

172 fit by the model (Figure 5). Of potential stock concern is that the NMFS longline survey is generally
173 under the expected survey abundance since 2010 (Figure 5), suggesting that information on larger fish in
174 the population added by this survey leads to a more pessimistic assessment of overall stock depletion (as
175 indeed shown by model 15.6A results). However, the model is not fully tuned, so such supposition may
176 be premature. However, it does highlight that **if the index is to be used, some evaluation of possible bias**
177 **in relation to Pacific cod, perhaps most importantly since 2010 is required.** The model that includes the
178 NMFS longline survey is able to fit the associated length compositions well.

179 **Before deciding to include the NMFS longline index and associated lengths in a proposed central EBS SS**
180 **model, an investigation into the properties of the EBS NMFS longline index in relation to Pacific cod in**
181 **particular should be done. The investigation should particularly examine possible bias in the index since**
182 **2010 as this appears to be influential on assessment results. If the index is found unlikely to be biased,**
183 **then I recommend inclusion in the model with additional sd estimated.**

184 Aleutian Islands

185 The overall fits by the AI model to lengths (Figure 6) and the abundance index appear reasonable.
186 Abundance index point estimates for 2004 and 2014 appear to most conflict with other information in the
187 AI model.

188 **Before deciding to include the NMFS longline index and associated lengths in a proposed central AI SS**
189 **model, an investigation into the properties of the AI NMFS longline index in relation to Pacific cod in**
190 **particular should be done. If the index is found unlikely to be biased, then I recommend inclusion in the**
191 **model with additional sd estimated.**

192 How should the various data sets be weighted?

193 For abundance index data, iterative reweighting to potentially allow additional index error was previously
194 an accepted procedure for many US and Australian stock synthesis assessments. Such iteration was done
195 manually, and more recently **the ability to internally estimate additional index error (via an additional sd)**
196 **has been added as an option to SS. Use of that option has become accepted practice for many recent**
197 **assessments. Estimation of additional index error is normally done for all indices included in a stock**
198 **assessment as (perhaps in my naive interpretation), the input variability usually only accounts for**
199 **measurement error and the process error component is unknown.**

200 **Input sample sizes for composition data have an influence on assessment results and it has also become**
201 **generally accepted practice for those sample sizes to more reflect the number of sampled fishing trips,**
202 **rather than the number of fish measured.**

203 Relative data weighting in stock assessments for composition data and the goal of standardized
204 approaches has been the subject of recent and ongoing research particularly in the US west-coast, and the
205 subject of a Center for the Advancement of Population Assessment Methodology (CAPAM) workshop in
206 La Jolla, CA in October of 2015 (<http://www.capamresearch.org/data-weighting/workshop>). **While there**
207 **has been some recent narrowing down of agreed procedures among US west-coast stock assessors, it has**
208 **also been recognized that it is not currently possible to recommend default procedures for composition**
209 **and conditional age-at-length (CAAL) data. There is agreement that the Francis weighting approach is**
210 **more appropriate in cases where the model is not correctly specified as it takes autocorrelation among**
211 **composition data into account. It is also agreed that for a correctly specified model, the McAllister-Ianelli**
212 **harmonic mean weighting method works well.** Both of these procedures have been extended from
213 marginal length or age composition data to conditional age-at-length (Francis A and B methods are
214 available for CAAL, with Francis B potentially preferred). **A possible further development that may**

215 provide a direction forward is using the Dirichlet multinomial likelihood (Thorson, 2014), although this
216 method will require review and implementation in SS before it may be used. Recent simulation work has
217 shown that the McAllister-Ianelli arithmetic mean procedure is inferior to other methods (Punt, In press).

218 What form (i.e., Stock Synthesis “pattern”) should be used for the selectivity functions?

219 SS provides a large number of selectivity pattern options (14 size and 12 age patterns excluding special,
220 discontinued and mirror – SS user manual v 3.24s). By far the most commonly used patterns in recent
221 stock assessments are logistic for simple asymptotic selectivity or the double-normal (most often size
222 pattern 24 or age pattern 20) where selectivity is allowed to be dome-shaped. The flexibility of the
223 double-normal is usually sufficient to account for the wide range of single-peaked shapes that may be
224 expected from a single fishing gear type. It is also possible to combine size and age selectivity patterns for
225 a fishery or survey and to have differential selectivity by sex to, for example, account for reduced
226 availability of older females in the population. To most easily account for “odd-shaped” selection that
227 may be due to, for example, a combined fishery composed of several gear types, SS provides an age
228 based selection pattern that generates an age-based random walk (age pattern 17).

229 Normally, fishery and survey selection is assumed to be primarily a length-based process as fishing gear
230 selection is usually size-dependent. However, selectivity in an assessment model combines gear
231 vulnerability with availability. Whether availability (e.g. due to migration, aggregation [e.g. for
232 spawning], schooling) is age- or length-based is a more difficult question, so although length-based
233 selection may be preferred for modeling, a case can still be made for age-based selectivity.

234 Generally, the selectivity pattern should be chosen (most likely from the options above) that has the
235 fewest parameters, and allows an acceptable fit to the available composition data (e.g. no bands at
236 particular lengths of significant length composition residuals). As surveys are designed to at least use the
237 same fishing gear throughout, a good reason to use more complex patterns than logistic or double-normal
238 would be required for those. If a fishery has fairly homogenous gear, a similar argument applies there as
239 well. In the case of a fishery with mixed gear types, an opportunity exists to use a less restricted pattern
240 shape, as provided by the age-based random walk. At present, I don't think a random-walk length-based
241 pattern is available, so selectivity in that case is restricted to being age-based.

242 Should the models be structured with respect to season?

243 It is usual practice for SS models to separate input data from surveys and fisheries that have demonstrably
244 different selectivity if data are available to do so. Normally, the minimum requirement to allow data
245 partitioning according to season, gear type or area is that a number of years of length or age composition
246 data that are believed to be representatively sampled are available within each partition. Partitioning of
247 composition data is only usually necessary if summary length/age compositions from comparable
248 partitions show obvious apparent selectivity differences. Partitioning may also be required for abundance
249 indices if different trends are observed by partition.

250 Models that specifically address the exploration of alternative structures regarding selectivity partitions
251 have been developed and were presented for the EBS, so the discussion here will be confined to models
252 from that region.

253 Simple examination of aggregated length data for the EBS shelf trawl survey, the slope survey, longline
254 fishery and NMFS and IPHC longline surveys (Figure 1) show a marked difference in the shelf trawl
255 survey to all of the others. Unfortunately, the trawl and pot fisheries were not included, but we know from
256 diagnostic output from model 11.5 that trawl fishery selectivity seems to be intermediate between the
257 trawl survey and longline fishery, and the pot fishery seems similar to the longline fishery (Figure 8).

258 Also notable is that the Jan-Apr trawl fishery lengths show a peak that is consistent with longline fisheries
259 during that period only, which corresponds to the spawning season. Conjecture has been made about
260 possible movement of larger fish from the NBS area, although another explanation may be the movement
261 of larger fish from waters targeted by the longline and pot fisheries into shelf trawl areas during the
262 spawning season. There is little information available from tagging and none that can address the question
263 of movement in and out of the NBS. The shelf trawl survey is made outside of the spawning season, and
264 at that time, less of the larger fish seem to be available on the shelf, although tagging of a small number of
265 fish does indicate apparent random movement of fish over the shelf during that time.

266 For modeling purposes, the model only requires that the composition of the fishery catches be adequately
267 accounted for each year, and the more important population abundance trends are taken from surveys (at
268 least for the models here). The difference in trawl fishery selection by season seems to be a feature that
269 can be addressed through seasonal model structure. This is done to some extent with model 11.5, but the
270 fit to the Jan-Apr trawl fishery length composition by that model is not particularly good (Figure 8). In
271 addition to gear/season partitioning, a large number of time blocks that allow selectivity to vary through
272 time have been used in model 11.5. It may be questioned whether such fine scale partitioning of the data
273 are supportable if partitioning and blocking first needs to be justified depending on whether prior data
274 examination or independent knowledge about changes in practices suggests that all of those partitions are
275 necessary, and that sufficient data are available within each to allow estimation of a different selectivity
276 pattern.

277 A new procedure for accounting for fishery selectivity has been proposed here in model 15.6 where an
278 age-varying random walk is used to characterize the selectivity for all combined fisheries (trawl, longline
279 and pots) each year. This procedure seems attractive given the high level of partitioning required for
280 model 11.5. If such a procedure can provide a means of accounting for total fishery removals each year
281 according to size/age, then it should be acceptable. Diagnostic plots for fishery lengths, both by year and
282 combined for model 15.6, show rather good fits to available data (all residuals are also within the range
283 -2.0 to 2.0). There is very little catch taken aged above about 8, so fixing selectivity above that age seems
284 reasonable.

285 As the proportion of trawl catch to longline has changed considerably over time, it would be expected that
286 large changes in the general pattern of selectivity would also be observed, that are somewhat evident in
287 the plot (Figure 9), but of possible concern. Is the amount of change consistent with the broad movement
288 of the fishery from trawl to longline over time?

289 Also of some concern is that the general fishery pattern for model 15.6 is dome-shaped, allowing the
290 model some flexibility to generate cryptic spawning biomass. This is also an area of on-going work, and
291 some diagnostics associated with it are in development or available from Github as additions to R4SS. At
292 present, the available code only works for 2 sex models, so cannot be applied here, but could be further
293 generalized to do so. The inclusion of surveys that are more directed towards the older fish in the
294 population help to alleviate cryptic biomass problems, and is therefore a further reason to consider the
295 addition of at least one longline survey to the base model.

296 I believe that options are only currently available in SS for a random walk by age for annual selectivity, as
297 used for model 15.6. If the same was done by length, more parameters would be required (if 1cm size
298 bins), or alternative bin patterns could be explored. Such a length-based exploration would be useful,
299 should such capability be available in SS.

300 As many current SS assessments grapple with highly partitioned fishery data, such a procedure has the
301 potential for resolving some of those problems also. I do not have previous personal experience with this
302 procedure, and am reluctant to agree on its use without a supporting simulation study that confirms its

303 equivalence or even superiority to a high degree of data partitioning. Such a study would be reasonably
304 easy to design and carry out. However, I am willing to agree that it seems to provide a good resolution to
305 the problem for the fishery selectivity in the EBS models.

306 Should the models be structured with respect to gear type?

307 As this question mostly relates to dealing with the fisheries and not surveys, the discussion under ToR
308 2.2.2c was generalized to address both season and gear type.

309 How much time variability should be allowed, and in which parameters?

310 The only population biological parameter allowed to vary with time in most SS stock assessments is
311 annual recruitment levels. Cumulative information on annual recruitment strength is provided fairly
312 directly by composition data, so the reasons especially for high peaks and troughs in recruitment are
313 usually apparent in the available data. It has also been recognized that other parameters are likely to vary
314 through time – in particular natural mortality, but also growth and maturity. For natural mortality it has
315 been considered difficult to estimate time trends in changes without strong independent estimates for
316 those changes, such as from ecosystem studies showing differences in predator abundance, and that time
317 trends in M are difficult to disentangle from other factors such as catch mis-specification (e.g. see
318 Brodziak et al., 2011). Allowing time variation in factors that directly affect productivity also lead to
319 questions about choice of appropriate time periods for the selection of management reference points, and
320 how to make appropriate stock projections.

321 Additional model parameters that may vary with time that are often dealt with using time-block methods
322 are fishery/survey selectivity and catchability. As already mentioned, for fisheries that are not associated
323 with an abundance index, a fairly freely estimated time-varying pattern (such as used for EBS model
324 15.6) may be acceptable if it suitably captures annual fishery removals by size/age. For surveys the
325 situation differs. Surveys are the most important source of abundance information for the model,
326 particularly because at least the gear selectivity can be maintained as a constant through time. Availability
327 (either by age or year) is another matter, but is usually treated as a source of additional random error. If a
328 true trend (or even a step) exists in either survey selectivity or catchability, then that survey is biased, and
329 the bias needs to be accounted for, or the survey truncated, split or discarded. Such a bias would ideally
330 be investigated and identified with a focused study and auxiliary data not necessarily used in the
331 assessment model. Adding annual time-variability to survey selectivity or catchability and finding that
332 trends are estimated may simply be providing a means for the model to trade trends in population
333 abundance to improve the fit to noisy composition data in preference to abundance indices. The reason
334 that such a model might result in trends in survey selectivity or catchability are not readily apparent from
335 standard input data sources, and may be difficult to diagnose. Results from estimation of annual
336 variability for the EBS trawl survey catchability in model 15.6 (Figure 10) do exhibit some runs in
337 residuals that may be of concern – particularly from 1993 to 1996. Time-changes in trawl survey
338 selectivity as estimated by the EBS model 15.6 shows very little change through time, suggesting that
339 time-variability in trawl survey selectivity as implemented is not required (Figure 11).

340 My own recommendation for now is that time variability should be allowed in a parameter when there is
341 an available reliable data source that fairly directly measures such a change, and that a trend exists in that
342 data source that needs to be captured by the assessment model. This situation only currently exists for
343 recruitment and fishery selectivity in the EBS model. It also provides some support to consider time
344 variability in weight-at-length or size-at-age if those data sets show considerable trends over time.

345 Others (e.g. Anders Nielsen, Jim Thorson) have proposed that a more appropriate way to deal with time
346 variability is to use mixed-effects models with time-varying “nuisance” variables such as recruitment

347 modeled as random effects. Improved solutions for time-varying parameters may be possible using all of
348 the currently available data sources, if/when SS RE becomes available.

349 What constraints, if any, should be placed on survey selectivity at older ages?

350 The models examined during the review for the EBS seem to fairly clearly demonstrate that the trawl
351 survey selectivity is dome-shaped. However, the possibility that the survey is in fact asymptotic has not
352 been eliminated. The extent of the survey dome-shape may, for example, be confounded with M. It may
353 be that different data sources are in conflict about the estimated value for M that can be diagnosed with a
354 Piner profile plot of likelihood components. Exploration of age-specific M (e.g. starting with a Lorenzen
355 function) could also be done.

356 A range of plausible alternative models should be explored, and the extent of the estimated dome
357 selectivity for the trawl survey examined for each to see if the dome is consistently required. However, as
358 the extent of the trawl survey dome is probably one of the major axes of uncertainty in the model at
359 present, it should remain freely estimated and informed by the available data in any chosen base model,
360 possibly with forcing more or less dome as sensitivity analyses in the final assessment.

361 What constraints, if any, should be placed on survey catchability?

362 Because of the history of the development and use of the trawl survey as an absolute index of abundance,
363 there remains some belief that there is sufficient information available to determine at least a plausible
364 acceptable range for survey q , and to some, that range could be perceived to be quite narrow. Much work
365 has been directed towards net avoidance and how that might be compensated by a q adjustment. I believe
366 that all major potential sources of error in survey q should at least be stated in an accessible document,
367 and errors in those dimensions at least be qualitatively examined and ranked. Those should include
368 avoidance and other gear-specific fish behavioral issues, and potential error in scaling the swept area
369 estimates to the population using assumptions about the population distribution during the survey by
370 depth and area, and also even the assumption of known stock boundaries. A qualitative evaluation such as
371 this would probably make it clear that the true error in q is reasonably high. It would also assist to
372 determine what priorities should be given to field studies that may be directed towards reduction of the
373 error in survey q and adjustments required to scale area swept biomass estimates to the total (available
374 given survey selectivity) population. An extension to a more quantitative evaluation of the potential errors
375 may also lead to a prior distribution for EBS shelf bottom trawl survey q that can be generally agreed, and
376 could then be used for modeling without much controversy. Without at least a comprehensive qualitative
377 evaluation of all major error sources, decisions about rejection of models that estimate q based on how
378 different the estimated q is from acceptable values remains difficult, and currently in the domain of
379 pragmatic judgment.

380 I believe that models that estimate the shelf bottom trawl survey q using a fairly non-informative prior (as
381 in model 15.6) should currently be preferred. Agreed bounds on prior survey q point estimates can be
382 used as one of the acceptance criteria for particular models. I personally have a fairly high tolerance for
383 those values (based however, on only a limited background knowledge for this particular survey), and am
384 comfortable with at least a factor of 2.0 (0.5 – 2.0 times the initial point estimates).

385 Should additional surveys be added to the models, q values for Pacific cod for those are less well
386 understood, and non-restrictive priors for those are preferable, with q estimated.

387 How should large gradients be dealt with in otherwise apparently converged models?

388 Large gradients are generally considered to be an indication of a problem. However, if the hessian can be
389 inverted and jitters also indicate convergence, then perhaps the problem is only minor. I do not have any
390 reason to doubt the explanation given in the EBS assessment document for why large gradients might
391 occur, but it does suggest to me that the implementation of age selectivity pattern 17 requires a closer
392 look to determine if the problem can be corrected (e.g. to determine whether it contains badly
393 behaved/non-differentiable “if” statements).

394 Anything else on which the reviewers care to comment

395 Retrospectives

396 Diagnosis of retrospective bias in stock assessments has received considerable past attention in the
397 literature and was also the subject of a BSAI/GOA working group in 2013 according to meeting
398 background information. Despite this attention, research is on-going, and means for diagnosis and
399 correction for retrospective patterns are not agreed. Several diagnostic measures are available including
400 Mohn’s ρ , the so-called Woods Hole ρ , and the RMSE method devised by the BSAI/GOA working group.
401 I am familiar with two rules of thumb that can be used to diagnose retrospective patterns that need to be
402 addressed in some way. The first and simplest is by Hurtado-Ferro et al. (2014) that says that “values of
403 Mohn’s ρ higher than 0.20 or lower than -0.15 for longer-lived species (upper and lower bounds of the
404 90% simulation intervals for the flatfish base case), or higher than 0.30 or lower than -0.22 for shorter-
405 lived species (upper and lower bounds of the 90% simulation intervals for the sardine base case) should
406 be cause for concern and taken as indicators of retrospective patterns.” The second by Brooks and Legault
407 (2015) from VPA assessments “is to plot the terminal year estimate of SSB(T) vs F(T) along with
408 bootstrap percentiles and compare that to the point estimate when SSB(T) and F(T) are adjusted by
409 ρ SSB,7 and ρ F,7, respectively” to see if the ρ -adjusted point estimate falls outside the bootstrap
410 percentiles on either axis - see Brooks and Legault (2015) for details. Brooks and Legault (2015) also
411 provide a procedure for adjustment of short-term projection results to account for substantial retrospective
412 patterns. Ideally, the diagnostics for a model acceptable for use for management advice should not show
413 significant retrospective bias. EBS model 11.5 and the initial AI SS models did show significant
414 retrospective bias (at least according to the Hurtado-Ferro et al. (2014) rule of thumb) that indicated that
415 results from those models are not reliable for use for management advice, and that improved alternative
416 models should be sought, or at least a projection correction may be required. Further model explorations
417 for both regions have found models that do not exhibit a strong retrospective bias, and on that basis would
418 be judged as improved models. Retrospective bias provides evidence for model mis-specification, but of
419 course, the lack of a retrospective bias does not prove that the model is correctly specified.

420 So-called Ianelli “squid plots” provide an additional useful means for looking at retrospective patterns in
421 annual recruitment deviations, but have potential application to any parameter allowed to deviate annually
422 in a model.

423 Catch uncertainty

424 As for many models, historical catch in particular is uncertain, and the best estimate of historical catch
425 has been made using assumptions that seem supportable. However, the construction of alternative
426 plausible historical catch scenarios would be useful for the determination of sensitivity of the model to
427 that uncertainty.

428 Steepness

429 Tier 3 methods by default assume a steepness value of 1.0. A requested run using a steepness value of 0.7
430 shows that EBS results are somewhat sensitive to the choice of steepness value, and this dimension of
431 uncertainty should be highlighted.

432 Regime change

433 A regime change in 1976-77 affecting log mean recruitment in EBS model 11.5 has been avoided in EBS
434 model 15.6 by starting the latter model after the regime change. Shifts in 1989 and 1999 have also been
435 suggested according to the ecosystem considerations in the assessment documentation. Regime change
436 was not examined at all during the review, but is another potential source of model uncertainty.

437 Inclusion of marginal age composition vs CAAL data

438 At present, both the EBS and AI enter age-at-length data as marginal age distributions. There has been a
439 gradual trend in stock assessments to make improved use of data from otoliths by entering the data into
440 models as conditional age-at-length. During the review the general wisdom of this approach was
441 questioned as it was mentioned that some recent assessments had reverted back to marginal age
442 distributions. A standard approach for dealing with age-at-length data currently seems to be unavailable.

1

2 *Reviewer 3: Jean-Jacques Maguire*

3 Executive summary

4 From what was discussed during the meeting and the documentation reviewed, there are no objective
5 reasons to reject the IPHC longline survey as an index of stock size, assuming it has been correctly put
6 together and calculated. The IPHC longline survey data should be thoroughly investigated. It should be
7 used in the assessment unless fatal flaws in the data, in the treatment of the data or in the survey
8 methodology are identified. Similar to the IPHC longline survey, there are no objective reasons to reject
9 the AFSC longline survey as an index of stock size. The AFSC longline survey should also be thoroughly
10 investigated and used in the assessment unless fatal flaws in the data, in the data treatment or in the
11 survey methodology are identified.

12 Regarding the form of the selectivity function, my preference would be to not allow too much flexibility
13 in selectivity changes over time and to not allow strange patterns (e.g. figures 2.1.3 in the Eastern Bering
14 Sea and 2A.11 and 2A.12 in the Aleutian Islands in the December 2015 SAFE report). If allowing these
15 strange patterns is a condition of getting a good fit or convergence, this would be a sign that something
16 else might be wrong. If allowed to change over time and age, the changes should be relatively smooth and
17 not result in peculiar patterns. The reason(s) for the apparent differences in selectivity between the IPHC
18 longline survey and the AFSC longline survey for lengths above 70cm should be further investigated.

19 It could be worth investigating further changes in growth (Figure 11), particularly with respect to the
20 implications for the assessment as growth changes may have an influence on fishing mortality and
21 population estimates.

22 In the Aleutian Islands area, it is unlikely that there is a single stock in the traditional understanding of the
23 concept. Simpler form of monitoring and management, in close cooperation with the industry and
24 possibly NGOs, could be a better way of protecting the resources and managing the fisheries.

25 One cannot model oneself out of lack of data, particularly for the Aleutian Islands assessment. Stock
26 Synthesis has so much flexibility that, given sufficient time, a skilled user can probably get almost any
27 stock trend from a dataset. Indices of abundance should be given more weight in the assessment than
28 length composition. Age composition, particularly from the commercial fishery, but also from surveys or
29 other indices of abundance can be very informative if analyzed appropriately. Information in the length
30 composition is at best indirect information on changes in stock size.

31 Analytical retrospective analyses are routinely done for both stocks. Historical retrospective, where there
32 are successive accepted assessments, is also informative and should be done to indicate how consistent
33 the assessments have been over time. Simpler models, e.g. like Robin Cook's or surplus production
34 models should be investigated. It is not necessary to go to Ensemble modeling, but looking at more than
35 one modeling framework might be informative.

36 Should data from the IPHC longline survey be used in either assessment?

37 During the review meeting it was not clear if the raw data received from the International Pacific Halibut
38 Commission (IPHC) had been treated appropriately to derive an index of stock size. Further work was
39 conducted by the AFSC survey unit during the meeting, and it seems that the series shown in the excel
40 spreadsheet "Survey index comparison (trawl surveys, longline surveys).xlsx" could be treated as an
41 index of stock size. The appropriateness of the data and how it was treated to calculate an index should be

42 further verified between now and the assessment meeting later in the year. What data were used and how
43 they were used should also be documented.

44 The IPHC longline survey has been conducted every year since 1997. The survey covers the Eastern
45 Bering Sea area well (Figure 1) but it is not clear if all stations were used in calculating a relative index or
46 if only those on the slope were used.

47 The main index of abundance used in the Eastern Bering Sea stock assessment is the AFSC shelf trawl
48 survey. The agreement between the AFSC shelf trawl survey and the IPHC longline survey is not very
49 good (Figure 2). This could be due to different size selectivities and / or inherent variability in the data.

50 The IPHC longline survey sample larger individuals (Figure 3) and a lag between the two indices would
51 therefore be expected. However, the sudden decreases in the relative index in 1999 and 2005, and
52 similarly sudden increase in the following year are unlikely to reflect real changes in stock sizes. These
53 anomalies warrant further investigations to try to identify what might cause them. If there are valid
54 reasons to exclude those two points, it might be possible to reconcile the IPHC longline and AFSC shelf
55 trawl survey time series taking into account that they sample different size groups.

56 Including longline surveys (IPHC or AFSC) in the assessment might alleviate concerns that the shelf
57 trawl survey samples poorly larger sizes, either because large Pacific cod are outside of the surveyed area
58 or because they are able to swim faster than the fishing gear and therefore escape capture. Including one
59 or more indices of stock sizes for larger fish sizes therefore has the potential to improve the assessment
60 and reduce the uncertainty in the population estimates of larger fish sizes.

61 Figure 4 shows the spawning stock biomass (SSB) trends for the Eastern Bering Sea Pacific cod stock
62 from various model configurations. I have not been able to find a model where only the IPHC longline
63 survey is added to the AFSC shelf trawl survey, but I think one was presented during the meeting. My
64 memory is that including data from the IPHC longline survey in the assessment implies lower terminal
65 year biomass than when only the shelf trawl survey is used. Figure 4, however, shows that when both the
66 IPHC longline survey and the AFSC longline survey are added (model 15.6B), the SSB estimates are the
67 lowest of the model considered. Adding the AFSC slope trawl survey (model 15.6C) implies essentially
68 identical results to adding the two longline surveys. Adding only the AFSC longline survey (model
69 15.6A) results in a SSB trend that is markedly different from those of the other models considered.

70 From what was discussed during the meeting and the documentation reviewed, there are no objective
71 reasons to reject the IPHC longline survey as an index of stock size, assuming it has been correctly put
72 together and calculated. Its influence on the assessment results, however, when used along with the AFSC
73 longline survey is puzzling.

74 The IPHC longline survey data should be thoroughly investigated. It should be used in the assessment
75 unless fatal flaws in the data, in the treatment of the data or in the survey methodology are identified.

76 Should data from the NMFS longline survey be used in either assessment?

77 The AFSC longline survey was developed as an index of abundance for sablefish, but recently, the results
78 have also been found useful as indices of abundance for roughey rockfish, blackspotted rockfish, and for
79 black halibut (aka Greenland turbot). Similar to the IPHC longline survey, the AFSC longline survey
80 started in 1997, but it is conducted during odd years in the Eastern Bering Sea and during even years in
81 the Aleutian Islands. In the Gulf of Alaska, the survey is conducted every year. There are few stations in
82 the Eastern Bering Sea and in the Aleutian Islands, but they do cover the expected area of distribution of
83 larger Pacific cod (Figure 5).

84 The agreement between the AFSC shelf trawl survey and the AFSC longline survey (Figure 6) is better
85 than between the AFSC shelf trawl survey and the IPHC longline survey. The AFSC longline survey,
86 similar to the IPHC longline survey, catch different sizes than the AFSC shelf trawl survey (Figure 7).

87 The agreement between the AFSC longline survey and the IPHC longline survey (Figure 8) is poor
88 overall. The two apparently anomalous points in the IPHC longline survey in 1999 and 2005 may explain
89 in part the discrepancy, but differences in the area surveyed, in the timing of the survey and slight
90 differences in the size composition may also play a role (keeping in mind that further work may be
91 needed to confirm that the index derived from the IPHC longline survey is appropriate).

92 Similar to the IPHC longline survey, there are no objective reasons to reject the AFSC longline survey as
93 an index of stock size. The AFSC longline survey should also be thoroughly investigated and used in the
94 assessment unless fatal flaws in the data, in the data treatment or in the survey methodology are
95 identified.

96 It is only by analyzing additional data that confidence will increase in the model results. Given the widely
97 different results that can be obtained with SS3 (Figure 4), and the volatility of some of those results, it is
98 not be possible to model oneself out of the uncertainty. Only careful examination and inclusion of
99 informative additional data will allow that.

100 The discussion above is based on examination of data from the Eastern Bering Sea surveys, but the
101 conclusions and recommendations also hold for the Aleutian Islands data and assessment.

102 How should the various data sets be weighted?

103 Stock Synthesis is a very flexible stock assessment framework. Giving different weights to the various
104 data sources, and depending on assumptions (e. g. fixed parameters), very different results in terms of
105 absolute stock size, but also sometimes in terms of trends, can be obtained (Figure 4). This can also occur
106 with other assessment frameworks, but because of SS3's flexibility, the problem is more severe.

107 Generally speaking, indices of abundance should be given more weight in the assessment than length
108 composition. Age composition, particularly from surveys or other indices of abundance can be very
109 informative if analyzed and used appropriately. Information in the length composition is at best indirect
110 information on changes in stock size and it may be misleading if substantial changes in growth occur over
111 time (Figure 11). In almost every stock where growth information is available by year, growth has been
112 found to vary with trends over time, sometimes quite considerably. SS3 does allow for time varying
113 growth, but without external information, it is unlikely to be able to estimate changes in growth correctly.

114 What form (i.e., Stock Synthesis “pattern”) should be used for the selectivity functions?

115 Selectivity is a very important parameter in any assessment framework. Changes in growth, natural
116 mortality, or fishing mortality can all be aliased as changes in selectivity. Several of the model
117 configurations examined during the review had very peculiar selectivity patterns. This was identified in
118 the pre-review material and in the presentations by the assessment team. Those were probably not real
119 and were likely due to sampling problems or aliasing other changes. My preference would be to NOT
120 allow too much flexibility in selectivity changes over time, and to NOT allow strange patterns (e.g.
121 figures 2.1.3 in the Eastern Bering Sea and 2A.11 and 2A.12 in the Aleutian Islands in the December
122 2015 SAFE report). If allowing these strange patterns is a condition of getting a good fit or convergence,
123 this would be a sign that something else might be wrong. If allowed to change over time and age, the
124 changes should be relatively smooth and not result in peculiar patterns.

125 During the meeting, Robin Cook noted that the ratio of catch at length in the longline commercial fishery
126 to the survey catch at length is reasonably constant above a certain length. A possible interpretation is that
127 domed selectivity estimated for the AFSC shelf survey could be an artifact. Data in the file "Long-term
128 sizecomp comparison (trawl surveys, longline surveys, longline fishery).xlsx", the ratio of the longline
129 commercial catch at length to the various population estimates at length from the surveys are plotted in
130 Figure 9. The ratio of the longline commercial catch at length to the ASFC shelf survey for the Eastern
131 Bering Sea is consistent with the observed size composition (Figure 10). The longline commercial fishery
132 catches very few Pacific cod less than 40cm and the ratio increases progressively from nearly zero at 40
133 cm to around 5-6 at 70 cm and the ratio is indeed relatively constant from 70cm or so. The ratio seems to
134 decrease above 100cm but this could easily be the result of low sample size. The ratio being relatively
135 constant at 70cm and above suggests that selectivity does not decrease at those sizes in the AFSC shelf
136 trawl survey, or that selectivity in the longline commercial fishery decreases at a similar rate. This is
137 unlikely but not impossible. The link between selectivity in the AFSC shelf trawl survey and selectivity in
138 the longline commercial fishery should be further investigated to guide modeling.

139 The longline surveys appear to have very low selectivity, lower than that of the commercial longline
140 fishery (Figure 9), for size less than the 55 cm for the IPHC longline survey and less than 40 cm for the
141 AFSC longline survey. Differences in selectivity between the two longline surveys may be due to the
142 differences in the sizes they catch (Figure 7). The ratio of the commercial longline catch at length to the
143 survey catch at length is near 1 for both surveys in the 60-70 cm range, but the ratios diverge thereafter.
144 The IPHC longline survey appears to have higher selectivity for the larger size than the commercial
145 longline fishery does, while the AFSC longline survey would have lower selectivity than the commercial
146 longline fishery. The AFSC slope trawl survey shows a pattern similar to the AFSC longline survey. The
147 reason(s) for the apparent differences in selectivity between the IPHC longline survey and the AFSC
148 longline survey for lengths above 70cm should be further investigated.

149 This is but a quick examination of what the data are telling us in terms of selectivity. Modeling results
150 would be expected to be consistent with those observations.

151 Should the models be structured with respect to season?

152 In both areas, there seem to be a strong seasonal pattern in the fishery. Therefore, where the data are
153 sufficient, it would be appropriate to structure the assessment model by season. However, for the Aleutian
154 Islands assessment, the data may not be sufficient to structure by season.

155 Should the models be structured with respect to gear type?

156 Bottom trawl and longline are the two main gear types in the fisheries. Their size selectivities are
157 expected to be different, and the models should definitely be structured with respect to gear type where
158 the data are sufficient to do so.

159 How much time variability should be allowed, and in which parameters?

160 Selectivity, catchability of the surveys, natural mortality, and growth could be allowed to vary over time
161 when there is independent information supporting that changes are happening. A change in the ratio
162 between total catch biomass and biomass estimate in the survey that could not be explained by changes in
163 management could be an indication that the catchability in the survey has changed. Changes in mesh sizes
164 in the trawl or hook size in the longline fishery could be an indication of a stepped change in selectivity.
165 Changes in the predator field or extreme weather events could be indications of changes in natural
166 mortality. Because most of these parameters are interlinked, great care should be taken in allowing them
167 to vary. Only those parameters where there is external information suggesting that changes are occurring

168 should be allowed to vary, probably one at a time to avoid incorrect interpretation. Because of the
169 flexibility in SS3 and because most of these parameters are interlinked allowing them to change may give
170 strange results, such as highly anomalous selectivity.

171 What constraints, if any, should be placed on survey selectivity at older ages?

172 Peculiar selectivity patterns have been identified as a problem in the presentations by the assessment
173 team. Based on the information in Figure 9, the selectivity for the AFSC shelf trawl survey in the Eastern
174 Bering Sea at ages corresponding to 70 cm and larger would be expected to be reasonably flat. For both
175 areas, sharp peak and valleys, unless based on external information, should be smoothed. As indicated
176 above, strange, irregular patterns should be constrained to be smoother.

177 What constraints, if any, should be placed on survey catchability?

178 In the mid 1980s survey catchability was estimated for cod and haddock on the Eastern Scotian Shelf and
179 in the Gulf of St. Lawrence in Eastern Canada. Catchability for cod at the time was about 0.5 and for
180 haddock it was close to 1. Vessels and gears have changed since and catchability estimates have also
181 changed and in some areas they are now estimated to be greater than 1. Survey catchability is a scaling
182 factor. In most assessment in the ICES area, survey catch per tow or catch per hour are used in the
183 assessments and survey catchability is not an issue. It is, however, good practice to check every now and
184 then if the assessment has the units more or less right by doing the areal expansion and comparing with
185 the population estimates in the assessment.

186 Survey catchability smaller than 1 are relatively easy to rationalize e.g. by fish swimming faster than the
187 net is towed, escaping above or below the net, or being more abundant in areas that are not surveyed.
188 Survey catchability greater than 1 would happen if there is herding or if fish density in the surveyed area
189 is expanded to areas where there are no fish, e.g. expanding flatfish density estimates from samples on
190 smooth flatfish habitat to rough hard substrate that are not sampled and where flatfish are not present.

191 Catchability and natural mortality are interlinked. Considerable work was done in the Gulf of St.
192 Lawrence in eastern Canada following the collapse of the groundfish stocks. Sinclair (20014) estimated
193 that natural mortality had likely increased for the southern Gulf of St. Lawrence stock. Subsequent stock
194 assessments (e.g. [http://www.dfo-mpo.gc.ca/csas-sccs/Publications/ResDocs-](http://www.dfo-mpo.gc.ca/csas-sccs/Publications/ResDocs-DocRech/2007/RES2007_033_B.pdf)
195 [DocRech/2007/RES2007_033_B.pdf](http://www.dfo-mpo.gc.ca/csas-sccs/Publications/ResDocs-DocRech/2007/RES2007_033_B.pdf) and [http://www.dfo-mpo.gc.ca/csas-sccs/Publications/ResDocs-](http://www.dfo-mpo.gc.ca/csas-sccs/Publications/ResDocs-DocRech/2007/RES2007_068_B.pdf)
196 [DocRech/2007/RES2007_068_B.pdf](http://www.dfo-mpo.gc.ca/csas-sccs/Publications/ResDocs-DocRech/2007/RES2007_068_B.pdf)) have used time varying natural mortality, but Canadian scientists
197 warned that "Estimation of M can be confounded by changes in survey catchability and fishery catch
198 reporting, and may be sensitive to assumptions and constraints applied in the ADAPT estimation
199 procedure." (http://www.dfo-mpo.gc.ca/csas/Csas/status/2007/SAR-AS2007_002_E.pdf). Therefore,
200 estimating catchability and natural mortality simultaneously would be challenging in the absence of
201 external information.

202 External information indicative of changes in catchability could be changes in gears in the surveys or
203 changes in predator abundance. Changes in catchability of pelagic species has been hypothesized to
204 explain apparent increases of small pelagics in Eastern Canada after the collapse of groundfishes but this
205 has been challenged. Catchability in longline surveys could occur if high prey abundance in the water
206 decreases the attractiveness of baited hooks.

207 This being said, Pacific cod appears to be a relatively well behaved species as far as trawl surveys are
208 concerned. Survey catchability estimates between 0.5 and 1.5 would not seem to be cause for concern.
209 The assessment team, the PDT and the SSC are concerned that catchability less than 1 imply very large
210 biomass estimates. As indicated above, I do not share that concern (within limits of course). Catchability

211 of the trawl survey in the Aleutian Islands area would be expected to be more uncertain than in the
212 Eastern Bering Sea area because bottom topography is likely rougher and more diverse in the Aleutian
213 Islands area than in the Eastern Bering Sea area.

214 How should large gradients be dealt with in otherwise apparently converged models?

215 Stock Synthesis User Manual version 3.24s, page 27, states: "When using more population length bins
216 than data bins, SS will run slower (more calculations to do), the calculated weights at age will be less
217 aliased by the bin structure, and you may or may not get better fits to your data.

218 While exploring the performance of models with finer bin structure, a potentially pathological situation
219 has been identified. When the bin structure is coarse (note that some applications have used 10 cm bin
220 widths for the largest fish), it is possible for a selectivity slope parameter or a retention parameter to
221 become so steep that all of the action occurs within the range of a single size bin. In this case, the model
222 will lose the gradient of the logL with respect to that parameter and convergence will be hampered. A
223 generic guidance to avoid this situation is not yet available."

224 I have no further advice on how to deal with large gradient than what is said in the Stock Synthesis User
225 Manual.

226 Changes in growth

227 For the Eastern Bering Sea, weights at age in the survey (from the preliminary assessment data file) show
228 trends over time that seem to be year-class specific. It could be worth investigating further changes in
229 growth (Figure 11), particularly with respect to the implications for the assessment as growth changes
230 may have an influence on fishing mortality and population estimates.

231 Recruitment index

232 For the Eastern Bering Sea, the population estimates in the AFSC shelf trawl survey seem to be
233 reasonably consistent for the first 3 age groups or so with reasonably good year-class tracking (Figure
234 12). If the AFSC shelf trawl survey for the Eastern Bering Sea is indeed following year-classes
235 reasonably well, it could provide at least 3 successive estimates of year-class size and this could be used
236 to obtain reasonably reliable estimates of year-class sizes.

237 In my cursory comparison of the AFSC shelf trawl survey length frequencies with the age frequencies in
238 the same survey, I got the impression that the smallest modal length group was sometimes aged as age 1
239 and in other cases as age zero. This should be verified.

240 Exploitation rate

241 The ratio of the commercial catch in tons to the survey biomass estimate in tons should be an indication
242 of exploitation rate (relative if the catch and survey biomass are not in the same units). Figure 13, using
243 data from run 15.6 for the Eastern Bering Sea shows the catch/survey ratio in biomass compared with the
244 fishing mortality estimate for the same model. The results suggest that fishing mortality in model 15.6
245 could be overestimated in recent years. Unless the catchability of the survey has changed, the results
246 below suggest that F has been lower than average since about 2007. The correlation between the
247 catch/survey ratio and F is low (0.009).

248 Reliability of total catch estimates

249 For the Aleutian Islands assessment model 15.7, there is reasonably good agreement between fishing
250 mortality estimates in the assessment and catch (Figure 14) except in the late 1980s and in 2010 when
251 fishing mortality estimates suggests that mortality has been higher. It might be worth investigating if
252 additional sources of mortality (e.g. increased M) occurred in those years. The correlation between F from
253 the assessment (model 15.7) and the ratio of catch to survey biomass is higher (Figure 15) for the
254 Aleutian Islands (0.47).

255 Stock structure

256 In the Aleutian Islands area, it is unlikely that there is a single stock in the traditional understanding of the
257 concept. Instead, a number of local spawning would be expected with limited mixing during spawning.
258 While these different spawning units may react similarly to changes in the environment and show similar
259 trends in recruitment, they are unlikely to form a single homogeneous biological unit. It is likely
260 impractical to do individual stock assessments for each of the individual units, and lumping all units into
261 a single assessment with indices of abundance for only a few of them may increase the risk to less
262 productive units. Simpler form of monitoring and management, in close cooperation with the industry and
263 possibly NGOs, could be a better way of protecting the resources and managing the fisheries.

264 Conclusions and Recommendations in accordance with the ToRs

265 For the IPHC longline survey, the appropriateness of the data and how it was treated to calculate an index
266 should be further verified between now and the assessment meeting later in the year. What data were used
267 and how they were used should also be documented. The apparent anomalies in 1999 and 2005 warrant
268 further investigations to try to identify what might cause them. If there are valid reasons to exclude those
269 two points, it might be possible to reconcile the IPHC longline and AFSC shelf trawl survey time series
270 taking into account that they sample different size groups. From what was discussed during the meeting
271 and the documentation reviewed, there are no objective reasons to reject the IPHC longline survey as an
272 index of stock size, assuming it has been correctly put together and calculated. The IPHC longline survey
273 data should be thoroughly investigated. It should be used in the assessment unless fatal flaws in the data,
274 in the treatment of the data or in the survey methodology are identified.

275 Similar to the IPHC longline survey, there are no objective reasons to reject the AFSC longline survey as
276 an index of stock size. The AFSC longline survey should also be thoroughly investigated and used in the
277 assessment unless fatal flaws in the data, in the data treatment or in the survey methodology are
278 identified.

279 The discussion above is based on examination of data from the Eastern Bering Sea surveys, but the
280 conclusions and recommendations also hold for the Aleutian Islands data and assessment.

281 With respect to weighting different data sets, indices of abundance should be given more weight in the
282 assessment than length composition. Age composition, particularly from surveys or other indices of
283 abundance can be very informative if analyzed and used appropriately. Information in the length
284 composition is at best indirect information on changes in stock size and it may be misleading if
285 substantial changes in growth occur over time (Figure 11).

286 Regarding the form of the selectivity function, my preference would be to NOT allow too much flexibility
287 in selectivity changes over time and to NOT allow strange patterns (e.g. figures 2.1.3 in the Eastern
288 Bering Sea, and 2A.11 and 2A.12 in the Aleutian Islands in the December 2015 SAFE report). If allowing
289 these strange patterns is a condition of getting a good fit or convergence, this would be a sign that

290 something else might be wrong. If allowed to change over time and age, the changes should be relatively
291 smooth and not result in peculiar patterns. The ratio of catch at length in the longline commercial fishery
292 to the survey catch at length being relatively constant at 70cm and above suggests that selectivity does not
293 decrease at those sizes in the AFSC shelf trawl survey, or that selectivity in the longline commercial
294 fishery decreases at a similar rate. This is unlikely but not impossible. The link between selectivity in the
295 AFSC shelf trawl survey and selectivity in the longline commercial fishery should be further investigated
296 to guide modeling. The reason(s) for the apparent differences in selectivity between the IPHC longline
297 survey and the AFSC longline survey for lengths above 70cm should be further investigated.

298 Where the data are sufficient, it would be appropriate to structure the assessment model by season.

299 Bottom trawl and longline are the two main gear types in the fisheries. Their size selectivity are expected
300 to be different and the models should definitely be structured with respect to gear type where the data are
301 sufficient to do so.

302 Selectivity, catchability of the surveys, natural mortality, and growth could be allowed to vary over time
303 when there is independent information supporting that changes is happening. Because most of these
304 parameters are interlinked, great care should be taken in allowing them to vary. Only those parameters
305 where there is external information suggesting that changes is occurring should be allowed to vary,
306 probably one at a time to avoid incorrect interpretation.

307 Based on the information in Figure 9, the selectivity for the AFSC shelf trawl survey in the Eastern
308 Bering Sea at ages corresponding to 70 cm and larger would be expected be reasonably flat. For both
309 areas, sharp peaks and valleys, unless based on external information, should be smoothed. As indicated
310 above, strange, irregular patterns should be constrained to be smoother.

311 I have no further advice on how to deal with large gradient than what is said in the Stock Synthesis User
312 Manual.

313 It could be worth investigating further changes in growth (Figure 11), particularly with respect to the
314 implications for the assessment as growth changes may have an influence on fishing mortality and
315 population estimates.

316 If the AFSC shelf trawl survey for the Eastern Bering Sea is indeed following year-classes reasonably
317 well, it could provide at least 3 successive estimates of year-class size and this could be used to obtain
318 reasonably reliable estimates of year-class sizes.

319 In my cursory comparison of the AFSC shelf trawl survey length frequencies with the age frequencies in
320 the same survey, I got the impression that the smallest modal length group was sometimes aged as age 1
321 and in other cases as age zero. This should be verified.

322 Figure 13, using data from run 15.6 for the Eastern Bering Sea shows the catch/survey ratio in biomass
323 compared with the fishing mortality estimate for the same model. The results suggest that fishing
324 mortality in model 15.6 could be overestimated in recent years.

325 For the Aleutian Islands assessment model 15.7, there is reasonably good agreement between fishing
326 mortality estimates in the assessment and catch (Figure 14) except in the late 1980s and in 2010 when
327 fishing mortality estimates suggests that mortality has been higher. It might be worth investigating if
328 additional sources of mortality (e.g. increased M) occurred in those years.

329 In the Aleutian Islands area, it is unlikely that there is a single stock in the traditional understanding of the
330 concept. Simpler form of monitoring and management, in close cooperation with the industry and
331 possibly NGOs, could be a better way of protecting the resources and managing the fisheries.

332 One cannot model oneself out of lack of data, particularly for the Aleutian Islands assessment. Stock
333 Synthesis has so much flexibility that, given sufficient time, a skilled user can probably get almost any
334 stock trend from a dataset. Indices of abundance should be given more weight in the assessment than
335 length composition. Age composition, particularly from the commercial fishery, but also from surveys or
336 other indices of abundance can be very informative if analyzed appropriately. Information in the length
337 composition is at best indirect information on changes in stock size. In almost every stock where growth
338 information is available by year, growth has been found to vary with trends over time, sometimes quite
339 considerably and this could very well be the case here for the Eastern Bering Sea (Figure 11). SS3 does
340 allow for time varying growth, but without external information, it is unlikely to be able to estimate
341 changes in growth correctly.

342 Analytical retrospective analyses are routinely done for both stocks. Historical retrospective, where the
343 successive accepted assessment are also informative and should be done to indicate how consistent the
344 assessments have been over time.

345 Simpler models, e.g. like Robin Cook's or surplus production models should be investigated. It is not
346 necessary to go to Ensemble modeling, but looking at more than one modeling framework might be
347 informative.

348