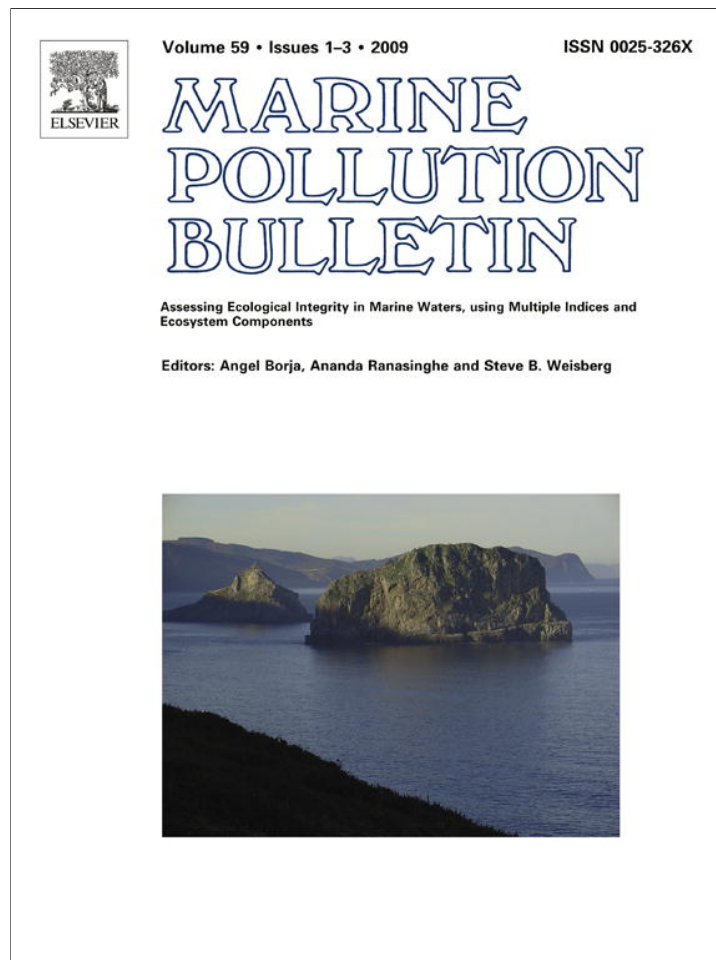


Provided for non-commercial research and education use.
Not for reproduction, distribution or commercial use.



This article appeared in a journal published by Elsevier. The attached copy is furnished to the author for internal non-commercial research and education use, including for instruction at the authors institution and sharing with colleagues.

Other uses, including reproduction and distribution, or selling or licensing copies, or posting to personal, institutional or third party websites are prohibited.

In most cases authors are permitted to post their version of the article (e.g. in Word or Tex form) to their personal website or institutional repository. Authors requiring further information regarding Elsevier's archiving and manuscript policies are encouraged to visit:

<http://www.elsevier.com/copyright>



Contents lists available at ScienceDirect

Marine Pollution Bulletin

journal homepage: www.elsevier.com/locate/marpolbul

Calibration and evaluation of five indicators of benthic community condition in two California bay and estuary habitats

J. Ananda Ranasinghe^{a,*}, Stephen B. Weisberg^a, Robert W. Smith^{*,†}, David E. Montagne^b, Bruce Thompson^c, James M. Oakden^d, David D. Huff^e, Donald B. Cadien^f, Ronald G. Velarde^g, Kerry J. Ritter^a

^a Southern California Coastal Water Research Project, 3535 Harbor Blvd, Suite 110, Costa Mesa, CA 92626, USA

^b P.O. Box 2004, Penn Valley, CA 95946, USA

^c San Francisco Estuary Institute, 7770 Pardee Lane, Oakland, CA 94621, USA

^d Moss Landing Marine Laboratory, Moss Landing, CA 95039, USA

^e Dept. of Fisheries, Wildlife and Conservation Biology, University of Minnesota, St. Paul, MN, USA

^f County Sanitation Districts of Los Angeles County, P.O. Box 4998, Whittier, CA 90607, USA

^g City of San Diego, Marine Biology Laboratory, 2392 Kincaid Rd., San Diego, CA 92101, USA

ARTICLE INFO

Keywords:

Bay and estuary benthos
Benthic community condition
Benthic index performance comparison
Benthic index combinations
Southern California
San Francisco Bay

ABSTRACT

Many types of indices have been developed to assess benthic invertebrate community condition, but there have been few studies evaluating the relative performance of different index approaches. Here we calibrate and compare the performance of five indices: the Benthic Response Index (BRI), Benthic Quality Index (BQI), Relative Benthic Index (RBI), River Invertebrate Prediction and Classification System (RIVPACS), and the Index of Biotic Integrity (IBI). We also examine whether index performance improves when the different indices, which rely on measurement of different properties, are used in combination. The five indices were calibrated for two geographies using 238 samples from southern California marine bays and 125 samples from polyhaline San Francisco Bay. Index performance was evaluated by comparing index assessments of 35 sites to the best professional judgment of nine benthic experts. None of the individual indices performed as well as the average expert in ranking sample condition or evaluating whether benthic assemblages exhibited evidence of disturbance. However, several index combinations outperformed the average expert. When results from both habitats were combined, two four-index combinations and a three-index combination performed best. However, performance differences among several combinations were small enough that factors such as logistics can also become a consideration in index selection.

© 2008 Elsevier Ltd. All rights reserved.

1. Introduction

Index-based approaches to summarizing data have facilitated the use of benthic infauna as indicators of sediment condition in marine and estuarine environments (Hyland et al., 1999; Bergen et al., 2000; Dauer et al., 2000; Summers, 2001; Hyland et al., 2003; Diaz et al., 2004; Borja and Dauer, 2008). While reducing complex biological data to a single value has disadvantages, the resulting indices remove much of the subjectivity associated with interpreting data. The indices also provide a simple means for communicating complex information to managers and for correlating benthic responses with stressor data (Dauer et al., 2000; Hale et al., 2004; Bilkovic et al., 2006).

There have been a number of approaches to creating benthic indices (Diaz et al., 2004; Borja and Dauer, 2008). Some integrate information at the community level and rely on parameters such as abundance, diversity, functional feeding groups, and depth beneath the sediment surface (Weisberg et al., 1997; Engle and Summers, 1999; Van Dolah et al., 1999; Diaz et al., 2004). Other indices focus on species composition, comparing sample composition to an expected species mix or quantifying the average pollution tolerance of species found at the site (Borja et al., 2000; Hawkins et al., 2000; Smith et al., 2001, 2003; Leung et al., 2005; Van Sickle et al., 2006). Although community-level approaches often include measures of sensitive and tolerant biota, these measures are usually based on just a few indicator organisms, while species composition indices include many taxa.

Despite the broad range of benthic index approaches, there have only been a few comparisons of performance among benthic indices (Ranasinghe et al., 2002; Labruno et al., 2006; Quintino et al., 2006; Borja et al., 2007; Zettler et al., 2007; Blanchet et al.,

* Corresponding author. Fax: +1 (714)755 3299.

E-mail address: AnandaR@sccwrp.org (J.A. Ranasinghe).

† Deceased.

2008; Borja et al., 2008; Puente et al., 2008; Teixeira et al., 2008). Most of these were limited to comparing just a few indices, did not compare community-level and species composition indices, or did not include a means to evaluate which performed best. As a result, there are no widely accepted generalizations about the relative efficacy of indices at these different levels of organization.

In this study, we evaluate the performance of five benthic indices that rely on different sets of community or species composition measures, comparing their site assessments to the professional judgment of nine benthic experts. Three of the indices were previously developed and applied in California bays, while the other two were developed in other habitats or geographic regions, but were considered to have potential for success. The five index approaches were (i) the Relative Benthic Index (RBI; Hunt et al., 2001), (ii) the Index of Biotic Integrity (IBI; Thompson and Lowe, 2004), (iii) the Benthic Response Index (BRI; Smith et al., 2001, 2003), (iv) the River Invertebrate Prediction and Classification System (RIVPACS; Wright et al., 1993; Van Sickle et al., 2006), and (v) the Benthic Quality Index (BQI; Rosenberg et al., 2004). The RBI and IBI are based on community measures, the BRI and RIVPACS on species composition, and the BQI on both. The comparisons were conducted in two ecologically and geographically distinct habitats: (a) the marine bays of southern California and (b) polyhaline San Francisco Bay. The objective was to evaluate the relative performance of these indices alone and in combination in each habitat, in relation to assessments by nine benthic experts.

2. Methods

The performance of the five benthic indices was evaluated in four steps:

- (i) Data for sampling sites in each of the two habitats were identified, acquired, and adjusted to create consistency across sampling programs.
- (ii) The five benthic indices were calibrated using a common set of data for all indices.
- (iii) Threshold values were selected for each index to assess benthic condition on a four-category scale.
- (iv) Performance of the indices, and all possible index combinations, was evaluated by applying them to independent data and comparing the condition assessments to that of nine benthic experts.

2.1. Data

Data from projects that collected benthic species abundance and sediment chemistry data synoptically from marine bays in southern California and polyhaline San Francisco Bay (Table 1)

Table 1
Data sources for calibration and validation samples.

Habitat (sampling methods)	Project	Period	Reference	No. of samples	
				Calibration	Validation
Southern California Marine Bays. (0.1-mm sieve; 0.1 m ² sample area)	Bight'98	1998	Ranasinghe et al. (2003)	107	5
	Bight'03	2003	Ranasinghe et al. (2007)	110	10
	San Diego TMDL	2001–2002	SCCWRP and SPAWAR (2004), Brown and Bay (2005)		4
	EMAP	1999	U.S. EPA (2004)	21	5
	Total			238	24
Polyhaline central San Francisco Bay. (0.5-mm sieve; 0.05 m ² sample area)	EMAP	2000	U.S. EPA (2004)	22	1
	BADA	1994–1997	Bay Area Dischargers Association (1994)	42	2
	BPTCP	1994, 1997	Hunt et al. (2001)	16	4
	RMP	1994–2000	Thompson et al. (1999)	45	4
	Total			125	11

were identified, acquired, evaluated for methodological consistency, normalized for units of measure, and assembled into a database. Data about habitat conditions such as depth, bottom water salinity, sediment grain-size distributions, and acute toxicity to amphipods were included, if available.

Only benthic data from samples sieved through the most frequently used screen sizes were included: 1-mm sieve data for southern California marine bays and 0.5-mm sieve data for polyhaline San Francisco Bay. Taxonomic inconsistencies among programs were eliminated by cross-correlating the species lists, identifying differences in nomenclature, and resolving discrepancies by consulting the taxonomists from each program. Species abundances were normalized to the most frequently occurring sample area by combining data from small samples or adjusting abundances to 0.1 m² in southern California marine bays and 0.05 m² in polyhaline San Francisco Bay.

Ninety percent of the available data was used to calibrate benthic indices while the remainder was set aside for evaluation (Table 1). Samples for evaluation were selected by ordering the data in each habitat by the mERMq (Long et al., 2000, 2006) and systematically selecting sites from within quartile groups in each habitat. While it is generally accepted that current models of benthic response do not discriminate between chemical contamination and other sources of stress (Borja et al., 2003), this approach ensured that a range of benthic conditions were represented in the calibration and evaluation data.

An additional subset of the calibration data was set aside to select index threshold values. Similar to selecting evaluation samples, the subset of 35 samples from southern California and 33 samples from San Francisco Bay was selected by ordering the calibration data in each habitat by the mERMq and systematically selecting sites within quartile groups in each habitat.

2.2. Benthic index calibration

All the indices, other than the BQI, have previously been calibrated, validated and used successfully in California (Hunt et al., 2001; Smith et al., 2001, 2003; Ranasinghe et al., 2004; Thompson and Lowe, 2004), although RIVPACS has been used only in freshwater streams (e.g., Rehn et al., 2007). The BQI was previously calibrated and used in Europe (e.g., Rosenberg et al., 2004). Our index calibration involved applying these previous calibration procedures to data from the southern California marine bays and polyhaline San Francisco Bay. Each index was calibrated separately for each habitat.

2.2.1. Benthic Response Index (BRI)

We calibrated the Benthic Response Index (BRI) using the methods of Smith et al. (2001, 2003), with slight variations in the first and third of their four steps. The first step in BRI calibration is iden-

tifying a disturbance (or pollution) vector in an ordination space to facilitate calculation of species tolerance scores based on the distribution of species abundances along the vector. The BRI (Smith et al., 2001) was originally developed offshore, where a well-understood gradient of point-source disturbance allowed a disturbance vector to be identified from *a priori*-selected disturbed and undisturbed sites. Such simple disturbance gradients do not exist in bays and estuaries because there are many types of disturbance, a number of contaminant sources and circulation patterns that often redistribute contaminants throughout the system. Therefore, the BRI disturbance vector was selected using the vector with the maximum value for $T = R_{MSR} - R_{NSP}$ where R_{MSR} is the Spearman rank correlation between vector position and the observation mean species range (MSR) and R_{NSP} is the Spearman rank correlation between vector position and the observation number of species (Table 2). The MSR quantifies the average species range along the disturbance vector for the species occurring at a site. The range for each species was calculated as the difference between the last and first occurrence on the disturbance gradient; the MSR for a site is the average of the ranges for the species occurring at that site. We identified the disturbance vector by creating test vectors in the ordination space using an optimizing algorithm and selecting the vector with the highest value for T . The R_{NSP} computations excluded observations toward the undisturbed end of the vector to prevent the use of observations that might be to the left of a Pearson–Rosenberg species diversity peak. Species diversity would be negatively correlated with the disturbance gradient to the right of the diversity peak, leading to the negative sign for R_{NSP} .

The second BRI calibration step was application of an optimization procedure to determine data transformations to be used in subsequent computations (see Smith et al., 2001, 2003). Tolerance scores were calculated for abundance transformations with exponents (e in the tolerance score equation) of 0, 0.25, 0.33, 0.5, and 1.0 in combination with BRI calculations using transformations with exponents (f in the BRI equation) of 0.25, 0.33, 0.5, and 1.0. The combination with the highest Spearman correlation between optimized index values and the disturbance vector was used in each habitat (Table 3).

The third BRI calibration step selects t , the maximum number of occurrences used for species tolerance score calculations. Where previous versions of the BRI optimized the same maximum number of occurrences for all species in a habitat, we customized values for each species with the objective of including low abundances in tolerance score calculations only if they contribute signal, rather than noise. The t value yielding the highest Spearman correlation between optimized index values and the disturbance vector was selected using another iterative optimization procedure. We used maximum occurrence values from iterations with Spearman correlations of 0.937 and 0.957 between the disturbance vector and the

Table 2

Spearman correlation coefficients between the vector in the ordination space selected to represent the disturbance gradient and (a) the mean species range and (b) the number of species in each habitat. The disturbance vector was selected by generating test vectors using an optimization procedure and selecting the vector that maximized the value of T . The mean species range is the average of the species ranges along the disturbance vector for the species occurring at each sampling site (see text).

	Southern California Marine Bays	Polyhaline Central San Francisco Bay
Spearman correlation with mean species range (R_{MSR})	0.9182	0.9007
Spearman correlation coefficient with number of species (R_{NSP})	−0.8457	−0.8632
$T = R_{MSR} - R_{NSP}$	1.7639	1.7639

Table 3

Optimum parameter values for exponents in the Benthic Response Index (BRI) equation for each habitat. The exponent f is used for index calculations, while e is used to develop species tolerance (p_i) values (see text and Smith et al., 2001, 2003).

	Southern California Marine Bays	Polyhaline Central San Francisco Bay
e	0.25	0.33
f	0	0
Spearman correlation coefficient between the optimized index and the disturbance vector	0.903	0.944

occurrence adjusted index values in southern California marine bays and polyhaline San Francisco Bay, respectively.

In the final step, pollution tolerance scores were calculated for species occurring in two or more samples in each habitat as the position of the weighted-average of the abundance distribution on the disturbance vector. Tolerance values were calculated for 460 species in southern California marine bays and 154 species in polyhaline San Francisco Bay. Higher BRI values are associated with higher pollution levels.

2.2.2. Benthic Quality Index (BQI)

We calibrated the Benthic Quality Index (BQI) for each habitat using the method of Rosenberg et al. (2004). First, for each sample in the calibration data, the expected number of species for a subset of 50 individuals was calculated as

$$ES50_k = \sum_{i=1}^s \left[1 - \frac{(N_k - N_{ki})!(N_k - 50)!}{(N_k - N_{ki} - 50)!N_k!} \right],$$

where s is the number of species in sample k , N_k is the total abundance of all species in sample k , and N_{ki} is the abundance of species i in sample k . Next, species tolerance scores, $ES50_{0.05i}$, were computed for species that were found in at least three samples in each habitat as the 5th percentile of the distribution of expected numbers of species for the samples in which the species occurred. Tolerance scores were calculated for 346 species in southern California marine bays and 132 species in polyhaline San Francisco Bay. Once species tolerance scores were calculated, the BQI value for each sample k was computed as

$$BQI_k = \left(\sum_i^n \left(\frac{A_i}{totA} ES50_{0.05i} \right) \right) (\log_{10}(S + 1)),$$

where n is the number of species in the sample with tolerance scores, A_i is the abundance of species i , $totA$ is the total abundance in the sample, and S is the number of species in the sample. Higher BQI values are associated with lower pollution levels.

2.2.3. Relative Benthic Index (RBI)

We calculated Relative Benthic Index (RBI) values following the method of Hunt et al. (2001). The RBI was first calibrated to each habitat by selecting negative and positive indicator taxa. Then, RBI values were calculated as the weighted sum of (a) four community parameters (total number of species, number of crustacean species, number of crustacean individuals, and number of mollusc species), and abundances of (b) three positive and (c) two negative indicator organisms. The negative indicator taxa selected for both habitats were oligochaeta and *Capitella capitata* complex, which have been used for this purpose in previous versions of the RBI. For positive indicator taxa, we followed the practice of selecting an amphipod, a bivalve, and a polychaete, which is typical of previous applications of the RBI. For southern California marine bays, we selected the amphipod *Monocorophium insidiosum*, the bivalve

Asthenothaerus diegensis, and the polychaete *Goniada littorea*. For polyhaline San Francisco Bay positive indicator taxa, we selected the amphipod *Sinocorophium heteroceratum*, the bivalve *Rochefortia* spp., and the polychaete *Prionospio lighti*. The RBI was scaled from 0 to 1.0 in each habitat by subtracting the lowest value and dividing by the range; thus 0 was the “worst” sample in the calibration data and 1 the “best.”

2.2.4. River Invertebrate Prediction and Classification System (RIVPACS)

The River Invertebrate Prediction and Classification System (RIVPACS) approach was calibrated following the methods of Wright et al. (1993) and Van Sickle et al. (2006). Cluster analysis was used to define site-groups of reference samples in the calibration data, based on the presence or absence of species occurring there. Discriminant function analysis of habitat variables at the site-groups was then used to build discriminant functions that can be used to classify future sampling sites into site-groups based on habitat variable values. Minimally impacted reference sites for this calibration were selected by eliminating samples with high toxicity (control-adjusted survival < 50%) to amphipods, one or more chemicals exceeding ERM concentrations (Long et al., 1995), three or more chemicals exceeding their ERL concentrations (Long et al., 1995) or from sites influenced by point-source discharges.

Several different habitat models explaining site groupings based on species abundances were explored in the southern California marine bays and polyhaline San Francisco Bay by altering the numbers of site groupings and by varying the habitat variables used to explain the groupings. Based on the proportion of variance explained, 12 and 4 site-group models based on latitude, longitude, and depth were selected for the southern California marine bays and San Francisco Bay, respectively. The probability of belonging to each of the site groups was calculated for each test site, based on the habitat variables. The site-group mean abundance for each taxon was then combined with the group probabilities to generate an expected taxon list specific to each test site. All permutations and combinations of numbers of groups and habitat variables were tested, and the combination with the greatest RIVPACS score improvement over an equivalent, non-predictive null model was selected (Van Sickle et al., 2005). Predictive improvement was quantified by calculating the reduction in root mean squared error (RMSE) of the predictive model (i.e., the model built using a discriminant function) from the null model. The chosen discriminant function model was then used to establish predictions for the species that would be expected to occur at reference sites in each group. The discriminant functions developed during calibration were used on the evaluation samples, first to identify the habitat site-group to which a sample belonged, and then to evaluate the observed species in relation to expectations for a minimally disturbed reference site. The difference between expected and observed assemblages measures the departure of the site from reference condition. For southern California marine bays, 457 species with >50% probability of occurring in reference samples were included in the predictive model, while 119 species were included for polyhaline San Francisco Bay. Summary statistics for the models are presented in Table 4. Based on a one to one ratio of modeled expected to observed (*O/E*) species present at validation sites they explained 89% and 96% of the variance, respectively.

2.2.5. Index of Biotic Integrity (IBI)

The Index of Biotic Integrity (IBI) approach developed by Thompson and Lowe (2004) was applied in San Francisco Bay without modification. The same approach was applied to the calibration data for the southern California marine bays. First, 22 candidate

Table 4

Summary statistics for RIVPACS predictive models (see Van Sickle et al., 2005). *O/E*: Observed to expected species ratio.

Statistic	Southern California Marine Bays	Polyhaline Central San Francisco Bay
<i>O/E</i> root mean squared error for predictive model based on validation sites	0.270	*
<i>O/E</i> standard deviation for null model (highest variability model)	0.434	0.451
<i>O/E</i> standard deviation for predictive model based on calibration sites	0.301	0.261
Predictive improvement over the null model	0.133	0.190
Standard deviation for calibration pseudoreplicate samples (least variability possible)	0.173	0.259

* Calibration data used for model development validation.

metrics were evaluated for suitability as indicators, based on criteria such as conforming to current conceptual models of benthic response to contamination and demonstrating measurable response to sediment contamination. Plots of candidate indicators vs. mERMq were examined, multiple regression analysis was conducted to evaluate the relationships between candidate IBI metrics and percent fines, TOC, and mERMq (independent variables), and four metrics were selected. Next, 59 reference samples were identified and reference ranges calculated for the four selected metrics as the maximum and minimum values for the reference samples. Reference sample selection was based on the same four criteria as Ranasinghe et al. (2004), including the absence of toxicity to amphipods. Table 5 presents the benthic assessment measures and reference ranges that were selected for each habitat. The assessment measures selected for southern California marine bays were based on the present study and reference ranges were established using the 59 designated reference samples. The measures and ranges for polyhaline San Francisco Bay are those of Thompson and Lowe (2004).

2.3. Index threshold scaling

All five index approaches were calibrated to the same four-category scale of benthic condition: (1) Unaffected – a community that would occur at a reference site for that habitat; (2) Marginal deviation from reference – a community that exhibits some indication of stress, but might be within measurement variability of reference condition; (3) Affected – a community that exhibits clear evidence of physical, chemical, natural, or anthropogenic stress; (4) Severely Affected – a community exhibiting a high magnitude of stress. Affected and severely affected communities are those believed to be showing clear evidence of disturbance, while unaffected and marginal communities do not. Disturbed communities could be due to the effects of one or more types of anthropogenic or natural stress while undisturbed communities likely indicate minimal stress of all types.

Three approaches were used to establish threshold values for each index and the threshold set that performed best with the evaluation samples was selected. The first, referred to as developer thresholds, was established by applying the principles used in the original index approach to the calibration data. Two other sets of thresholds were established by applying statistical optimization methods to compare index values and benthic condition categories.

For the BRI, the developer thresholds were based on reductions in the numbers of species along the disturbance gradient. Thresh-

Table 5

IBI assessment measures and reference ranges for each habitat. The index value for a sample is the number of assessment measures with values outside the reference range (see Thompson and Lowe, 2004).

Southern California Marine Bays				Polyhaline Central San Francisco Bay			
Assessment measure	Reference Range			Assessment measure	Reference Range		
	Min.	Max.	Mean		Min.	Max.	Mean
Number of taxa (per 0.1 m ² sample)	13	99	48.5	Number of taxa (per 0.05 m ² sample)	21	66	40.4
Molluscan taxa (per 0.1 m ² sample)	2	25	10.6	Amphipod taxa (per 0.05 m ² sample)	2	11	5.3
<i>Notomastus</i> sp. abundance (per 0.1 m ²)	0	59	2.7	Total abundance (per 0.05 m ² sample)	97	2931	905.7
Sensitive taxa (%)	9.0	47.1	26.9	<i>Capitella capitata</i> abundance (per 0.05 m ²)	0	13	2.0

olds were established at index values along the disturbance gradient where the number of species declined to 95%, 75% and 25% of the reference species pool. These thresholds are equivalent to those established for the southern California mainland shelf by Smith et al. (2001) because similar reductions in numbers of species accompanied the changes in community structure and function on which those thresholds were based (see Smith et al., 2003).

The BQI developer thresholds were three equally spaced thresholds along the index range, following the approach of Rosenberg et al. (2004). RBI developer thresholds were based on the distribution of index values, following Hunt et al. (2001). Reference thresholds were selected to segregate clusters of stations with high RBI values, high values for community parameters, and the presence of at least two of the three positive indicator taxa. The threshold differentiating between disturbed and undisturbed areas (i.e., between Marginal and Affected) was designated as the minimum RBI value where all three positive indicator taxa were found; 0.26 was selected in polyhaline San Francisco Bay because *P. lighti* first occurred at this RBI value. The Reference–Marginal threshold was selected at a mode of first occurrence for 18–20 species in the southern California marine bay calibration data; when a number of species have their first station of occurrence around a certain RBI value, it probably indicates a combination of factors that represent a significant change in habitat quality. Because there was no obvious mode in first stations of occurrence for San Francisco Bay, the threshold between Moderate and Severely Affected was chosen at 0.10, the RBI value of the first station of occurrence of the positive indicator species *S. heteroceratum*.

For the RIVPACS approach, developer thresholds were set at 0.5, 1.0 and 2.0 standard deviations of the calibration score mean on either side of an observed to expected (*O/E*) ratio of 1.0. Benthic condition is considered to deteriorate when the *O/E* ratio deviates from 1.0.

For the IBI, the threshold development process of Thompson and Lowe (2004) was used. Sample IBI values were evaluated graphically and statistical comparisons of IBI values and sediment contamination (mERMq) in disturbed and undisturbed samples were used to evaluate whether the assessment results reflected significant differences in sediment contamination. In southern California, sites with no IBI measures outside a reference range were considered Reference, sites with only one measure outside a reference range were considered Marginal, sites with two measures outside the ranges were considered Affected, and sites with three or four measures outside their ranges were considered Severely Affected. In San Francisco Bay, sites with no measures or only one measure outside a reference range were considered Reference, sites with two measures outside their reference ranges were considered Marginal, sites with three measures outside their reference ranges were considered Affected, and sites with four measures outside their reference ranges were considered Severely Affected (Thompson and Lowe, 2004).

Two non-developer sets of thresholds were selected for each indicator, based on consensus condition categories assigned by four benthic experts to the 68 site subset of the calibration data.

One optimization technique was based on maximizing the weighted kappa statistic (Cohen, 1960, 1968), which measures agreement between indicator and consensus categories beyond that expected by chance. Weights were based on the linear weighting scheme of Cicchetti and Allison (1971), which give “partial credit” according to severity of disagreement. The second set of thresholds was based on maximizing agreement between indicator and consensus categories, with no weighting factors. To find the optimal set of thresholds in each case, weighted kappa statistics and percent agreement were computed for all possible sets of triplicate thresholds.

2.4. Evaluation of index performance

Index performance was assessed by comparing index results to the consensus assessment of nine benthic experts that were given species abundances, together with habitat, depth, salinity and sediment grain-size information for 35 sites (Weisberg et al., 2008) that were not used in index development or calibration. Site identity data beyond habitat (southern California or San Francisco Bay) were withheld from the experts. The experts were asked to (1) rank the sites in each habitat from best to worst condition and (2) classify each site on the four-category scale of benthic condition to which the benthic indices were calibrated.

Index condition rank order was evaluated against the average expert rank order using Spearman rank correlation coefficients based on index values for all 35 evaluation samples, except in the case of the IBI. For the IBI, in San Francisco Bay, index values were available for only 5 (of 11) evaluation samples that met Thompson and Lowe, 2004) assemblage criteria for IBI calculation. Associations among the five indices were also evaluated, and compared to associations among the experts, using Spearman rank correlation coefficients. All the index values used for the index condition rank order evaluation were used in this analysis.

Condition category assessments by the benthic indices, and by all possible index combinations, were compared to the consensus expert condition assessment in three ways:

1. Status classification accuracy, the accuracy with which an index differentiated benthos identified by the nine experts as disturbed (Affected or Severely Affected categories) from benthos identified as undisturbed (Reference or Marginal categories). This mimics the evaluation approach used in most previously published benthic indicator development efforts.
2. Categorical classification accuracy with respect to the four categories established for index calibration (Reference, Marginal, Affected or Severely Affected). This is more challenging than status classification because it requires finer discrimination of the same benthic responses among a larger number of categories.
3. Bias in category designation; the sum of differences between index (or index combination) category and the consensus categorical classification of the experts when categories are expressed numerically (Reference = 1, Severely Affected = 4).

Positive bias indicates a tendency to score samples as more disturbed than the expert consensus, while negative bias indicates a tendency to score samples as less disturbed. Larger absolute values indicate stronger bias.

Index combinations were evaluated as the median of the numeric categories (Reference = 1, Severely Affected = 4). If the median for the indices in a combination fell between categories, it was rounded to the higher effect category. Comparisons to the experts were performed for each of the three threshold approaches associated with each index, with the best-performing thresholds used when combining indices. Developer thresholds were selected for the BRI and IBI, kappa-optimized thresholds for the BQI and RBI, and category-optimized thresholds for the RIVPACS index. Status and category classification accuracy and category bias were calculated for 32 of the 35 evaluation samples. Two southern California samples and a San Francisco Bay sample were excluded because the experts were almost evenly split as to the condition of the sites.

3. Results

Spearman correlation coefficients between index condition ranks and the average expert ranks for the 35 evaluation samples ranged from 0.70 to 0.89 (Table 6). The strongest correlation coefficient for an index (0.89) was slightly stronger than the weakest correlation coefficient for an expert in polyhaline San Francisco Bay (0.88) and slightly weaker than the weakest expert (0.90) in the southern California marine bays. All the Spearman correlations were highly significant ($p < 0.01$), except for the IBI in San Francisco Bay.

Spearman correlation coefficients among index values were all substantially lower than correlations among the experts in southern California bays (Table 7). The highest inter-index value was

Table 6

Spearman rank correlation coefficients between index condition ranks and average expert condition rankings for evaluation samples. The average, maximum and minimum correlations for the benthic experts are presented below to provide context for the index correlations.

Index	Spearman rank correlation coefficient	
	Southern California Marine Bays ($n = 24$; $p < 0.0001$)	Polyhaline Central San Francisco Bay ($n = 11$; $p < 0.01$ except ‡; $n = 5$; NS)
BQI	0.89	0.89
BRI	0.88	0.77
IBI	0.70	0.71 [‡]
RBI	0.82	0.87
RIVPACS	0.84	0.82
Expert minimum	0.90	0.88
Expert mean	0.95	0.95
Expert maximum	0.98	0.99

Table 7

Spearman rank correlation coefficients among indices for the evaluation samples. For the same samples, mean, maximum and minimum values among the experts are presented at the bottom of Table 6.

	BRI	IBI	RBI	RIVPACS
<i>Southern California Bays ($n = 24$, $p \leq 0.001$)</i>				
BQI	0.78	0.63	0.73	0.78
BRI		0.64	0.63	0.77
IBI			0.70	0.78
RBI				0.75
<i>Polyhaline San Francisco Bay ($n = 11$, $p \leq 0.01$; except for IBI where $n = 5$, $p > 0.05$)</i>				
BQI	0.85	0.71	0.90	0.94
BRI		0.71	0.78	0.77
IBI			0.71	0.71
RBI				0.91

0.78 and the lowest inter-expert correlation was 0.90. Although inter-index correlations in San Francisco Bay were also generally lower than inter-expert correlations, the distributions overlapped. The three highest inter-index correlations (0.94, 0.91 and 0.90) were greater than the lowest inter-expert correlation (0.88) but smaller than the mean (0.95). There was no pattern in correlations among community measure and species composition indices.

Index condition categories were evaluated for the 32 category evaluation samples. In the southern California marine bays, the BRI and RIVPACS indices performed best, with 95.5% and 90.9% correct status classification, 63.6% and 68.2% correct category classification, and low bias (Table 8). Their status classification accuracy was higher than three of the nine experts and the BRI tied with two others. Status classification accuracy for the BRI was slightly higher than the average expert (94.4%), but not as high for RIVPACS. The RIVPACS category classification accuracy was higher than, and the BRI equal to, the lowest expert. None of the other indices had a status classification accuracy higher than the lowest expert but, except for the RBI, all were higher than 75%, which has been used as a standard for indices developed in other estuarine systems (e.g., Engle and Summers, 1999; Van Dolah et al., 1999). In polyhaline San Francisco Bay, at 100%, status classification accuracy for all five indices was the same as the three highest experts. All five indices had higher category classification accuracy than the weakest expert, but only the BQI was higher than the average expert.

When there were differences, indices based on species composition almost always had higher classification accuracy both for status and for four-category assessments than indices based only on community measures. In southern California marine bays, the BRI, RIVPACS and BQI, which are based on species composition, had status classification accuracy of 95.5%, 90.9% and 81.8%. The IBI and RBI, which are based on community measures, had status classification accuracy of 77.3% and 72.7%, respectively (Table 8). Four-category classification accuracy was 68.2%, 63.6% and 63.6% for the species composition based RIVPACS, BRI, and BQI, and 54.5% for the community measure based RBI and IBI. Category bias was also lower for RIVPACS and the BRI than for either of the community measure based indices. In polyhaline San Francisco Bay, category classification accuracy for the species composition based RIVPACS and BQI was 80.0% and 90.0%, and 70.0% and 75.0% for the community measure based RBI and IBI, respectively. The category classification accuracy for the BRI here was 70.0%, which was the only instance where accuracy for a species composition based index was lower than any community measure based index.

Index combinations generally performed better than individual indices, and combinations of three or more indices generally performed better than combinations of two. In southern California marine bays, 10 combinations of three or more indices achieved the highest status classification accuracy of 95.5% (Table 8). One of these combinations, #29, had the highest four-category classification accuracy of 81.8%. The accuracy for this four-index combination of the BRI, BQI, IBI and RIVPACS was the same as the accuracy of 81.8% for the average expert. Another six of these combinations were in second place for category classification accuracy at 77.3%. In polyhaline San Francisco Bay, the percentage of index combinations with category classification accuracy of 80% or higher increased from 40% for single indices to 50%, 80%, 100% and 100% for combinations of two, three, four and five indices.

When results for both habitats were combined, the three-index combinations that performed best were #24, a three-index combination of the BRI, RBI, and RIVPACS, #26, a four-index combination of the BRI, the RBI, the IBI and RIVPACS, and #29, a four-index combination of the BRI, the BQI, the IBI and RIVPACS. These combinations had the highest status classification accuracy (96.9%), the highest category classification accuracy (81.3%) and low bias (2, 4, and 4, respectively). These combinations outperformed the aver-

Table 8

Classification Accuracy and Bias for Indices and Index Combinations. Classification accuracy is presented for “undisturbed” vs. “disturbed” status and four condition categories. Each of 32 category evaluation samples was assessed into one of four numeric categories by the index or index combination and compared with consensus categories from an independent assessment by nine benthic experts. Bias is the sum of differences between index combination and consensus categories; positive values indicate a tendency to score samples as more disturbed than the expert consensus, while negative values indicate a tendency to score samples as less disturbed. The categories were 1: Reference; 2: Marginal; 3: Affected; 4: Severely Affected. Categories 1 and 2 were considered “undisturbed” and 3 and 4 as “disturbed.” Index results were combined as the median of the numeric categories; if the median fell between categories, it was rounded to the higher effect category. Results for the benthic experts are presented to provide context.

No. of indices	#	Measure	Southern California Marine Bays (n = 22)			Polyhaline Central San Francisco Bay (n = 10)		
			Category accuracy (%)	Category bias	Status accuracy (%)	Category accuracy (%)	Category bias	Status accuracy (%)
One	1	BQI	63.6	7	81.8	90.0	-1	100.0
	2	BRI	63.6	-2	95.5	70.0	-1	100.0
	3	IBI	54.5	-8	77.3	75.0	-1	100.0
	4	RBI	54.5	7	72.7	70.0	3	100.0
	5	RIV	68.2	2	90.9	80.0	0	100.0
Two	6	BQI, BRI	59.1	7	86.4	90.0	1	100.0
	7	BQI, IBI	59.1	5	81.8	90.0	-1	100.0
	8	BQI, RBI	50.0	10	77.3	70.0	3	100.0
	9	BQI, RIV	63.6	10	77.3	80.0	0	100.0
	10	BRI, IBI	72.7	0	90.9	70.0	-1	100.0
	11	BRI, RBI	59.1	7	86.4	70.0	3	100.0
	12	BRI, RIV	68.2	6	90.9	90.0	1	100.0
	13	IBI, RBI	45.5	6	72.7	70.0	3	100.0
	14	IBI, RIV	68.2	2	90.9	80.0	0	100.0
	15	RBI, RIV	50.0	10	77.3	70.0	3	100.0
Three	16	BRI IBI RBI	77.3	-1	95.5	80.0	2	100.0
	17	BQI BRI IBI	72.7	0	95.5	80.0	0	100.0
	18	BQI BRI RBI	72.7	4	86.4	90.0	1	100.0
	19	BQI BRI RIV	72.7	2	95.5	80.0	0	100.0
	20	BQI IBI RBI	68.2	5	86.4	70.0	1	100.0
	21	BQI IBI RIV	77.3	1	95.5	80.0	0	100.0
	22	BQI RBI RIV	68.2	5	86.4	80.0	0	100.0
	23	BRI IBI RIV	68.2	-3	95.5	80.0	0	100.0
	24	BRI RBI RIV	77.3	1	95.5	90.0	1	100.0
	25	IBI RBI RIV	77.3	1	95.5	70.0	1	100.0
Four	26	BRI IBI RBI RIV	77.3	3	95.5	90.0	1	100.0
	27	BQI IBI RBI RIV	68.2	5	86.4	80.0	0	100.0
	28	BQI BRI RBI RIV	72.7	6	86.4	90.0	1	100.0
	29	BQI BRI IBI RIV	81.8	4	95.5	80.0	0	100.0
	30	BQI BRI IBI RBI	72.7	6	86.4	90.0	1	100.0
Five	31	All	77.3	3	95.5	80.0	0	100.0
Expert Consensus		Minimum	63.6	0	86.4	60.0	0	90.0
		Average	81.8	-0.2	94.4	83.3	0.56	94.4
		Maximum	95.5	+4, -3	100.0	100.0	+4, -2	100.0

age expert for status classification, but were outperformed by four of the nine experts for categorical classification. All three of the best-performing combinations include a mixture of community measures and species composition indices.

4. Discussion

Indices that include measures of species composition generally outperformed indices that include only community measures. This is consistent with Weisberg et al. (1997), who found that relative dominance of pollution-tolerant and pollution-sensitive species were the metrics in their index that had the best relationship to pollution gradients. Pearson and Rosenberg (1978) suggest that the initial benthic response to low levels of stress is a shift in species composition, with shifts in community metrics, such as loss of species richness and biomass, manifesting at later stages of stress. Thus, indices based on community metrics should be more effective at differentiating sites subject to high levels of stress, but less effective at differentiating sites with low to intermediate levels of stress that are more typical of the estuarine sites encountered in California.

Combinations of indices consistently outperformed individual indices. Each of the indices relies on a subset of metrics used by experts. Generally, these metrics correlate, but there are circum-

stances when they can differ considerably, such as when the presence of a large filter feeder reduces species richness and abundance, or when only a few individuals of a few sensitive species occur. Use of multiple indices incorporates a larger number of metrics and presumably balances the occasional erratic behavior of individual metrics. In addition, some of the indices showed biases, with the RBI assessing samples as more disturbed than the experts and the IBI behaving the opposite. Use of multiple indices apparently balances out those biases.

Conclusions about relative performance of indices are reliant upon proper implementation of the index approaches. Our study team included the original developers of the index approach, or investigators who had previously published applications of these indices in other habitats, for four of the five indices evaluated. The team had less experience with the BQI, but this method involves the least amount of developer judgment in its calibration. One indication that our results reflect successful implementation was the high classification accuracy for discriminating among undisturbed and disturbed benthic community status for all of our indices. Our range of 72.7–100% classification accuracy achieved for the individual indices compares favorably with the average status classification accuracy of 85% that Weisberg et al. (1997) achieved for seven Chesapeake Bay habitats, the 85% that Van Dolah et al. (1999) achieved in the best of his four southeast-

ern USA estuaries, and the 76% that Engle and Summers (1999) achieved for Gulf of Mexico estuaries.

One factor that may have led to our slightly higher validation success was our approach to validation. Validation has historically been conducted by using chemical and toxicological exposure measures to identify sites of supposedly extreme condition (Borja and Dauer, 2008; Weisberg et al., 2008). Here, we used the professional judgment of benthic ecologists that reviewed benthic macrofaunal data alone to establish a validation site's condition (Weisberg et al., 2008), minimizing the likelihood of incorrect classifications due to reliance on predictions from exposure data. Use of expert judgment reduces false undisturbed classifications of sites affected by unmeasured chemicals or physical disturbance. It also avoids false disturbed site designations due to contaminants that are measured in chemical analysis but are tightly bound to sediments and unavailable *in situ* to benthic organisms (Batley et al., 2005).

Using expert judgment to classify sites for index validation has the additional advantage of allowing evaluation of index performance at sites experiencing intermediate levels of disturbance. This cannot be conducted using exposure measures to classify validation sites, as there is no expectation of a linear relationship between biological responses and chemical exposure. Assessment of intermediate conditions is a more difficult, but more relevant, assessment challenge for benthic indices. Interestingly, the indices matched expert opinion for the intermediate sites as well as they did for sites of more extreme condition.

The level of agreement among experts provides a benchmark for evaluating index performance. Historically, index developers have deemed an index successful if it correctly identifies 75–80% of sites with extreme exposure conditions (Van Dolah et al., 1999). However, since indices are intended to reproduce the experience of experts in interpreting benthic data using an objective, repeatable, transparent tool, a better evaluation benchmark is whether an index ranks and classifies sites with levels of correlation and accuracy comparable to that among experts. In this study, none of the individual indices achieved this mark, but several index combinations did.

References

- Batley, G.E., Stahl, R.G., Babut, M.P., Bott, T.L., Clark, J.R., Field, L.J., Ho, K.T., Mount, D.R., Swartz, R.C., Tessier, A., 2005. Scientific underpinnings of sediment quality guidelines. In: Wenning, R.J., Batley, G.E., Ingersoll, C.G., Moore, D.W. (Eds.), *Use of Sediment Quality Guidelines (SQGs) and Related Tools for the Assessment of Contaminated Sediments*. Society of Environmental Toxicology and Chemistry, Pensacola, pp. 39–119.
- Bay Area Dischargers Association, 1994. Local effects monitoring program, quality assurance project plan. Bay Area Clean Water Association, Oakland, CA.
- Bergen, M., Cadien, D.B., Dalkey, A., Montagne, D.E., Smith, R.W., Stull, J.K., Velarde, R.G., Weisberg, S.B., 2000. Assessment of benthic infaunal condition on the mainland shelf of Southern California. *Environmental Monitoring and Assessment* 64, 421–434.
- Bilkovic, D.M., Roggero, M., Hershner, C.H., Havens, K.H., 2006. Influence of land use on macrobenthic communities in nearshore estuarine habitats. *Estuaries and Coasts* 29, 1185–1195.
- Blanchet, H., Lavesque, N., Ruellet, T., Dauvin, J.C., Sauriau, P.G., Desroy, N., Desclaux, C., Leconte, M., Bachelet, G., Janson, A.L., Bessineton, C., Duhamel, S., Jourde, J., Mayot, S., Simon, S., de Montaudouin, X., 2008. Use of biotic indices in semi-enclosed coastal ecosystems and transitional waters habitats – implications for the implementation of the European Water Framework Directive. *Ecological Indicators* 8, 360–372.
- Borja, A., Dauer, D.M., 2008. Assessing the environmental quality status in estuarine and coastal systems: comparing methodologies and indices. *Ecological Indicators* 8, 331–337.
- Borja, A., Franco, J., Perez, V., 2000. A marine Biotic Index to establish the ecological quality of soft-bottom benthos within European estuarine and coastal environments. *Marine Pollution Bulletin* 40, 1100–1114.
- Borja, A., Muxika, I., Franco, J., 2003. The application of a Marine Biotic Index to different impact sources affecting soft-bottom benthic communities along European coasts. *Marine Pollution Bulletin* 46, 835–845.
- Borja, A., Josefson, A.B., Miles, A., Muxika, I., Olsgard, F., Phillips, G., Rodriguez, J.G., 2007. An approach to the intercalibration of benthic ecological status assessment in the North Atlantic ecoregion, according to the European Water Framework Directive. *Marine Pollution Bulletin* 55, 42–52.
- Borja, A., Dauer, D.M., Diaz, R., Llanso, R.J., Muxika, I., Rodriguez, J.G., Schaffner, L., 2008. Assessing estuarine benthic quality conditions in Chesapeake Bay: a comparison of three indices. *Ecological Indicators* 8, 395–403.
- Brown, J., Bay, S.M., 2005. Temporal assessment of chemistry, toxicity and benthic communities in sediments at the mouths of Chollas Creek and Paleta Creek, San Diego Bay. Southern California Coastal Water Research Project, Draft Report, Westminster, CA.
- Cicchetti, D.V., Allison, T., 1971. A new procedure for assessing reliability of scoring EEG sleep recordings. *American Journal of EEG Technology* 11, 101–109.
- Cohen, J., 1960. A coefficient of agreement for nominal scales. *Educational and Psychological Measurement* 20, 37–46.
- Cohen, J., 1968. Weighted Kappa nominal scale agreement with provision for scale disagreement or partial credit. *Psychological Bulletin* 70, 213–220.
- Dauer, D.M., Ranasinghe, J.A., Weisberg, S.B., 2000. Relationships between benthic community condition, water quality, sediment quality, nutrient loads, and land use patterns in Chesapeake Bay. *Estuaries* 23, 80–96.
- Diaz, R.J., Solan, M., Valente, R.M., 2004. A review of approaches for classifying benthic habitats and evaluating habitat quality. *Journal of Environmental Management* 73, 165–181.
- Engle, V.D., Summers, J.K., 1999. Refinement, validation, and application of a benthic condition index for Northern Gulf of Mexico estuaries. *Estuaries* 22, 624–635.
- Hale, S.S., Paul, J.F., Heltshe, J.F., 2004. Watershed landscape indicators of estuarine benthic condition. *Estuaries* 27, 283–295.
- Hawkins, C.P., Norris, R.H., Hogue, J.N., Feminella, J.W., 2000. Development and evaluation of predictive models for measuring the biological integrity of streams. *Ecological Applications* 10, 1456–1477.
- Hunt, J.W., Anderson, B.S., Phillips, B.M., Tjeerdema, R.S., Taberski, K.M., Wilson, C.J., Puckett, H.M., Stephenson, M., Fairey, R., Oakden, J.M., 2001. A large-scale categorization of sites in San Francisco Bay, USA, based on the sediment quality triad, toxicity identification evaluations, and gradient studies. *Environmental Toxicology and Chemistry* 20, 1252–1265.
- Hyland, J.L., Van Dolah, R.F., Snoots, T.R., 1999. Predicting stress in benthic communities of southeastern U.S. estuaries in relation to chemical contamination of sediments. *Environmental Toxicology and Chemistry* 18, 2557–2564.
- Hyland, J.L., Balthis, W.L., Engle, V.D., Long, E.R., Paul, J.F., Summers, J.K., Van Dolah, R.F., 2003. Incidence of stress in benthic communities along the US Atlantic and Gulf of Mexico coasts within different ranges of sediment contamination from chemical mixtures. *Environmental Monitoring and Assessment* 81, 149–161.
- Labrone, C., Amouroux, J.M., Sarda, R., Dutrieux, E., Thorin, S., Rosenberg, R., Gremare, A., 2006. Characterization of the ecological quality of the coastal Gulf of Lions (NW Mediterranean). A comparative approach based on three biotic indices. *Marine Pollution Bulletin* 52, 34–47.
- Leung, K.M.Y., Bjorgesæter, A., Gray, J.S., Li, W.K., Lui, G.C.S., Wang, Y., Lam, P.K.S., 2005. Deriving sediment quality guidelines from field-based species sensitivity distributions. *Environmental Science and Technology* 39, 5148–5156.
- Long, E.R., MacDonald, D.D., Smith, S.L., Calder, F.D., 1995. Incidence of adverse biological effects within ranges of chemical concentrations in marine and estuarine sediments. *Environmental Management* 19, 81–97.
- Long, E.R., MacDonald, D.D., Severn, C.G., Hong, C.B., 2000. Classifying probabilities of acute toxicity in marine sediments with empirically derived sediment quality guidelines. *Environmental Toxicology and Chemistry* 19, 2598–2601.
- Long, E.R., Ingersoll, C.G., MacDonald, D.D., 2006. Calculation and uses of mean sediment quality guideline quotients: a critical review. *Environmental Science and Technology* 40, 1726–1736.
- Pearson, T.H., Rosenberg, R., 1978. Macrobenthic succession in relation to organic enrichment and pollution of the marine environment. *Oceanography and Marine Biology: An Annual Review* 16, 229–311.
- Puente, A., Juanes, J.A., Garcia, A., Alvarez, C., Revilla, J.A., Carranza, I., 2008. Ecological assessment of soft bottom benthic communities in northern Spanish estuaries. *Ecological Indicators* 8, 373–388.
- Quintino, V., Elliott, M., Rodrigues, A.M., 2006. The derivation, performance and role of univariate and multivariate indicators of benthic change: case studies at differing spatial scales. *Journal of Experimental Marine Biology and Ecology* 330, 368–382.
- Ranasinghe, J.A., Frithsen, J.B., Kutz, F.W., Paul, J.F., Russell, D.E., Batiuk, R.A., Hyland, J.L., Scott, K.J., Dauer, D.M., 2002. Application of two indices of benthic community condition in Chesapeake Bay. *Environmetrics* 13, 499–511.
- Ranasinghe, J.A., Montagne, D.E., Smith, R.W., Mikel, T.K., Weisberg, S.B., Cadien, D.B., Velarde, R.G., Dalkey, A., 2003. Southern California Bight 1998 Regional Monitoring Program: VII. Benthic Macrofauna. Southern California Coastal Water Research Project, Westminster, CA.
- Ranasinghe, J.A., Thompson, B., Smith, R.W., Lowe, S., Schiff, K.C., 2004. Evaluation of benthic assessment methodology in southern California bays and San Francisco Bay. Southern California Coastal Water Research Project, Technical Report 432, Westminster, CA.
- Ranasinghe, J.A., Barnett, A.M., Schiff, K.C., Montagne, D.E., Brantley, C., Beegan, C., Cadien, D.B., Cash, C., Deets, G.B., Diener, D.R., Mikel, T.K., Smith, R.W., Velarde, R.G., Watts, S.D., Weisberg, S.B., 2007. Southern California Bight 2003 Regional Monitoring Program: III Benthic Macrofauna. Southern California Coastal Water Research Project Authority, Costa Mesa, CA.
- Rehn, A.C., Ode, P.R., Hawkins, C.P., 2007. Comparisons of targeted-riffle and reach-wide benthic macroinvertebrate samples: implications for data sharing in

- stream-condition assessments. *Journal of the North American Benthological Society* 26, 332–348.
- Rosenberg, R., Blomqvist, M., Nilsson, H.C., Cederwall, H., Dimming, A., 2004. Marine quality assessment by use of benthic species-abundance distributions: a proposed new protocol within the European Union Water Framework Directive. *Marine Pollution Bulletin* 49, 728–739.
- Southern California Coastal Water Research Project, Space and Naval Warfare Systems Center San Diego (SCCWRP and SPAWAR), 2004. Sediment assessment study for the mouths of Chollas and Paleta Creek, San Diego. Phase I Draft Report. San Diego Regional Water Quality Control Board, Commander Navy Region Southwest, City of San Diego, San Diego, CA.
- Smith, R.W., Bergen, M., Weisberg, S.B., Cadien, D.B., Dalkey, A., Montagne, D.E., Stull, J.K., Velarde, R.G., 2001. Benthic response index for assessing infaunal communities on the southern California mainland shelf. *Ecological Applications* 11, 1073–1087.
- Smith, R.W., Ranasinghe, J.A., Weisberg, S.B., Montagne, D.E., Cadien, D.B., Mikel, T.K., Velarde, R.G., Dalkey, A., 2003. Extending the southern California Benthic Response Index to assess benthic condition in bays Southern California Coastal Water Research Project, Technical Report 410, Westminster, CA.
- Summers, J.K., 2001. Ecological condition of the estuaries of the Atlantic and gulf coasts of the United States. *Environmental Toxicology and Chemistry* 20, 99–106.
- Teixeira, H., Salas, F., Borja, A., Neto, J.M., Marques, J.C., 2008. A benthic perspective in assessing the ecological status of estuaries: the case of the Mondego estuary (Portugal). *Ecological Indicators* 8, 404–416.
- Thompson, B., Lowe, S., 2004. Assessment of macrobenthos response to sediment contamination in the San Francisco Estuary, California, USA. *Environmental Toxicology and Chemistry* 23, 2178–2187.
- Thompson, B., Anderson, B.S., Hunt, J.W., Taberski, K.M., Phillips, B.M., 1999. Relationships between sediment contamination and toxicity in San Francisco Bay. *Marine Environmental Research* 48, 285–395.
- U.S. Environmental Protection Agency, 2004. National Coastal Condition Report II. U.S. Environmental Protection Agency, Office of Research and Development, EPA-620/R-03/002, Washington, DC.
- Van Dolah, R.F., Hyland, J.L., Holland, A.F., Rosen, J.S., Snoots, T.R., 1999. A benthic index of biological integrity for assessing habitat quality in estuaries of the southeastern USA. *Marine Environmental Research* 48, 269–283.
- Van Sickle, J., Hawkins, C.P., Larsen, D.P., Herlihy, A.T., 2005. A null model for the expected macroinvertebrate assemblage in streams. *Journal of the North American Benthological Society* 24, 178–191.
- Van Sickle, J., Huff, D.D., Hawkins, C.P., 2006. Selecting discriminant function models for predicting the expected richness of aquatic macroinvertebrates. *Freshwater Biology* 51, 359–372.
- Weisberg, S.B., Ranasinghe, J.A., Schaffner, L.C., Diaz, R.J., Dauer, D.M., Frithsen, J.B., 1997. An estuarine benthic index of biotic integrity (B-IBI) for Chesapeake Bay. *Estuaries* 20, 149–158.
- Weisberg, S.B., Thompson, B., Ranasinghe, J.A., Montagne, D.E., Cadien, D.B., Dauer, D.M., Diener, D.R., Oliver, J.S., Reish, D.J., Velarde, R.G., Word, J.Q., 2008. The level of agreement among experts applying best professional judgment to assess the condition of benthic infaunal communities. *Ecological Indicators* 8, 389–394.
- Wright, J.F., Furse, M.T., Armitage, P.D., 1993. RIVPACS: a technique for evaluating the biological water quality of rivers in the UK. *European Water Pollution Control* 3, 15–25.
- Zettler, M.L., Schiedek, D., Bobertz, B., 2007. Benthic biodiversity indices versus salinity gradient in the southern Baltic Sea. *Marine Pollution Bulletin* 55, 258–270.