (1)  
 228 Southern Southeast waters; (2) Lynn Canal and Stephens Passage; (3) Chatham and Icy Straits;  
 229 and (4) Northern Outside waters. Samples were apportioned fairly equally across area. In  
 230 addition, these areas provided both geographic coverage and geospatial separation. Pedigree  
 231 samples were separated into strata based on the four creeks used in the project: Fish, Prospect,  
 232 Admiralty, and Sawmill creeks. Care was taken to distribute readings evenly among readers,  
 233 across areas, and by time. Samples were distributed among readers equally because we have  
 234 observed that when the LCM was heavily weighted by one individual, it performed poorly.

235 We have also observed that “reader drift” can occur over time as readers observe more marks  
 236 and sometimes altered their initial perception of a mark pattern (intra-rater reliability). To ensure  
 237 that the LCM analyses included this potential scenario, we assigned readers samples from across  
 238 the entire study period.

239 A critical assumption for both the LCM estimates of reader ability to detect a mark and *kappa*  
240 agreement values (see below) is that readings are independent, meaning that the reading of each  
241 otolith by a reader is independent of any other reading by the same reader and independent of  
242 readings by other readers for a given otolith. To support these assumptions, otolith first and  
243 second reads are provided to readers in random order by box. Another assumption is that  
244 individual accuracy rates are known to be greater than the error rates (Blick and Hagen 1998).  
245 Historically, reader agreement associated with mark recoveries conducted during the commercial  
246 sockeye fishery exceed 95%, so we believe this assumption is likely valid for the MTA Lab.

### 247 3) *Kappa*

248 The *kappa* statistic (Fleiss 1981) is frequently used to test inter-rater reliability. Rater reliability  
249 represents the extent to which the data collected in a study represent the variables measured. The  
250 *kappa* statistic provides examination of overall agreement between readers as well as agreement  
251 by specific thermal mark and an associated standard error (Fleiss 1981). Individual *kappa*  
252 statistics can be calculated for each category and pooled from different trials. Traditionally,  
253 inter-rater reliability was measured as percent agreement, calculated as the number of agreement  
254 scores divided by the total number of scores. Cohen (1960) critiqued the use of percent  
255 agreement due to its inability to account for chance, thus percent agreement tends to be higher  
256 when a category being rated has a high probability of occurrence. He introduced the Cohen's  
257 *kappa* (1960), which is chance corrected or accounts for the possibility that raters guess on some  
258 variables due to uncertainty.

259 *Kappa* is calculated by correcting the observed agreement for the degree of agreement expected  
260 by chance alone ( $P_o = (n_{HH} + n_{WW})/n$ ). Overall *kappa* is weighted and is defined as:

$$\hat{\kappa}_w = \frac{P_o - P_e}{1 - P_e} \quad (3)$$

262 where  $P_e$  is the proportion of expected agreement =  $(n_H n_H + n_W n_W)/n^2$  (Cohen 1960; Blick  
263 and Hagen 1998; Fleiss 1981). The weighted version of *kappa* has the same properties discussed  
264 above, but it is adjusted by giving lower weight to disagreements over marks with small numbers  
265 and full weight to disagreements over marks where agreement is high (Hagen et al. 1995). This  
266 better reflects agreement on what is marked and unmarked and reduces the influence of mark  
267 identifications with only one or two otoliths. Overall  $\hat{\kappa}$ , which assesses overall agreement  
268 between readers, is a weighted average of individual  $\hat{\kappa}$  for each individual thermal mark  
269 identified and is equal to the sum of the individual  $p_o - p_e$  (i.e., the sum of the numerators of the  
270 individual  $\hat{\kappa}$ ) divided by the sum of the individual  $1 - p_e$  differences (i.e., the sum of the  
271 denominators of individual  $\hat{\kappa}$ , Fleiss 1981).

272 The standard error for  $\hat{\kappa}_w$  is estimated by:



273 
$$SE(\hat{\kappa}_w) = \frac{\sqrt{A+B-C}}{(1-p_e)\sqrt{n}} \quad (4)$$

274 where

275 
$$A = \sum_{i=1}^n p_{ij} [1 - (p_i + p_j) + (1 - \hat{\kappa}_w)]^2, \quad (5)$$

276 
$$B = (1 - \hat{\kappa}_w)^2 \sum \sum p_{ij} (p_i + p_j)^2, \quad (6)$$

277 and

278 
$$C = [\hat{\kappa}_w - p_e(1 - \hat{\kappa}_w)]^2 \quad (7)$$

279 for readers  $i$  and  $j$  who have read  $n$  samples.

280 Although *kappa* is a commonly used inter-rater reliability statistical test, it has limitations.  
 281 Judgments about what level of *kappa* is acceptable are often questioned. As in most correlation  
 282 statistics, *kappa* values range from -1 to +1, where  $\hat{\kappa}_w = 1$  indicates complete agreement and  $\hat{\kappa}_w$   
 283 = -1 indicates complete disagreement. If observed agreement is greater than or equal to chance  
 284 agreement,  $\hat{\kappa}_w \geq 0$ , and if observed agreement is less than or equal to chance alone,  $\hat{\kappa}_w \leq 0$   
 285 (Landis and Koch 1977). Landis and Koch (1977) suggested that  $\hat{\kappa}_w > 0.61$  indicates substantial  
 286 agreement beyond chance. Values between 0.41 and 0.60 represent moderate agreement, and  $\hat{\kappa}_w$   
 287 < 0.40 represent slight to poor agreement (Landis and Koch 1977). Although Landis and Koch  
 288 (1977) interpreted a *kappa* score of 0.41 as acceptable, this might be considered too lenient for a  
 289 project like AHRP.

290 At the MTA Lab, we use *kappa* to ascertain amount of agreement among marks between readers.  
 291 Overall *kappa* among a suite of marks can be high (>0.80), but sometimes *kappa* scores for  
 292 individual marks can be low (<0.50). This occurs for a variety of reasons: 1) the mark was rarely  
 293 observed in a sample, usually older-aged fish; 2) otoliths were over- or underground; 3) mark  
 294 application was incomplete or differed among incubation groups, causing recovering to be  
 295 challenging; and 4) duplication of mark patterns among brood years required that otoliths be  
 296 aged to differentiate between years. Once we have determined why errors occurred, we  
 297 determine whether a higher proportion of the sample need to be second read or whether we need  
 298 to have some samples re-examined to determine whether marks were missed (i.e.; mount right  
 299 side of otolith and examine for thermal mark by a third reader). In the last instance, we work  
 300 with staff to improve thermal mark identification proficiency.

301 Thermal marks with poor *kappa* values are examined and discussed among readers during each  
 302 year of the project. They are also targeted for study prior to each project year. If a sample has a  
 303 poor overall *kappa* value, then those otoliths are examined further to determine the cause (i.e.;

304 multiple poor marks or a sample coordination errors). *Kappa* values are archived on the local  
305 network.

306 Because *Kappa* is an index, it is important to remember that interpretation can be affected by the  
307 values of the underlying parameters (Blick and Hagen 1998). Thus, direct comparison of  $\hat{\kappa}$   
308 across populations with different underlying proportions is not appropriate. Although agreement  
309 measures may be subject to some ambiguity, they are useful in monitoring results for potential  
310 errors and pinpointing areas for the Lab to re-examine.

## 311 **Discussion**

312 Fisheries research often requires that trained individuals classify data according to a strict but  
313 somewhat subjective set of rules. In many situations, there is no standard available with which  
314 to confirm classifications, and it is necessary to apply some other method to determine the  
315 accuracy of the determinations. Distinguishing thermal-marked fish from wild fish is a good  
316 example of this type of problem because: 1) most thermal-marked salmon do not receive a  
317 secondary mark, so cross-validation is not possible; and 2) the ability to read otoliths for thermal  
318 mark presence and identification requires training and experience because natural variation in  
319 growth rings observed in chum salmon otoliths can appear similar to thermal mark patterns. In  
320 the absence of samples of known origin, it is common to collect multiple, independent  
321 observations of the same samples and assume that percent agreement among readers serves as a  
322 proxy for read accuracy. Agreement indices (matrices and *kappa*) are easy to compute and  
323 indicate read discrepancies in mark recovery and identifications. For the AHRP project, these  
324 QA/QC methods provide additional direction for validation of reader accuracy and precision.  
325 They also provide some quantitative indication of reader accuracy.

326 In addition, we use the agreement measures described above to highlight results in need of closer  
327 examination and suggest potential areas for critical review. When agreement measures indicate  
328 that results require evaluation, we examine the data to determine whether we need to: 1) conduct  
329 additional reader training when an individual is under- or over-grinding and missing marks, 2)  
330 read samples a third time by another independent reader when marks are especially difficult to  
331 discern, and 3) examine potential issues in greater detail during the next season's training period  
332 if a particular mark or brood year is expected to return.

333 Although these indices are fairly easy to calculate and are useful indicators of reading problems,  
334 it is important to remember that some of these indices are not directly comparable. It is difficult  
335 to compare *kappa* statistics across populations with different underlying proportions. Because of  
336 this, even when a suite of *kappas* is consistent, it may not be clear how reader  
337 agreement/disagreement influences the contribution estimate. In addition, these indices do not  
338 provide inferences about the relative ability of one reader over another to determine a particular  
339 set of patterns. Latent class models, however, provide readily interpretable qualities that can be  
340 easily calculated. Classification accuracies or errors provide direct, meaningful parameters,

341 unlike the use of an index of agreement alone. In addition, LCMs provide estimates of hatchery  
342 proportions ( $p$ ).

343 We feel that the procedures described above provide a combination of approaches to provide a  
344 comprehensive examination of error rates and accuracy of reads conducted in the MTA Lab.  
345 The matrices and *kappa* statistics point out areas for review, and the LCM provides direct,  
346 meaningful parameters that can be compared from year-to-year.

347 **Questions for AHRP Science Panel**

- 348 1) Are the methods presented here adequate for assessing accuracy of detecting the presence of  
349 a hatchery (thermal) mark?  
350 2) Are the methods presented here adequate for assessing the accuracy of identifying hatchery-  
351 specific marks?

352 **AHRP Review and Comments**

353 *This technical document has been reviewed.*

354 This document covers some of the long and well established procedures used by the Alaska  
355 Department of Fish and Game, Mark Tag and Age Lab for thermal mark recovery. There were  
356 no comments from the AHRG.

357 This document is acceptable to the AHRG.

358  
359  
360  
361  
362  
363  
364  
365  
366  
367  
368  
369  
370  
371  
372  
373  
374  
375  
376  
377  
378  
379  
380  
381  
382  
383

## References

Blick, D. J., and P. T. Hagen. 1998. The use of agreement measures and latent class models to assess the reliability of thermally marked otolith classifications. *NPAFC Doc 370*:1-15.

Campana, S. E. 1983. Calcium deposition and otolith check formation during periods of stress in coho salmon, *Oncorhynchus kisutch*. *Comparative Biochemistry and Physiology* 75A.(2):215-220.

Cohen, J. 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement* XX(1):37-46.

Fleiss, J. L. 1981. *Statistical methods for rates and proportions* 2nd edition. John Wiley and Sons, New York, NY.

Gwet, K. L. 2014. *Handbook of inter-rater reliability*, 4th edition. Advanced Analytics, Gaithersburg, MD.

Haddon, M. 2001. *Modelling and quantitative methods in fisheries*. Chapman and Hall, New York, NY.

Hagen, P., K. Munk, B. W. Van Alen, and B. White. 1995. Thermal mark technology for inseason fisheries management: a case study. *Alaska Fishery Research Bulletin* 2(2):143-155.

Hui, S. L., and S. D. Walter. 1980. Estimating the error rates of diagnostic tests. *Biometrics* 36:167-171.

Josephson, R., B. A. Agler, K. F. Van Kirk, and D. S. Oxman. 2006. A proposal to simplify the thermal mark code notation. *NPAFC Doc 944*:1-4.

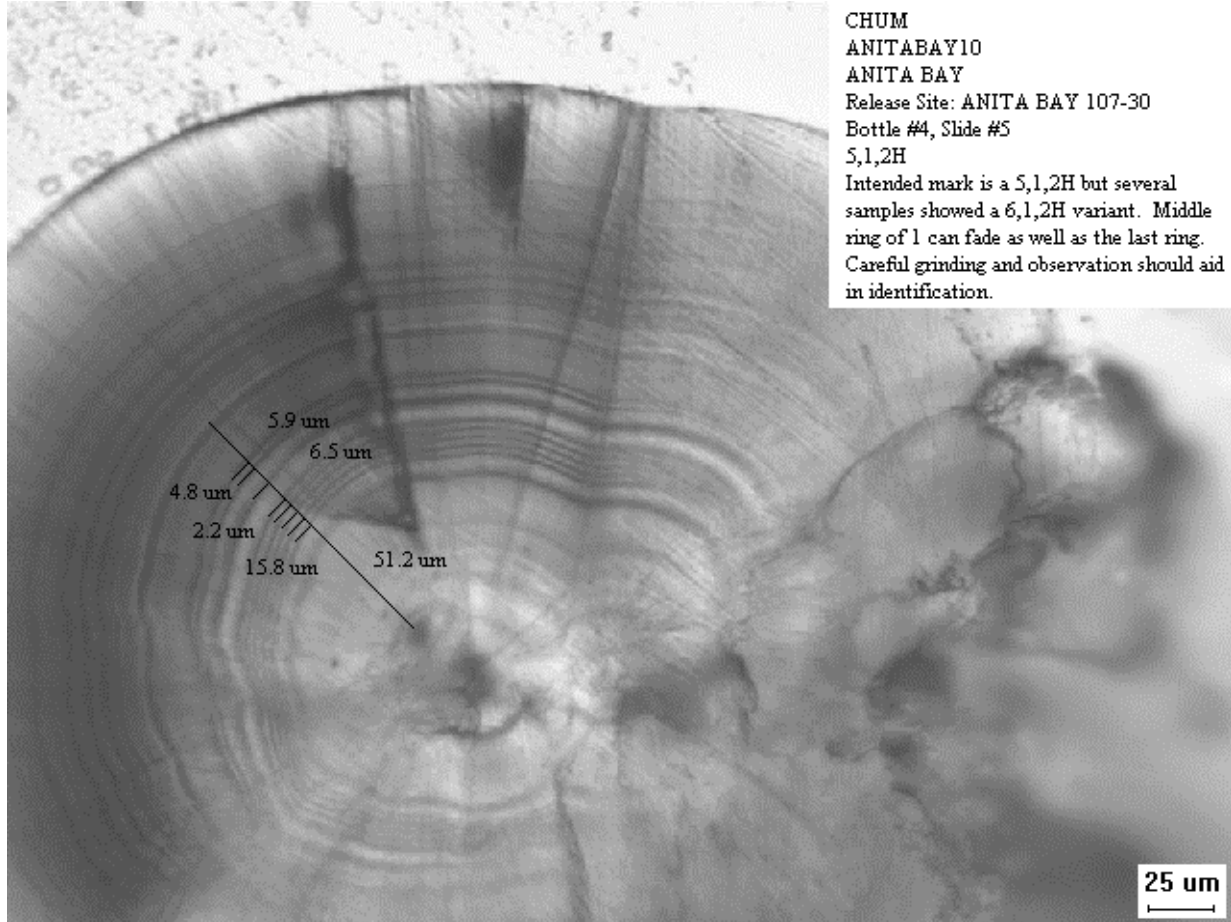
Landis, J. R., and G. G. Koch. 1977. The measurement of observer agreement for categorical data. *Biometrics* 33:159-174.

Vercessi, L. 2014. Alaska salmon fisheries enhancement program 2013 annual report. Alaska Department of Fish and Game, Anchorage.

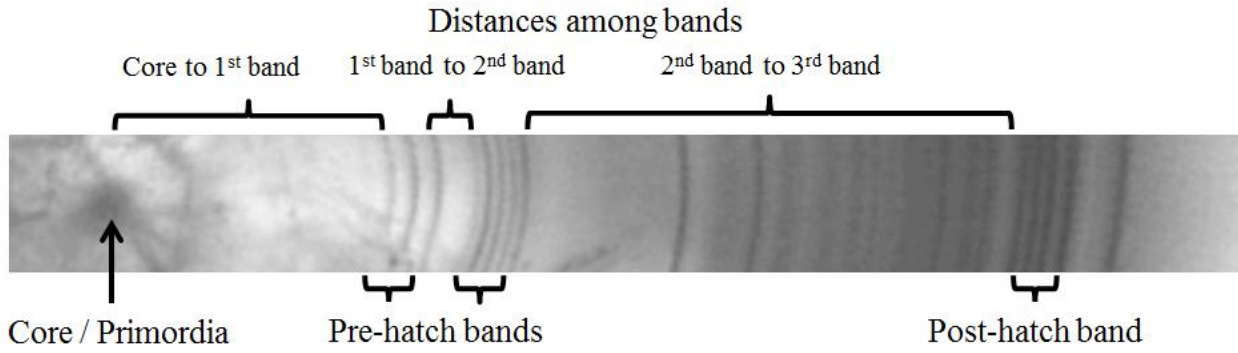
Volk, E. C., S. L. Schroder, and K. L. Fresh. 1990. Inducement of unique otolith banding patterns as a practical means to mass-mark juvenile Pacific salmon. *American Fisheries Society Symposium* 7:203-215.

Walter, S. D. 1984. Measuring the reliability of clinical data: the case for using three observers. *Revue d'épidémiologie et de santé publique* 32(3-4):206-211.

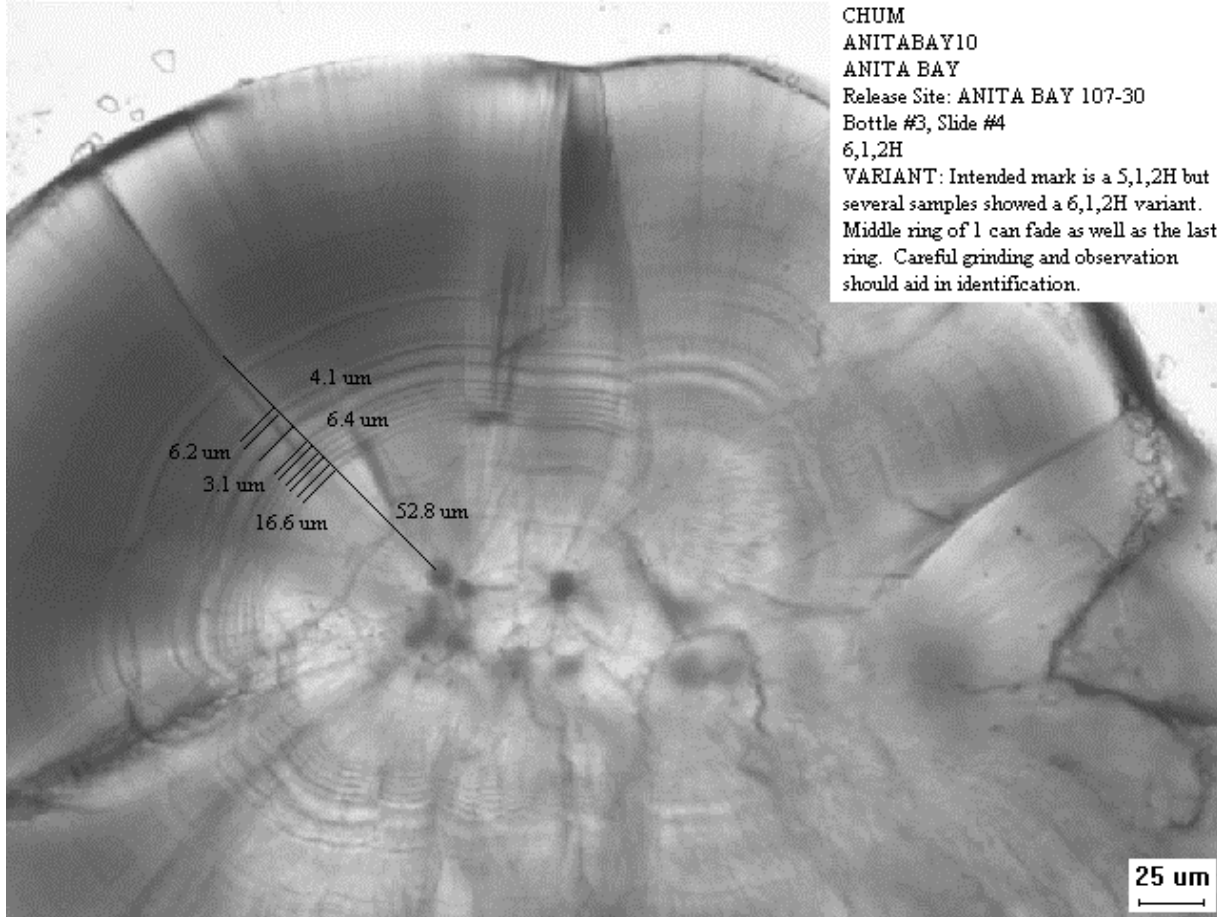
## Figures



385  
386 Figure 1. Image of a thermal mark reference specimen. From Southern Southeast Regional  
387 Aquaculture Association (SSRAA), this brood year 2010 mark (ANITABAY10) has a hatch  
388 code of 5,1,2H. The code indicates that the first band from the otolith's core contains five dark  
389 dark rings, then there is a space, followed by a band with one ring, followed by another space and a  
390 final band with 2 rings prior to the hatch mark (blurry, wide, dark area). Annotated  
391 measurements on the transect line include distance from otolith core (primordia) to first band,  
392 width of first band, space between first and second bands, and average distance between rings in  
393 each band. All thermal mark images are available online through the North Pacific Anadromous  
394 Fish Commission (NPAFC) Working Group on Salmon Marking (WGOSM) website:  
395 <http://wgosm.npafc.org/MarkSummary.asp>  
396

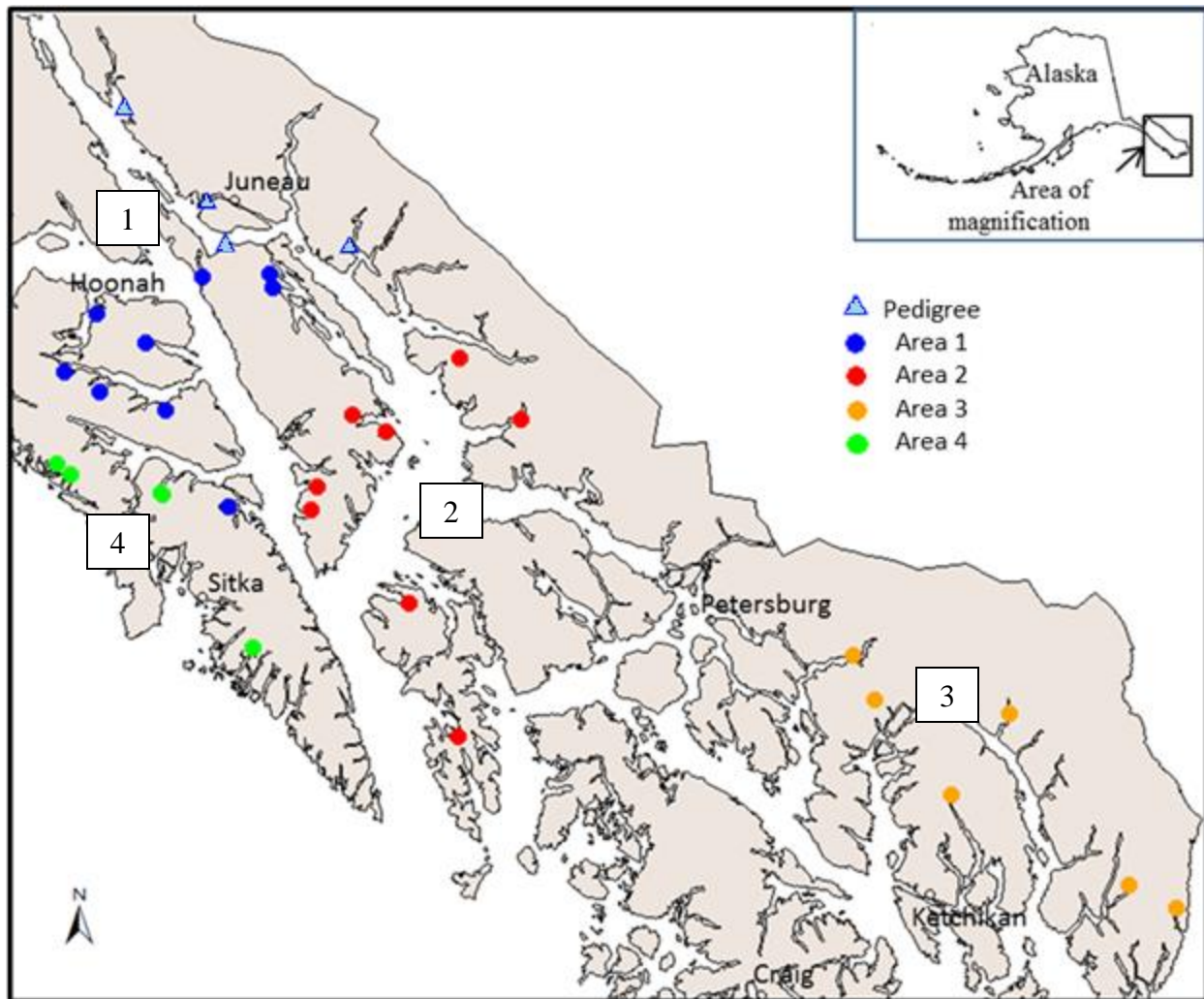


397  
 398 Figure 2. Thermal mark image with measurements shown in the Mark Characteristic Report.  
 399 This figure shows a 3,5nH4 mark with a pre- and post-hatch mark. Thus this mark has two  
 400 bands prior to hatch (the first with three rings and the second with 5 rings) and one band after the  
 401 hatch containing 4 rings. The individual rings are the dark lines in each band, and in the second  
 402 band, the spacing among the rings is narrower than that in the other bands so the 5 is followed by  
 403 an “n.”  
 404  
 405  
 406



407  
 408  
 409  
 410  
 411  
 412  
 413  
 414

Figure 3. Image of a thermal mark variant of a reference specimen. This figure shows another image of Figure 1, thermal mark ID ANITABAY10. This fish, assigned a target thermal mark of 5,1,2H, which indicates that the first band from the otolith core contains five dark rings, a space, then a band with one ring, a space, and a band with two rings followed by the hatch mark (the blurry, wider, dark area). Instead, this otolith shows a 6,1,2H or a variant, meaning that the first band has six rings instead of the planned five rings.



415

416 Figure 4. Four strata used for assessing the accuracy of thermal mark readings of chum salmon  
 417 otoliths recovered from streams in Southeast Alaska during 2013 and 2014 for the Alaska  
 418 Hatchery Research Project.

419

420

421

422

423

424



425

### Tables

426 Table 1. Example matrix comparing thermal mark reader agreement. Row and column names  
 427 represent potential thermal marks identified by each reader (TM1 through TM6), otoliths  
 428 classified as wild, and otoliths classified as unreadable (ND). The number of otoliths where both  
 429 readers agree is in bold font along the diagonal between the row and columns.

1 <sup>st</sup> Reads	2 <sup>nd</sup> Reads								Total
	TM 1	TM 2	TM 3	TM 4	TM 5	TM 6	Wild	ND	
TM 1	<b>0</b>	1							1
TM 2	1	<b>12</b>							13
TM 3			<b>34</b>						34
TM 4			2	<b>9</b>			6		11
TM 5					<b>26</b>				26
TM 6						<b>4</b>			4
Wild			3				<b>357</b>	1	358
ND							1	<b>3</b>	4
Total	1	13	36	9	26	4	358	4	451

430

431