



JRC TECHNICAL REPORTS

Coastal waters North East Atlantic geographic intercalibration group

Benthic invertebrate fauna ecological assessment methods

Van Hoey, G., Bonne, W., Muxica, I., Josefson, A., Borgersen, G., Rygg, B., Borja, A., Philips, G., Miles, A., Dubois, S., Desroy, N., Buchet, R., Ximenes, M.C., O'Beirn, F., Witt, J., Heyer, K., van Loon, W., Ruiter, H., Neto, J., Marques, J.C., Garcia, E., Puente, A., Salas Herrero, F

2019

This publication is a Technical report by the Joint Research Centre (JRC), the European Commission's science and knowledge service. It aims to provide evidence-based scientific support to the European policy-making process. The scientific output expressed does not imply a policy position of the European Commission. Neither the European Commission nor any person acting on behalf of the Commission is responsible for the use which might be made of this publication.

Contact information

Name: Fuensanta Salas Herrero
Address: Via Enrico Fermi 2749, 21027 Ispra (VA), Italy
Email: Fuensanta.Salas-Herrero@ec.europa.eu
Tel.: +39 0332 78 5701

JRC Science Hub

<https://ec.europa.eu/jrc>

JRC115475

EUR 29640 EN

PDF ISBN 978-92-79-99232-2 ISSN 1831-9424 doi:10.2760/16318

Luxembourg: Publications Office of the European Union, 2019

© European Union, 2019

The reuse policy of the European Commission is implemented by Commission Decision 2011/833/EU of 12 December 2011 on the reuse of Commission documents (OJ L 330, 14.12.2011, p. 39). Reuse is authorised, provided the source of the document is acknowledged and its original meaning or message is not distorted. The European Commission shall not be liable for any consequence stemming from the reuse. For any use or reproduction of photos or other material that is not owned by the EU, permission must be sought directly from the copyright holders.

All images © European Union 2019, except:

Cover page, Ria de Aveiro, Aveiro, Portugal. © Heliana Teixeira, Guimarães, Portugal

How to cite: Van Hoey, G., Bonne, W., Muxica, I., Josefson, A., Borgersen, G., Rygg, B., G., Borja, A., Philips, G., Miles, Al., Dubois, S., Desroy, N., Buchet, R., Ximenes, MC., O'Beirn, F., Witt, J., Heyer, K., van Loon, W., Ruitter, H., Neto, J., Marques, JC., Garcia, E., Puente, A., Salas Herrero., F. Coastal waters North East Atlantic Geographic Intercalibration Group. Benthic invertebrate fauna ecological assessment methods; EUR 29640 EN, Publications Office of the European Union, Luxembourg, 2019, ISBN 978-92-79-99232-2, doi:10.2760/16318, JRC115475

Contents

Foreword	1
Acknowledgment	2
Abstract	3
1. Introduction	4
PART A- Common type NEA 1/26.....	5
A2. Description of national assessment methods	5
A2.1 Methods and required BQE parameters	6
A2.2 Sampling and data processing	11
A2.3 National reference conditions.....	11
A2.4 National boundary setting	18
A2.5 Results of WFD compliance checking	25
A3. Intercalibration feasibility checking	27
A3.1 Typology	27
A3.2 Pressures addressed.....	27
A3.2.1 Sample level.....	27
A3.2.2 Higher level comparison.....	30
A.3.3 Assessment concept.....	31
A.4 Collection of intercalibration dataset and benchmarking	35
A.4.1 Dataset description	35
A.4.2 Data acceptance criteria.....	35
A.4.2.1 General multivariate analyses	37
A.4.2.2 Multivariate analyses of the benthic univariate parameters	39
A4.3 Common benchmark	41
A.4.3.1 Benchmark standardization	42
A5. Comparison of methods and boundaries	46
A.5.1 Intercalibration option and common metrics	46
A.5.2 Results of the regression comparison.....	47
A.5.2.1 Sample level comparison.....	47
A.5.2.2. Higher level comparison (+ BEQI, Belgium).....	48
A.5.3 Comparability criteria	50
A.5.3.1 Sample level comparison.....	50
A.5.3.2 Higher level comparison (BEQI, Belgium)	51

A.6 Final results to be included in the EC	53
A.7 Ecological characteristics	54
A.7.1 Description of reference or alternative benchmark communities	54
A.7.2 Description of good status communities	54
PART B- Common type NEA 3/4	58
B.2 Description of national assessment methods	58
B.2.1 Methods and required BQE parameters	58
B.2.2 Sampling and data processing	60
B.2.3 National reference conditions.....	60
B.2.4 National boundary setting	63
B.2.4 Results of WFD compliance checking	65
B.3 Feasibility checking	66
B.3.1 Typology	66
B.3.2 Pressures addressed.....	66
B.3.3 Assessment concept	68
B.4 Collection of intercalibration dataset and benchmarking	70
B.4.1 Dataset description	70
B.4.2 Data acceptance criteria.....	71
B.4.3 General multivariate analyses	71
B.4.4 Common benchmark	73
B.4.5 Benchmark standardization	73
B.5 Comparison of methods and boundaries	75
B.5.1 Intercalibration option and common metrics	75
B.5.2 Results of the regression comparison.....	75
B.5.3 Comparability criteria	75
B.6 Final results to be included in the EC	77
B.7 Ecological characteristics	77
B.7.1 Description of reference or alternative benchmark communities	77
B.7.2 Description of good status communities	78
PART C- Common type NEA 7	79
C.2 Description of national assessment methods	79
C.2.1 Methods and required BQE parameters	80
C.2.2 Sampling and data processing	81
C.2.3 National reference conditions.....	82

C.2.4 National boundary settings.....	83
C.2.5 Results of WFD compliance checking	85
C.3 Feasibility checking	86
C.3.1 Typology.....	86
C.3.2 Pressures addressed.....	86
C.3.3 Assessment concept.....	89
C.4. Collection of intercalibration dataset and benchmarking.....	90
C.4.1 Dataset description	90
C.4.2 Common benchmarking or reference conditions	91
C.5 Comparison of methods and boundaries	92
C.5.1 Intercalibration option and common metrics	92
C.5.2 Results of regression comparision	92
C.5.3 Comparability criteria	93
C.6 Final results to be included in the EC.....	93
C.7 Ecological characteristics	94
C.7.1 Description of reference or alternative benchmark communities.....	94
C.7.2 Description of good status communities	94
PART D-Type NEA 5	95
D.2 Description of national assessment methods	95
D.2.1 Methods and required BQE parameters.....	95
D.2.2 Sampling and data processing	96
D.2.3 National reference conditions	97
D.2.4 National boundary settings.....	97
D.2.5 Results of WFD compliance checking	98
D.3 Feasibility checking	99
D.3.1 Typology.....	99
D.3.2 Pressures addressed	99
D.4 Ecological characteristics	99
D.4.1 Description of reference or alternative benchmark communities.....	99
D.4.2 Description of good status communities.....	99
References	101
Annex 1. Common type NEA 1/26: Sampling and data processing information	104
Annex 2. Common type NEA 1/26: Alternative benchmark approach (based on biotic variables)	119

Annex 3. Common type NEA 3/4. Pressures 133

List of abbreviations and definitions 136

Key Terms: 136

 Abbreviations: 137

List of figures 138

List of tables 140

Foreword

The European Water Framework Directive (WFD) requires the national classifications of good ecological status to be harmonised through an intercalibration exercise. In this exercise, significant differences in status classification among Member States are harmonized by comparing and, if necessary, adjusting the good status boundaries of the national assessment methods.

Intercalibration is performed for rivers, lakes, coastal and transitional waters, focusing on selected types of water bodies (intercalibration types), anthropogenic pressures and Biological Quality Elements. Intercalibration exercises are carried out in Geographical Intercalibration Groups - larger geographical units including Member States with similar water body types - and followed the procedure described in the WFD Common Implementation Strategy Guidance document on the intercalibration process (European Commission, 2011).

The Technical report on the Water Framework Directive intercalibration describes in detail how the intercalibration exercise has been carried out for the water categories and biological quality elements. The Technical report is organized in volumes according to the water category (rivers, lakes, coastal and transitional waters), Biological Quality Element and Geographical Intercalibration group. This volume addresses the intercalibration of the Coastal Waters-North East Atlantic Benthic Invertebrates Fauna ecological assessment

Acknowledgments

Part of this work (Intercalibration for common types NEA 1/26 and NEA 3/4) for the WFD intercalibration obligation of the Member States has received funding within the framework of JPI Oceans through the JPI Oceans pilot action "Joint funding of the Scientific Intercalibration exercise for the EU Water Framework Directive (WFD) coastal and transitional waters in the North-East Atlantic". This work is the effort of many benthic experts of the Member States, which deliver the necessary information and feedback.

This exercise has been supported by JRC and a review panel formed by recognized experts on IC procedure and benthic invertebrates element.

Thanks to all contributors during the three phases of the intercalibration. Recording all contributors to this intercalibration work in the NEA-GIG is impossible, but at least the current responsible persons for each Member State are mentioned.

Abstract

The European Water Framework Directive (WFD) requires the national classifications of good ecological status to be harmonised through an intercalibration exercise. In this exercise, significant differences in status classification among Member States are harmonized by comparing and, if necessary, adjusting the good status boundaries of the national assessment methods.

Intercalibration is performed for rivers, lakes, coastal and transitional waters, focusing on selected types of water bodies (intercalibration types), anthropogenic pressures and Biological Quality Elements. Intercalibration exercises are carried out in Geographical Intercalibration Groups - larger geographical units including Member States with similar water body types - and followed the procedure described in the WFD Common Implementation Strategy Guidance document on the intercalibration process (European Commission, 2011).

The Technical report on the Water Framework Directive intercalibration describes in detail how the intercalibration exercise has been carried out for the water categories and biological quality elements. The Technical report is organized in volumes according to the water category (rivers, lakes, coastal and transitional waters), Biological Quality Element and Geographical Intercalibration group.

This report gives a description of the intercalibration of the different benthic assessment approaches for in coastal waters in the North East Atlantic Geographical Intercalibration Group (NEA-GIG) for types NEA 1/26 (Exposed or sheltered, euhaline, shallow waters), NEA 3/4 (Wadden sea type) and NEA 7 (Deep fjordic and sea loach systems). The benthic assessment approaches of nine European Member States (Belgium, Germany, Denmark, France, Ireland, the Netherlands, Portugal, Spain and the United Kingdom) and Norway are intercalibrated. In Spain, the competent authorities for the WFD application are the regions ('autonomous communities'); therefore for the benthic assessment methods three regions have been considered: Basque Country, Andalusia and Cantabria (no information on Galicia or Asturias). Part D of the report describes the Germany assessment approach for the type NEA 5. This type is not shared with the rest of the Members States, and therefore, the Intercalibration is not possible

1. Introduction

The intercalibration in the NEA-GIG region for coastal waters has a long history. In the first phase, a pioneering intercalibration exercise was executed, which showed a high consistency between the different benthic assessment approaches of United Kingdom, Spain (Basque Country), Denmark and Norway on a common benthos dataset (Borja et al., 2007). In the second phase, when the intercalibration guidelines were developed, a re-run of the analyses of the coastal waters of phase I following the new comparability criteria was expected. However, this process could not be completed in phase II for several reasons. The main recommendation from the Review Panel on the intercalibration exercise for the coastal waters in the NEA-GIG region was that additional analyses should be done (including all methods and all Member States) to further refine the comparability (Davies, 2012). Currently, further clarifications/justification should be compiled to confirm the comparability of the NEA-GIG benthic assessment approaches. Therefore, in phase III, under the form of a JPI oceans pilot action (<http://www.jpi-oceans.eu/intercalibration-eu-water-framework-directive>) and with the support of the Joint Research Centre (JRC), this process has been executed. The objectives of this action are:

- WFD method compliance documentation check, explanations of the justifications for assessment methods including specific parameters, reference conditions and the boundary setting procedure. Also to check or improve pressure-response relationships (1st and 2nd phase results are available).
- Provide an alternative benchmarking clarification, trying to take regional biological differences and sampling protocol differences into account, based on already available data or validated expert judgment.
- Check and improve comparability analysis (1st and 2nd phase results are available).
- Prepare and compile finalized intercalibration technical report from the several existing current reports.

This report compiles all the latest information regarding the benthic assessment approaches, boundary- and reference settings for each Member State and common dataset characteristics. Specific analyses were conducted to demonstrate the pressure-response relationships of the benthic assessment approaches, detect possible bio-geographical differences in the common dataset, perform an alternative benchmark delineation and the comparability analyses following the intercalibration guidelines (Guidance document 14: guidance document on the intercalibration process 2008-2011).

PART A- Common type NEA 1/26

A2. Description of national assessment methods

Within the NEA-GIG region for coastal waters, 10 benthic assessment approaches were defined (Table 1). A benthic assessment approach consists of an indicator algorithm, boundary settings and a reference setting approach. Some Member States used the same indicator algorithm (e.g. m-AMBI), but were considered as a separate benthic assessment approach due to different boundary or reference settings. Only United Kingdom/Ireland and the Basque Country (BC)/Cantabria (C) in Spain share the same benthic assessment approach, the IQI and m-AMBI respectively. Each benthic assessment approach is considered as a separate method in the intercalibration exercise. RAT method is applied on rocky substratum and is not comparable with the rest of assessment methods.

Table 1. Overview of the national assessment methods

Member State		Method		WISER database	Included in this IC exercise
Belgium	BE	Benthic Ecosystem Quality Index	BEQI	X	Yes
Germany	DE	Multivariate AZTI's Marine Biotic Index	m-AMBI ¹	X	Yes
Denmark	DK	Danish Quality Index	DKI	X	Yes
France	FR	Multivariate AZTI's Marine Biotic Index	m-AMBI ²	X	Yes
Ireland	IE	Infaunal Quality Index	IQI	X	Yes
Netherlands	NL	Benthic Ecosystem Quality Index 2	BEQI2	X	Yes
Norway	NO	Norwegian Quality Index	NQI	X	Yes
Portugal	PT	Benthic Assessment Tool	BAT	X	Yes
Portugal	PT	Rocky Shore Assessment Tool	RAT		No
Spain (Basque Country; Cantabria)	ES-BC/C	Multivariate AZTI's Marine Biotic Index	m-AMBI	X	Yes
Spain (Andalusia)	ES-A	Benthic Opportunistic polychaetes/amphipods index	BO2A	x	yes
United Kingdom	UK	Infaunal Quality Index	IQI	X	yes

¹m-AMBI method, but other reference and boundary settings.

²m-AMBI method, but other reference settings.

A2.1 Methods and required BQE parameters

The current intercalibration exercise is based on the latest versions of the indicator algorithms (Table 2). The EQR values determined for the samples within the common dataset are re-calculated based on those algorithms. The metric values (e.g. Shannon diversity, AMBI, S, etc.) were determined based on the latest version of the common benthic dataset, which was made available by Angel Borja (the NEA-GIG benthos lead in phase II). The metric AMBI is now determined in the same way for all benthic assessment approaches, which was not the case for the previous exercises (Borja et al., 2007). These recalculations have led to slightly different EQR values for the samples of the common dataset compared to the previous analyses. The advantage of this is that the analyses were standardized, transparent and are repeatable in time. The WFD requires the inclusion of certain metrics within the national assessment method for benthic invertebrates, which are summarized for each Member State in Table 3.

Table 2. Overview of the algorithms of the NEA-GIG benthic invertebrate indicators for intercalibration.

MULTIMETRIC		
BEQI (Belgium)	EQR=average (EQR species+ EQR density+ EQR similarity)	(Van Hoey et al., 2007) http://www.beqi.eu
DKI¹ (Denmark)	$\left(\frac{1 - \frac{AMBI}{7} + \left(\frac{H'}{Hmax} \right)}{2} \right) * \left(\frac{\left(1 - \frac{1}{N} \right) + \left(1 - \frac{1}{S} \right)}{2} \right)$	(Borja et al., 2007)
NQI² (Norwegian)	$\left(0.5 * \left(1 - \frac{AMBI}{7} \right) + \left(0.5 \frac{SN}{2.7} * \frac{N}{N + 5} \right) \right);$	(Rygg, 1985 and 2002)
IQI (UK, Ireland)	$IQI_{v,IV} = \left(\left(0.38 * \left(\frac{1 - (AMBI/7)}{1 - (AMBI_{Ref}/7)} \right) \right) + \left(0.08 * \left(\frac{1 - \lambda'}{1 - \lambda'_{Ref}} \right) \right) + \left(0.54 * \left(\frac{S}{S_{Ref}} \right)^{0.1} \right) - 0.4 \right) / 0.6$	Philips et al., 2014
BEQI2¹ (The Netherlands)	EQR (ecotope) = 1/3 * [Sass / Sref] + 1/3 * [H'ass / H'ref] + 1/3 * [(6 - AMBI _{ass}) / (6 - AMBI _{ref})]	Van Loon et al., 2015
BO2A (Andalusia [Spain])	$BO2A = \log_{10}((fAO^3 / (fA^3 + 1)) + 1) //$ $EQR_{BO2A} = (\log(2) - BO2A_{measured}) / (\log(2) - BO2A_{reference}).$	Dauvin & Ruellet, 2007
RAT (Portugal)	$RAT (EQR) = (ES10_B + 2 * BENTIX_B + 2 * BENTIX) / 5$	Vinagre et al., 2017

MULTIVARIATE

M-AMBI

**(Basque
[Spain],
Cantabria
[Spain],
France,
Germany)**

Factor analysis: S, AMBI, Shannon diversity index¹

(Borja et al., 2004 and
Muxika et al., 2007)
<http://ambi.azti.es>

**BAT
(Portugal)**

Factor analysis⁴: AMBI, Margaleff diversity index, Shannon
diversity index¹

Teixeira et al., 2009;
Marques et al., 2009

¹DKI, BEQI2, m-AMBI, BAT: Shannon diversity: log base 2.

²NQI: $SN = \frac{\ln(S)}{\ln(\ln(N))}$; 2.7 is the ref value for SN

³fAO = frequency opportunistic annelid (fpo = frequency opportunistic polychaeta + fo = frequency oligochaeta) and fA = frequency amphipods

⁴Factor analysis BAT in *Statgraphics Plus 5.1 (rotation=varimax)*

The BEQI assessment approach does not allow a calculation of EQR values at sample level, due to the fact that it acts on habitat or water body level (Van Hoey et al., 2007; 2013). For the calculation of BEQI EQR values, a set of samples need to be considered for the assessment. Therefore, a separate comparison of the BEQI approach at higher level with the other benthic assessment approaches is executed (see separate intercalibration in phase I). Therefore, the samples of the other Member States are grouped per 10 (ideally), but at least to 5, to allow a BEQI calculation. The grouping of the samples is done, based on the fact that they are from the same site and same time (or time period). The EQR values of the pooled samples are based on the average value of the individual sample EQR's. The BEQI assessment approach determines the difference between a set of assessment and reference samples and classifies this according to the five classes of the WFD. The set of reference samples needs to be country/area/habitat specific; for this reason, the set of benchmark samples per country out of the common dataset is used as the set of reference samples. In this way, the principle of the BEQI approach is intercalibrated with the other benthic assessment approaches.

The RAT assessment approach is no comparable with the rest of methods because is applied on rocky substratum

Table 3. Overview of the metrics included in the national assessment methods

Member state	Full BQE method	Taxonomic composition	Abundance	Disturbance sensitive taxa	Diversity	Bio-mass	Taxa indicative of pollution	Combination rule of metrics
Belgium	Yes	Yes, species composition by means of Bray Curtis similarity	yes	As species composition without pre-classifying species in classes.	Yes, number of species	Yes	Specific opportunistic species	Average of the four parameters
Denmark	Yes	Not strictly – only as groups (5) of different sensitivity	Abundance is included as correction factor and relative abundance of different sensitivity groups and proportional abundance in Shannon Wiener index	5 sensitivity classes (AMBI)	Yes, number of species and Shannon wiener index	No	Specific opportunistic species	Algorithm
Netherlands	Yes	Not strictly – only as groups (5) of different sensitivity	As relative abundance of different sensitivity groups and proportional abundance in Shannon Wiener index	5 sensitivity classes (AMBI)	Yes, number of species and Shannon Wiener index	No	Specific opportunistic species	Average of 3 univariately normalized indicator EQR scores
Norway	Yes	Not strictly – only as groups (5) of	Species abundance as correction factor ($N_{tot}/N_{tot}+5$) and relative abundance of	5 sensitivity classes (AMBI)	Yes, number of species	No	Specific opportunistic species	Weighted algorithm: 50% AMBI and 50%

Member state	Full BQE method	Taxonomic composition	Abundance	Disturbance sensitive taxa	Diversity	Bio-mass	Taxa indicative of pollution	Combination rule of metrics
		different sensitivity	different sensitivity groups					number of species-abundance
Portugal-BAT method	Yes	Not strictly – only as groups (5) of different sensitivity	As relative abundance of different sensitivity groups and proportional abundance in Shannon Wiener index and Margalef index	5 sensitivity classes (AMBI)	Yes, Shannon Wiener index and Margalef index	No	Specific opportunistic species	Factorial analyses, calculating vectorial distances to reference conditions
Portugal-RAT method	Yes	Not strictly – only as groups of different sensitivity	As relative abundance of different sensitivity	Sensitivity classes (BENTIX)	Yes, Hulbert index	Yes	Specific opportunistic species	Algorithm
Spain (Basque Country; Cantabria); France, Germany	Yes	Not strictly – only as groups (5) of different sensitivity	As relative abundance of different sensitivity groups and proportional abundance in Shannon Wiener index	5 sensitivity classes (AMBI)	Yes, number of species and Shannon Wiener index	No	Specific opportunistic species	Factorial analyses, calculating vectorial distances to reference conditions
Spain (Andalusia)	No	Not strictly – only as groups (2) of	As relative abundance of opportunistic	2 sensitivity classes	No	No	Specific opportunistic species	No combination

Member state	Full BQE method	Taxonomic composition	Abundance	Disturbance sensitive taxa	Diversity	Bio-mass	Taxa indicative of pollution	Combination rule of metrics
		different sensitivity	polychaetes and amphipods	(sensitive or tolerant)				
United Kingdom / Ireland	Yes	Not strictly – only as groups (5) of different sensitivity	As relative abundance of different sensitivity groups and proportional abundance in Simpson index	5 sensitivity classes (AMBI)	Yes, number of species and Simpson index	No	Specific opportunistic species	Weighted algorithm: AMBI for 38%; Simpson for 8% and number of species 54%

A2.2 Sampling and data processing

The method of taking the benthic samples and processing for the WFD Monitoring within the different Member States is outlined in detail in annex 1. The information is extracted from the online WISER project database, which compiles all information regarding WFD assessment methods (version of Birk et al., 2010; <http://www.wiser.eu/results/method-database/>) excepting for the RAT method. This database is subjected to change: an update will probably be made in the near future.

In the case of the RAT method (Portugal), the description of sampling and data processing is as follows:

- Sampling time and frequency; summer, once each evaluation cycle
- Sampling method; quadrat technique 12x12 cm, 4 replicates per intertidal zone (upper, mid and lower), in total 12 replicates per site
- Data processing; biomass (g AFDW m⁻²) and density (ind m⁻²) estimated per species within each replicate
- Identification level; species level

A2.3 National reference conditions

The determination of the reference conditions is a delicate subject (Van Hoey et al., 2010; Birk et al., 2013). The ecological status in the WFD has to be measured as a deviation from a reference condition. These reference conditions need to correspond to largely undisturbed (=‘near-pristine’) conditions (no or minor impact from human activities). Indeed, the lack of appropriate reference sites or robust historical datasets is one of the major problems addressed in the intercalibration exercises and in setting the good ecological status boundaries (Borja et al., 2007; 2009). Scientists are faced with virtual lack of undisturbed sites along the European coasts and estuaries, and historical data are not easily accessible (Borja et al., 2004). Reference settings will need to be based on clear stressor-response relationships, a knowledge of the ‘naturalness’ of the system; and expert judgment may also have a role to play (Van Hoey et al., 2010). As summarized in Table 4, all Member States used the best available information (e.g. areas with least disturbed conditions) and their expert judgment to delineate appropriate reference values for their metrics.

The reference values used to calculate the EQR values for each sample in the common dataset are listed in Table 5. Those reference values were considered appropriate values for the samples of the subtidal soft-sediment habitats within the common dataset by each Member State. Those values were applied per benthic assessment approach on the entire common dataset.

Table 4. Overview of the methodologies used to derive the reference conditions for the national assessment methods included in the IC exercise

Member State	Type and period of reference conditions	Number of reference sites	Location of reference sites	Reference criteria used for selection of reference sites
Belgium ¹	Expert knowledge, Historical data, Least Disturbed Conditions. Data period 1994-2012 Habitat-specific	No reference sites; the reference data per habitat is selected out of the available benthos data collected over the period 1994-2012.	No reference sites	The most appropriate data for each benthic habitat of the BPNS as reference is based on the following selection criteria: <ul style="list-style-type: none"> - The data must be collected in the period 1994-2012 on the BPNS. - Data collected in areas where a certain human activity (dredge disposal, sand extraction, wind-farm construction) can disturb the natural variability of the benthic characteristics were excluded. - To have a good temporal and spatial coverage of samples within the reference dataset, we tried to have a balanced sampling (similar number of samples) over the years and within the areas of the BPNS.
Germany	Expert knowledge, Historical data, Least Disturbed Conditions; reference time: 1959 up to now. Habitat-specific	subtidal coast: 17	different sites Wadden Sea of Lower Saxony	The communities at the sites had to correspond with description of the reference community description referring to a certain habitat.
Denmark	Least Disturbed Conditions (Sites the least impacted - farthest from impact source);	Depends on type, but typically 5-50 sites	n.a.	Reference community and impact factor close to background.

Member State	Type and period of reference conditions	Number of reference sites	Location of reference sites	Reference criteria used for selection of reference sites
	Recent data from least impacted sites. Surface water type-specific			
France	Expert knowledge, Historical data, Least Disturbed Conditions. Data period : 1995-2006 Habitat-specific	Bretignolles_S Morlaix1_S SSMF06_S (Rade de Cherbourg)	Channel & Atlantic	Expert knowledge and least disturbed conditions. The reference conditions for 3 metric component M-AMBI were defined by habitat type, based on recent data (last decade) collated on sites of French Atlantic and English Channel coasts, in particular data collected as part of the French benthic network (REBENT: http://www.rebent.org/).
Netherlands	Historical data for 1991-2006; (b) AMBI(ref) estimated as the 1 percentile value; theoretical bad values: S(bad) = 0; H'(bad) = 0; AMBI(bad) = 6. (c) Statistical modelling for S(ref) and H'(ref): 99 percentile of S and H' for large ecotope dataset (highest indicator value which is robustly not an outlier).	Not true reference sites, but least disturbed sites can be selected if necessary.	The Wadden Coast and Wadden Sea are less impacted areas, compared to the Dutch Coast and Voordelta coastal zones.	Not applicable because coastal waters in The Netherlands are always subject to at least some level of anthropogenic impact. However, least disturbed samples from distinct sampling locations can be selected based on expert judgment using information on pressures at the sampling locations.
Norway	Recent data from least impacted sites	n/a	Outer coast of Skagerrak, southern Norway.	Reference sites were selected by the following criteria: Deeper than 5m, limited fresh water influence (> 1km from nearest estuary) and of sufficient distance (based on expert judgment) from any known pollution sources, such as large cities or industrial activity.

Member State	Type and period of reference conditions	Number of reference sites	Location of reference sites	Reference criteria used for selection of reference sites
Portugal-BAT method	Existing near-natural reference sites, Expert knowledge, Historical data, Least Disturbed Conditions; Data period Outer Minho (CW NEA1) – Feb and Jul 2006; Praia do Garrao (CW NEA26) - Apr and Nov 2006. Habitat-specific	14 sites (7 H/G and G/M sites, 2 historical data sites, 5 outfalls data)	Outer Minho (CW NEA1) – Reference site High-Good; Praia do Garrao (CW NEA26) - Reference site High-Good	Reference condition samples were identified as being from least disturbed conditions, selected on the basis of a) unimpacted sites; and b) from impact gradient study control sites. Reference condition values for Margalef, Shannon-Wiener and AMBI were identified from the data. Data was used from sites with low levels of natural disturbance and outliers (e.g. those with anomalously high taxa numbers in contrast to the remaining data) were identified according to expert judgment and excluded.
Portugal-RAT method	Least Disturbed Conditions (Sites the least impacted - farthest from impact source)	n/a	Less impacted site	Reference condition samples were identified as being from least disturbed conditions, Reference condition values for the indices included in the methods were identified from the data (For more details, see Vinagre, 2017)
Spain (Basque Country, Cantabria region)	Expert knowledge, Historical data, Modeling (extrapolation of model results); period 1995-2005. Habitat-specific	no specific number	Basque Country	Virtual locations, see: Muxika, I., A. Borja, J. Bald, 2007. Using historical data, expert judgment and multivariate analysis in assessing reference conditions and benthic ecological status, according to the European Water Framework Directive. Marine Pollution Bulletin, 55: 16-29.
Spain (Andalusia)	Least disturbed conditions. Habitat-specific	No real reference sites, only a benchmark site	In front of the Doñana National Reserve (site code: 51C0090, water body wise code:	Lowest impact of urban and industrial sewage and lowest amount of agriculture and urban land use.

Member State	Type and period of reference conditions	Number of reference sites	Location of reference sites	Reference criteria used for selection of reference sites
			510001, aprox. Coordinates (DD ETRS89: longitude - 6.601, latitude 36.879	
United Kingdom/ Ireland	Expert knowledge, Least Disturbed Conditions and Modeling (extrapolation of model results); Data from 1979 to 2003. Habitat-specific	No reference sites; >1000 sites from UK and Ireland are used for setting reference conditions		Reference condition samples were identified as being from least disturbed conditions, selected on the basis of a) expert judgement and b) from impact gradient study control sites. Reference condition values for AMBI, Simpsons and taxa number were identified from the data. Data was used from sites with low levels of natural disturbance and outliers (e.g. those with anomalously high taxa numbers in contrast to the remaining data) were identified according to expert judgement and excluded.

¹Changed compared to the WISER input, based on Van Hoey et al., 2014 report.

Table 5. Overview of the reference values per benthic characteristics used in the intercalibration exercise.

REFERENCE VALUES	Sample surface (m ²)	Number of taxa	Shannon (H' log2)	SN	Hulbert	Simpson	Margalef	AMBI	Density (Ind/m ²)	BENTIX	Biomass (gWW/m ²)	Bray Curtis similarity
Belgium	1 ⁷	153							2517.8		642.7	1
Germany	0.1	34	3.65					0.597				
Denmark	0.1 ²		5					0				
France	0.9 ⁴	58	4					1				
Ireland	0.1	68				0.97		0.96 ³				
Netherlands	0.078	31 ¹	3.8 ¹					0.01				
Norway	0.1			0.27								
Portugal-BAT method	0.1		4.1				5	0				
Portugal-RAT method	Quadrat technique 12X12				7					6		
Spain (Basque Country, Cantabria)	0.3 ⁵	42	4					1				

REFERENCE VALUES	Sample surface (m ²)	Number of taxa	Shannon (H' log2)	SN	Hulbert	Simpson	Margalef	AMBI	Density (Ind/m ²)	BENTIX	Biomass (gWW/m ²)	Bray Curtis similarity
Spain (Andalusia)												
United Kingdom ⁶	0.1	68				0.97		0.96 ³				

¹The values for the Netherlands are based on the combined reference value for the three Dutch coastal zones together.

²It is from circa 0.1m² obtained by pooling 6-7 smaller (ca 0.013m²)

³ $(1 - (\text{AMBI}/7)) = 0,96$

⁴9 replicates of 0.1 m²

⁵3 replicates of 0.1m²

⁶ these values are specifically set for fully marine subtidal muddy sand/sandy mud sediments from 0.1 m² grabs, sieved at 1 mm and using 2004 truncation rules

⁷The reference values are generated for each sample surface (from 0.1m² to max reference sample surface) based on a randomization procedure of the reference dataset for each boundary. The values shown in the table are those that generate an EQR value of 1. The values for the good/moderate boundary are 56 for number of taxa, 0.48 for Bray Curtis similarity, 179.9 and 6089.8 for biomass and 1182 and 7835 for density. In this report, the values for the muddy fine sand habitat for a sample surface of 1m² are reported.

A2.4 National boundary setting

The reported information in WISER regarding the boundary setting procedure for each Member State is summarized in Table 7. Most Member States reported that they take the boundaries from phase I intercalibration (Borja et al., 2007; 2009) and no specific approach for H/G or G/M boundary was reported in WISER. The boundary values used in the intercalibration for the different assessment approaches were summarized in Table 6.

Table 6. The boundary values (High/good and Good/moderate) for the different assessment approaches as used in the intercalibration exercise. BC: Basque Country, C: Cantabria, A: Andalusia.

National Method	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	the Netherlands	Germany	France	Spain (A)	Belgium
H/G	0,80	0,75	0,77	0,72	0,79	0,78	0,85	0,77	0,83	0,80
G/M	0,60	0,64	0,53	0,63	0,58	0,58	0,70	0,53	0,50	0,60
M/P	0,40	0,44	0,38	0,40	0,44	0,38	0,40	0,38	0,40	0,40
P/B	0,20	0,24	0,20	0,20	0,27	0,18	0,20	0,20	0,20	0,20

The Portuguese method for rocky substratum is not included in the current IC exercise, as it is not possible the comparison with the rest of methods based on soft bottom; its boundaries are H/G boundary= 0.80; G/M boundary=0.60

Table 7. Explanations for national boundary setting of the national methods included in the IC exercise

Member State	Type of boundary setting	Specific approach for H/G boundary	Specific approach for G/M boundary	BSP: method tested against pressure
Belgium	Equidistant division of the EQR gradient. The moderate/poor and poor/bad reference value were determined by equal scaling (respectively 2/3 and 1/3 of the good/moderate reference value).	The boundary setting procedure is based on the output of the randomization procedure of the reference dataset. The reference value for the high/good boundary is determined based on the median value (number of species, similarity) or the 25th and 75th percentile (density, biomass) out of	The boundary setting procedure is based on the output of the randomization procedure of the reference dataset. The reference value of the good/moderate boundary is determined based on the 5th percentile (number of species, similarity) or on the 2.5th and 97.5th percentile (density, biomass) out of the	

Member State	Type of boundary setting	Specific approach for H/G boundary	Specific approach for G/M boundary	BSP: method tested against pressure
		the permutation distribution.	permutation distribution of each parameter of the reference dataset.	
Germany	<p>Boundaries taken over from the intercalibration exercise (Borja et al., 2007¹). Calibrated against pre-classified sampling sites.</p> <p>The boundary setting procedure is in line with the WFD's normative definitions.</p>			The boundaries were additionally adjusted by the assessment of expert judgment (Heyer 2007). The m-AMBI relates to pressures of sediment enrichment, eutrophication and hazardous substances (Muxika et al. 2007).
Denmark	Equidistant division of the EQR gradient. Using discontinuities in the relationship of anthropogenic pressure and the biological response.		Usually, the border between good and moderate EcoQS (G/M) is determined as some deviation from a reference situation. Reference data, however, are difficult to find. An alternative procedure is described to estimate the G/M border, not requiring reference data. Threshold values, where faunal structure deterioration	

Member State	Type of boundary setting	Specific approach for H/G boundary	Specific approach for G/M boundary	BSP: method tested against pressure
			<p>commences, were identified from non-linear regressions between indices and impact factors. Index values from the less impacted side of the thresholds were assumed to come from environments of Good and High EcoQS, and the 5th percentile of these data was defined as the G/M border.</p>	
France	<p>Boundaries taken over from the intercalibration exercise (Borja et al., 2009) and calibrated against pre-classified sampling sites</p>			See: Borja et al., 2009.
Netherlands		<p>The Good/Moderate boundary of 0.58 is primarily derived from the initial G/M boundary for sheltered coastal waters (Wadden Sea), which was estimated</p>		

Member State	Type of boundary setting	Specific approach for H/G boundary	Specific approach for G/M boundary	BSP: method tested against pressure
		using expert judgment and set at 0.58.		
Norway	National boundaries (Molvær et al., 1993) adjusted following the intercalibration exercise (Borja et al., 2007)			
Portugal-BAT method	Boundaries taken over from the intercalibration exercise.			<p>AMBI ecological group proportions were established for samples over a pressure gradient (urban treated outfall). Initially, equidistant class boundaries were set and each AMBI EG proportion was calculated for i) the overall status and ii) the lower and upper quartiles of the data in each status. Where the AMBI EG proportions did not conform to those interpreted from the WFD Normative Definitions, the status boundary was adjusted towards the quartile that gave a more accurate</p>

Member State	Type of boundary setting	Specific approach for H/G boundary	Specific approach for G/M boundary	BSP: method tested against pressure
				representation. Boundaries were further optimized during Intercalibration Phase I.
Portugal-RAT method	Equidistant division of the EQR gradient			
Spain (Basque Country, Cantabria region)	Boundaries taken over from the intercalibration exercise (Borja et al., 2007)			Borja et al., 2009 & others.
Spain (Andalusia)	Using paired metrics approach, using the frequency of opportunistic annelid and the frequency of amphipods as metrics. Moderate/Status: amphipod frequency (except Jassa) less than 0.45, and opportunistic polychaete frequency higher than 0.55 - Poor/bad Status: amphipod frequency (except Jassa) less than 0.28, and opportunistic	Dauvin & Ruellet (2007) use the limits of the AMBI index (Borja et al., 2000) proposed by Borja et al.(2004) to theoretically calibrate BOPA limits: High/Good Status: amphipod frequency (except Jassa) between 1 and 0.80, and opportunistic polychaete frequency between 0 and 0.20.	Good/Moderate Status: amphipod frequency (except Jassa) less than 0.80, and opportunistic polychaete frequency higher than 0.20.	Yes, quantitative; The methods relates to a pressure gradient of eutrophication (nutrient and organic matter enrichment and discharges).

Member State	Type of boundary setting	Specific approach for H/G boundary	Specific approach for G/M boundary	BSP: method tested against pressure
	polychaete frequency higher than 0.72.			
United Kingdom/Ireland	Boundaries taken over from the intercalibration exercise (Borja et al., 2007 ¹).			<p>AMBI ecological group proportions were established for samples over a sewage sludge disposal pressure gradient. Initially, equidistant class boundaries were set and each AMBI EG proportion was calculated for i) the overall status and ii) the lower and upper quartiles of the data in each status. Where the AMBI EG proportions did not conform to those interpreted from the WFD Normative Definitions, the status boundary was adjusted towards the quartile that gave a more accurate representation. Boundaries were further optimized during Intercalibration Phase I.</p>

A2.5 Results of WFD compliance checking

Table 8. WFD compliance checking criteria.

Compliance criteria	Compliance checking conclusions
1. Ecological status is classified by one of five classes (high, good, moderate, poor and bad).	Yes, for all benthic assessment approaches
2. High, good and moderate ecological status are set in line with the WFD's normative definitions (Boundary setting procedure)	Yes, for all benthic assessment approaches
3. All relevant parameters indicative of the biological quality element are covered (see Table 1 in the IC Guidance). A combination rule to combine parameter assessment into BQE assessment has to be defined. If parameters are missing, Member States need to demonstrate that the method is sufficiently indicative of the status of the QE as a whole.	All Member States included the relevant parameters (see Table 3), except Spain-Andalusia. They do not include a diversity parameter (2011-12-16technical_report_NEA_CW_invertebrates_ES(AN)_Dec2011). A combination rule to combine parameter assessment is defined by all benthic assessment approaches.
4. Assessment is adapted to intercalibration common types that are defined in line with the typological requirements of the WFD Annex II and approved by WG ECOSTAT	Yes, for all Member States (see Table 9 and Table 10)
5. The water body is assessed against type-specific near-natural reference conditions	No. Alternative benchmark conditions (based on a "least disturbed condition" criteria) had to be defined due to the absence of near-natural reference conditions in the intercalibrated type.
6. Assessment results are expressed as EQRs	Yes, for all benthic assessment approaches (see Table 3).
7. Sampling procedure allows for represent-tative information about water body quality/ecological status in space and time	In most cases, the monitoring is considered as representative by the Member State itself (see annex 1). This aspect is not confirmed by specific, standardized analyses to test their representativeness. Sampling procedures are outlined in general, but not linked with the running WFD monitoring programs.

<p>8. All data relevant for assessing the biological parameters specified in the WFD's normative definitions are covered by the sampling procedure</p>	<p>Yes, for all benthic assessment approaches. The sampling procedure defined by each Member State allows the collection of species-abundance data (see annex 1), which is necessary to calculate all metrics of the different benthic assessment approaches.</p>
<p>9. Selected taxonomic level achieves adequate confidence and precision in classification</p>	<p>Yes, for all benthic assessment approaches, with some difference in taxonomic detail per Member State, but sufficient comparability (see annex 1). Taxonomy between Member States datasets is standardized for intercalibration purposes.</p>

There can be concluded that all compliance criteria were met for the benthic assessment approaches of Belgium, Germany, Denmark, France, United Kingdom/Ireland, the Netherlands, Norway, Portugal and Spain (Basque and Cantabria region) (Table 8). Only, the benthic assessment approach of the Andalusia region does not meet the requirements of compliance criteria N°3, due to the lack of a diversity parameter within their approach. However, a scientific justification for this is presented in their separate intercalibration document (2011-12-16technical_report_NEA_CW_invertebrates_ES(AN)_Dec2011) and accepted by review panel.

A3. Intercalibration feasibility checking

A3.1 Typology

In the NE Atlantic, seven basic intercalibration types have been agreed upon. In this report the type NEA1/26 is taken into account (see outline of characteristics in Table 9).

Table 9. NEA GIG Intercalibration Type NEA1/26

New Type ID	Name	Salinity	Tidal range (m)	Depth (m)	Current velocity (knots)	Exposure	Mixing	Residence time
CW - NEA1/26	Exposed or sheltered, euhaline, shallow	Fully saline (> 30)	Mesotidal (1 - 5)	Shallow (< 30)	Medium (1 - 3)	Exposed or sheltered	Fully mixed	Days

The types above occur in Member States' waters as detailed below in Table 10, and compromise all NEA-GIG countries except Sweden.

Table 10. Member States sharing types

Type	BE	DE	DK	ES	FR	IE	NL	NO	PT	SE	UK
CW - NEA1/26	X	X	X	X	X	X	X	X	X		X

For benthic invertebrates, all classification schemes intercalibrated relate only to the soft sediment infauna component. RAT method based on rocky substratum is not included in the current IC exercise. Differences occur in the reference conditions for the types; these are specific for the habitat type, and for some Member States (NL and DE), sometimes even specific for the water body. However, the basic metrics in each country's benthic assessment approach remain the same.

A3.2 Pressures addressed

A3.2.1 Sample level

All methods can show in one or another way, a certain response to certain pressures (Table 7). For benthic indicators also an abundant number of papers and reports are available that shows their pressure-response relation (e.g. Borja et al., 2009; Josefson et al., 2009; Fitch et al., 2014; and others). Therefore, it can be concluded that the response of a certain benthic assessment approach is slightly different from pressure to pressure type and from area to area. Unfortunately, no combined analyses has been made regarding the pressure-response relationship of the 10 benthic assessment approaches of the NEA-GIG region on a certain pressure dataset. Therefore, rather than summarizing the available literature regarding this subject, the pressure-response of the different benthic assessment

approaches is tested on a large pressure dataset out of the common dataset. This allows to a uniform comparison of the responses of the different benthic assessment approaches, instead of different independent comparisons.

An appropriate dataset for this exercise was the Garroch Head sewage sludge disposal ground data set of the UK (provided by Marine Scotland), which is a very large dataset (180 samples) that is already standardized for IC purposes and with accompanying quantitative pressure information (organic and metal pollution concentrations) available. The elements (nitrogen, carbon, copper, zinc, lead and chromium) are correlated with each other and are the explanatory variables for the pollution gradient at Garroch Head. In the further analyses and figures, Copper is used as proxy for the pollution gradient at Garroch Head, due to the fact that it shows the highest correlation with the benthic assessment approaches (Table 11).

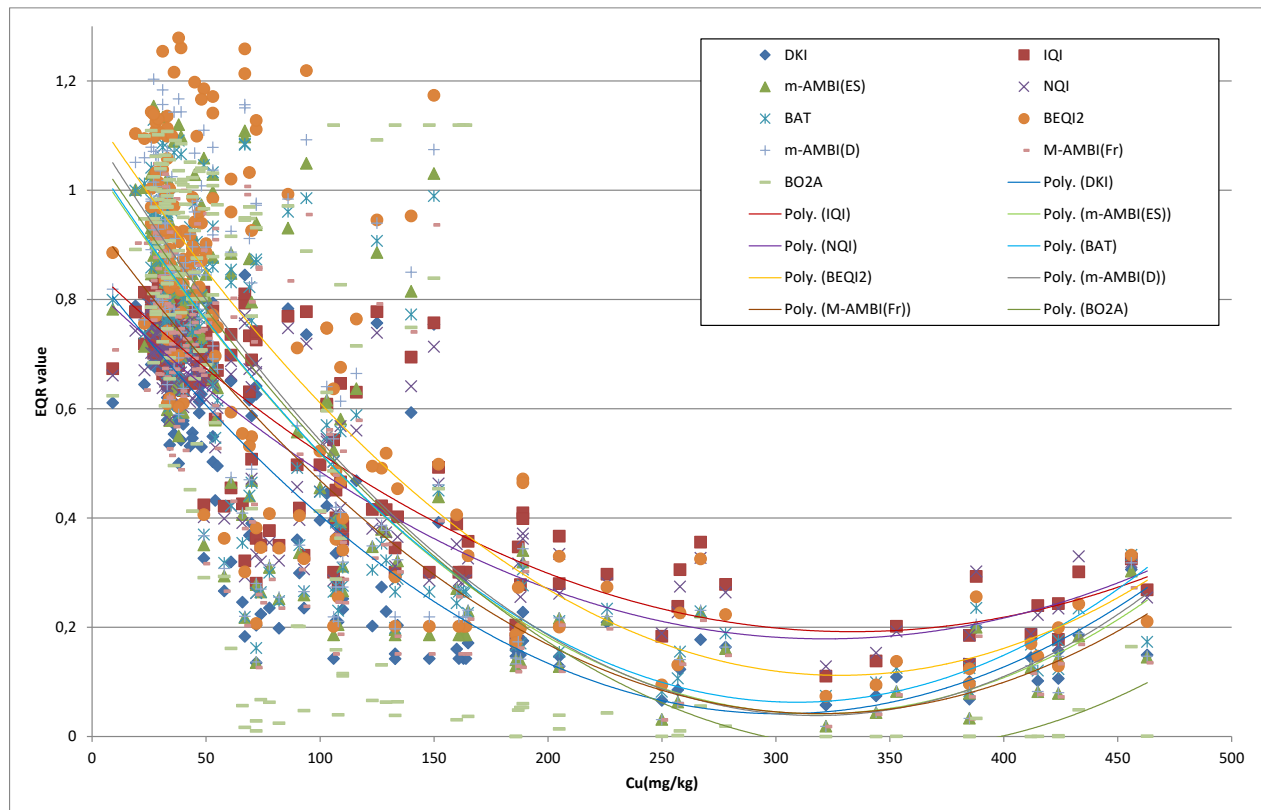


Figure 1. Correlation plot with trend line (polynomial 2nd order) between the different assessment approaches and Cu(mg/kg).

The different benthic assessment approaches shows no linear relation with the pollution gradient (copper), but a shift in benthic characteristics from 50-150mg/kg Cu (Figure 1). All benthic assessment approaches shows a clear and similar response to the pressure. Same, non-linear patterns in benthic characteristics against a metal pollution gradient were shown in the study of Josefson et al. (2009). All benthic assessments show a very similar correlation value with the pressure (Table 11). The highest correlation (cf Draftmans; Primer software) value is obtained with the IQI (UK/ROI) and the lowest with the BO2A (Spain, Andalusia).

Table 11. Draftmans plot correlation factors between benthic assessment approaches and organic and metal pollution parameters.

	Denmark	UK/ROI	Spain (BC, CC)	Norway	Portugal	Netherlan	Germany	France	Spain (AC)
N (%)	-0,681	-0,728	-0,692	-0,717	-0,684	-0,686	-0,693	-0,691	-0,580
Cu (mg/kg)	-0,729	-0,787	-0,732	-0,777	-0,728	-0,720	-0,735	-0,729	-0,672
Zn (mg/kg)	-0,704	-0,754	-0,710	-0,743	-0,704	-0,699	-0,712	-0,707	-0,632
Pb (mg/kg)	-0,621	-0,660	-0,636	-0,656	-0,633	-0,630	-0,638	-0,635	-0,572
C (%)	-0,701	-0,768	-0,719	-0,754	-0,708	-0,717	-0,720	-0,718	-0,628
Cr (mg/kg)	-0,692	-0,729	-0,696	-0,723	-0,694	-0,685	-0,699	-0,694	-0,624

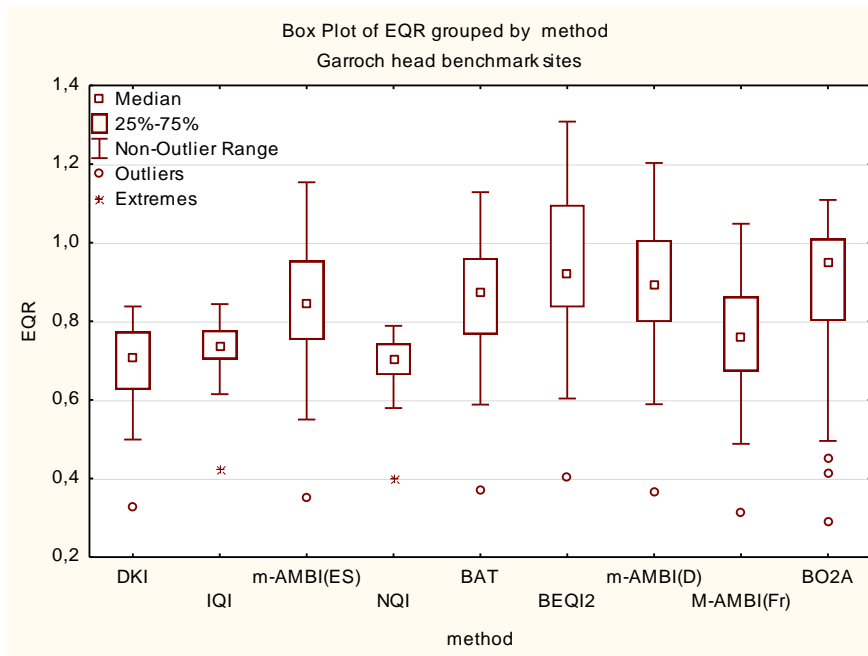


Figure 2. Box-Whisker plot of the EQR values of the benthic assessment approaches for the classification of the Garroch head benchmark sites.

The samples with a copper concentration less than 50mg/kg seem to represent non-disturbed conditions and could be used as benchmark sites (least disturbed samples). The box-whisker plot (Figure 2) gives a distinct visualization of the differences between the EQR values of the different benthic assessment approaches for these benchmark sites. Some approaches were more similar to each other than others. The median EQR values of the benchmark sites were a little bit lower for the DKI, IQI and NQI, which can be related to their higher reference values compared to the other approaches (Table 5). The BO2A shows the highest median EQR values for the benchmark sites. The values of m-AMBI (Fr) are in between. The m-AMBI (ES), BAT, BEQI 2 and m-AMBI(DE) EQR values were more or less similar for these benchmark sites. The differences of the EQR values of the benchmark sites were significantly different between the m-AMBI (ES), BAT, BEQI 2, m-AMBI(DE) and DKI, IQI, NQI and m-AMBI(Fr) (Kruskal-Wallis mean rank test) (Table 12). The IQI was not significantly different with the NQI, DKI and m-AMBI(Fr). The NQI was significantly different with all other approaches, except the IQI and DKI. The DKI is also significantly different with all other approaches, except the IQI and NQI. This to illustrate that there were differences in the benthic assessment approaches in the classification of the samples under similar pressure conditions. This benchmark aspect is further analyzed in point 4.3 below.

Table 12. Kruskal-Wallis p levels (multiple comparisons of mean ranks) by comparison the EQR values of each approach for the Garroch head benchmark sites.

	DKI	IQI	m-AMBI(ES)	NQI	BAT	BEQI2	m-AMBI(DE)	m-AMBI(Fr)	BO2A
DKI		1,000	0,000	1,000	0,000	0,000	0,000	0,009	0,000
IQI	1,000		0,000	1,000	0,000	0,000	0,000	1,000	0,000
m-AMBI(ES)	0,000	0,000		0,000	1,000	0,160	1,000	0,004	1,000
NQI	1,000	1,000	0,000		0,000	0,000	0,000	0,003	0,000
BAT	0,000	0,000	1,000	0,000		0,658	1,000	0,000	1,000
BEQI2	0,000	0,000	0,160	0,000	0,658		1,000	0,000	1,000
m-AMBI(DE)	0,000	0,000	1,000	0,000	1,000	1,000		0,000	1,000
m-AMBI(Fr)	0,009	1,000	0,004	0,003	0,000	0,000	0,000		0,000
BO2A	0,000	0,000	1,000	0,000	1,000	1,000	1,000	0,000	

A3.2.2 Higher level comparison

The samples of the Garroch head are grouped in sets of samples from the same location and same time period to allow a BEQI comparison. The reference dataset are the samples which are characterized by a copper content of less than 50 mg/kg. A similar trend of the benthic assessment approaches in relation to copper is found as on sample level (Figure 3). The EQR values decreased with increasing copper value. The BEQI approach shows a similar pattern as the other approaches.

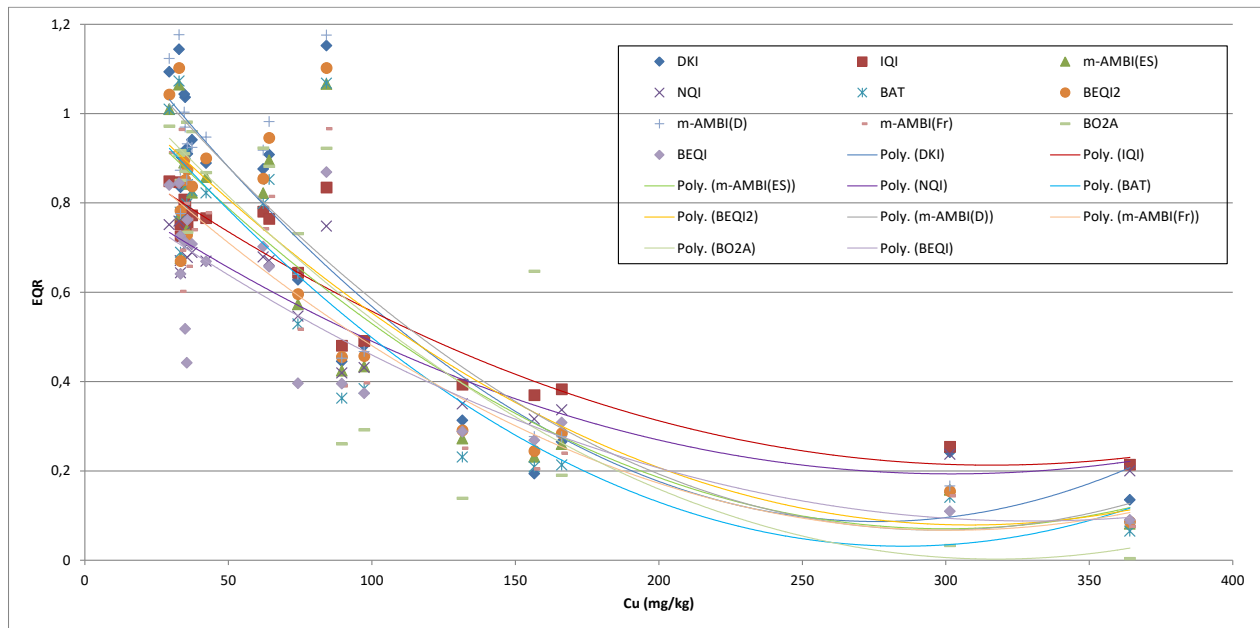


Figure 3. Correlation plot with trend line (polynomial 2nd order) between the different assessment approaches and Cu(mg/kg) for the set of pooled samples.

Table 13. Draftmans plot correlation factors between benthic assessment approaches and copper for the pooled samples.

	DKI	IQI	m-AMBI(BC, Q)	NQI	BAT	BEQI2	m-AMBI(D)	m-AMBI(Fr)	BO2A	BEQI
Cu	-0,810	-0,886	-0,817	-0,875	-0,808	-0,813	-0,823	-0,814	-0,828	-0,805

The correlation between the copper concentration and the EQR values of the benthic assessment approaches are all high and comparable (Table 13). The BEQI shows the lowest correlation; the IQI the highest.

RAT methods was not included in the analyses shown above. It was compared directly against anthropogenic disturbance pressure, against EQR values estimated by MarMAT at site level for the same sampling occasion (Fig. 4). The correlation between the 2 methods was high and significant.

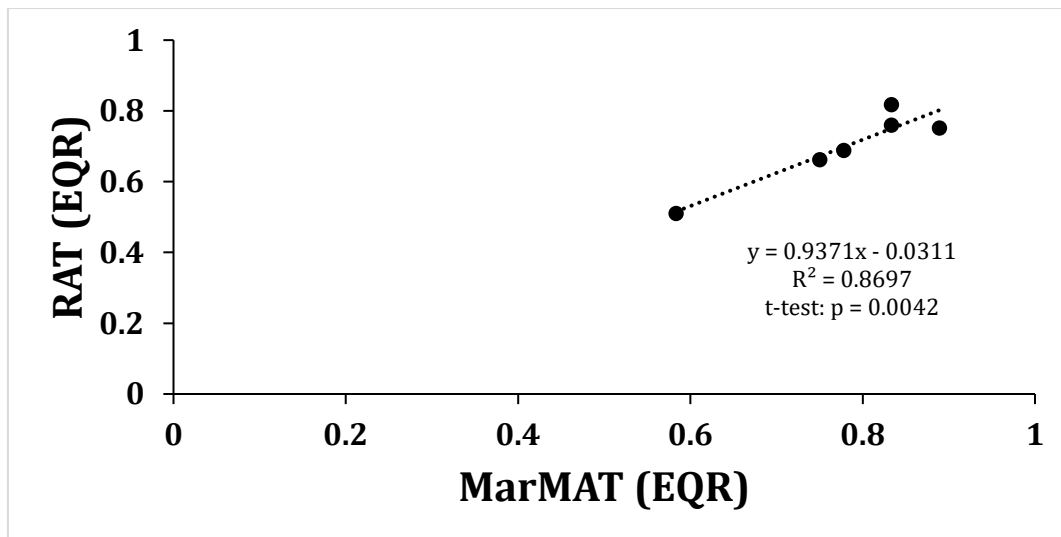


Figure 4. RAT method-pressure relationships

A.3.3 Assessment concept

Do all national methods follow a similar assessment concept?

The benthic assessment approaches within the NEA-GIG region are very similar, except the BEQI and the RAT method. Based on included metrics (parameters) and algorithms, those benthic assessment approaches can be grouped in 4 groups, as outlined in Table 14. The difference in the methodology of calculation of the BEQI (sample aggregation a prior to assessment), compared to the others (at samples level), led to the need for a separate comparability test. This comparability test is executed on aggregated set of samples out of the common dataset.

On the other hand, the RAT method is specific for rocky substratum, so is not possible the intercalibration with the rest of the assessment methods.

Table 14. The different types of benthic assessment approaches.

Method	Assessment concept	Remarks

Method group A: m-AMBI, BEQI2	<p>These approaches consist of similar parameters (AMBI, number of species and Shannon wiener), but a different algorithm (factorial analyses [m-AMBI] versus simple algorithm [BEQI2].</p> <p>The assessment is performed on sample level.</p>	
Method B: IQI, DKI, NQI, BAT	<p>These approaches consist of different parameters (AMBI, number of species, Shannon wiener, Simpson, Margaleff or abundance) and a different algorithm (factorial or simple algorithm).</p> <p>The assessment is performed on sample level.</p>	The simple algorithm differences are based on a different weighing of the parameters or using it as a correction factor (e.g. abundance)
Method C: BEQI	<p>Algorithm including number of species, abundance, (biomass), species composition (Bray-Curtis Similarity)</p> <p>The assessment is performed on habitat level (sample are aggregate prior to assessment).</p>	Difference in community characteristics, use of species composition index instead of a sensitive taxa proportion index.
Method D: BO2A	Based on the abundance of opportunistic polychaetes and amphipods; no diversity parameter.	Not fully WFD compliant
Method E: RAT	This method consists of different parameters (BENTIX and Hulbert index)	Boundaries calculated for rocky substratum

Conclusion

Is the Intercalibration feasible in terms of **assessment concepts**?

No identical approaches for the assessment, because they differs in their parameters or algorithm.

The majority of benthic assessment approaches (method type A, B and D) can be intercalibrated on sample level. The BEQI approach (Method type C) needs to be intercalibrated separately on an aggregated set of samples (habitat/ water body level), because this approach does not generate EQR values per sample. Therefore, this method is compared separately with the other assessment approaches on a higher level. In the case of RAT method is applied on rocky substratum, therefore is not possible comparability analyses with the rest of methods

Theoretical behavior of the different benthic assessment approaches

To better understand and illustrate the differences between the different assessment approaches, a test was run to show the dependency of the metrics (parameters) within each algorithm on the overall EQR score and the behavioral response of the different algorithms. This was done by running analyses on a fictive benthic dataset, where some metrics were gradually changed and others were kept fixed. Some of those theoretical samples do not occur in nature, but this exercise was intended to increase the insights into the different

algorithms of the benthic assessment approaches. The BO2A is not included, because it has no similar metrics compared to the other approaches.

As visualized in **Error! Reference source not found.5**, the different concepts show each some particularities, which can be summarized as follows:

- The approaches DKI, m-AMBI and BEQI2 shows a linear trend, when all metric values were slowly increased, whereas the NQI and IQI shows a more parabolic trend (decrease in EQR more strongly when low metric values were obtained). This type of pattern is related to the metric 'number of species' in both approaches.
- The behavior of the IQI is more complex. A decrease in number of species is buffered due to the transformation of the metric within the IQI, because the EQR values tend to decrease very slowly, except when low species numbers were reached. The IQI shows the highest dependency from the AMBI and the lowest for the Simpson.
- The DKI approach shows a linear pattern with increasing parameters, except for number of species (parabolic trend). This can be related to the correction factor $(1-1/S)$ in the algorithm, when the number of species (5-10) are low.
- The EQR values obtained by the m-AMBI approach seem to be most influenced by changes in the metric AMBI and less by the diversity parameters (number of species, Shannon wiener).
- The BEQI2 approach is equally dependent on the metrics, which is related to the equal weight that is given to those metrics within the algorithm.

It is obvious that those differences between the algorithms of the benthic assessment approaches are partly responsible for the variation in the scoring of the samples in the common dataset.

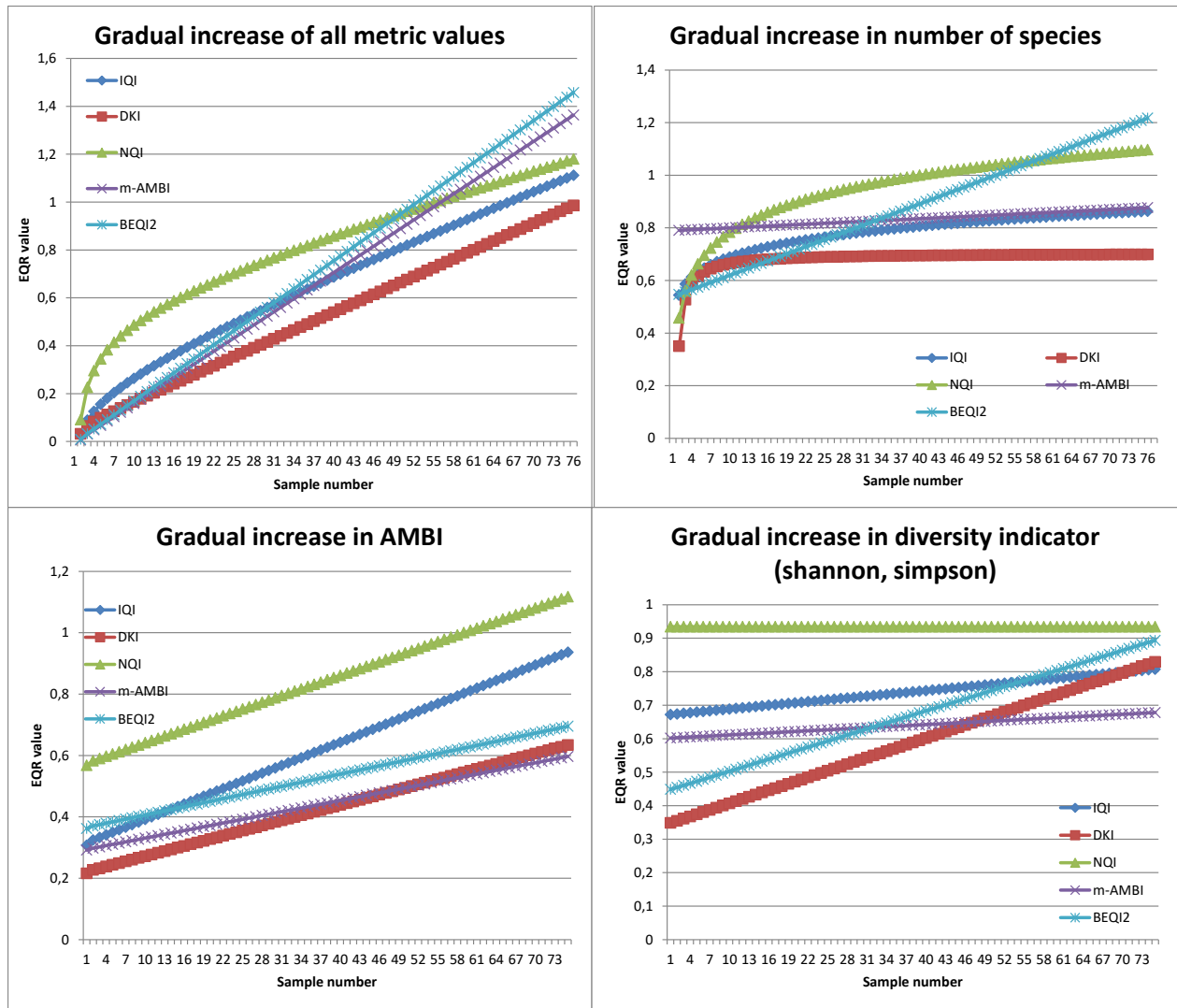


Figure 5. Changes in EQR values on fictive datasets, to show the metric dependency and behavioral response of the algorithm.

A.4 Collection of intercalibration dataset and benchmarking

A.4.1 Dataset description

The benthic dataset of phase I is used for the intercalibration, because there was no time foreseen in this action for collecting new data. Data from Portugal, the Netherlands and France was added in a later stage (Borja et al., 2009) (not included in the publication of Borja et al., 2007). The Danish data set was not from the NEA1/26 type, but the data came from the Kattegat but with rather similar physical characteristics. Nevertheless, the methods used within this type (NEA8a/9/10) were already intercalibrated. Therefore, they supplied new data, which include some NEA1/26 type data. According to the advice of JRC (Fuensanta Salas Herrero), only data of the NEA 1/26 type will be used for the further analyses. A part of the samples of Ireland were excluded (e.g., Clew Bay), due the incomparable sampling size [small]. These were the small modifications done on the common dataset in comparison to phase I. An overview of the metadata information of the final common NEA-GIG, type 1/26 benthic dataset is given in Table 15.

The NEA-GIG intercalibration dataset consists of 656 samples taken from Portugal to Norway. Most of the data originates from time series (samples at certain station sampled in time) or some from spatial monitoring (mainly the Belgian ones). There were 838 different taxa recorded in the entire database, which were constructed based on the 2004 UK taxonomical truncation rules.

A.4.2 Data acceptance criteria

All NEA-GIG Member States have delivered data for the intercalibration exercise. Nevertheless, the Spanish data is only from the Basque Country, because no data from the regions Andalusia and Cantabria was immediately available.

To explore the common intercalibration dataset for benthic macro-invertebrates, we performed some standard multivariate analyses. This to evaluate the following aspects:

- to check for outliers (samples very different from the rest and showing a problem)
- If there were regional or sub-regional differences between the samples
- If different benthic communities could be detected, which can be related to different physical habitats (sedimentology).
- If there is any pattern in the data that justifies the delineation of sub-types for benchmarking

Table 15. Sample description of data submitted by Member States, from the NEA-GIG for the intercalibration exercise. VV=van Veen grab; HC=Haps core; DG= Day grab; BC=Box core; SMI=Smith-McIntyre

Country	Location Code		Sample method	Sample size	Number of stations	Period	Replicates per station	Samples submitted	Depth (m)	Sediment	
B	Belgium	BGP	Station P2	VV	0,1026	1	1995	1	1	6,7	Sand (97%)
B	Belgium	BHA	Stations Habitat1999	VV	0,1026	37	1999	1	37	5-15	Sand (85%)-Mud(15%)
B	Belgium	BHA	Stations Habitat2000	VV	0,1026	12	2000	1	12	5-15	
B	Belgium	BMA	Stations Marebass	VV	0,1026	1	2000	1	1	13,8	Sand(30%)-Mud(70%)
B	Belgium	BMO	Stations M&OD	VV	0,125	6	1996	1	6	14,2	Sand(>99%)
B	Belgium	BOP	Station O&P	VV	0,125	17	1994,2	1	17	3,3	Sand(>97%)
B	Belgium	BSU	Subtidale stations	VV	0,1026	58	2002	1	58	5-10	Sand(>93%)
DK(NS)	Denmark	Jammerb	Jammerbugten	HC	0,1*	3	1995	3	3	4-10	Fine sand
DK(NS)	Denmark	Skagerra	Skagerrak	HC	0,1*	3	2004	3	15	8-20	Fine sand
D	Germany	VOR	NS2 Vortrapptief	VV	0,1	1	1987-2004	3-5	64	13	Sand (94%)
NL	the	Ems-	Ems-Wadden coast	BC	0,078	6	2000-2003	1	24	<20	Muddy sand
NL	the	Holland	Holland coast	BC	0,078	5	2000-2003	1	20	<20	Muddy sand
NL	the	Voordelt	Voordelta	BC	0,078	4	2000-2003	1	16	<20	Muddy sand
PT	Portugal	E	Ericeira	SMI	0,1	9	2001	1	9	10-30	Very fine sand
PT	Portugal	FF	Figueira da Foz	SMI	0,1	3	2002	1	3	10-30	Very fine sand
Fr	France	MORWI	Bay of Vilaine	SMI	0,1	5	1992	3	15	<30	muddy fine sand
Fr	France	QUIW	Bay of Quiberon	SMI	0,1	8	2004	3-5	34	<30	muddy fine sand
UK	UK-	HAR	Harwich	DG	0,1	3	2004	5	15	6,4	Mud(85,3%)
UK	UK-	LIV	Liverpool Bay	DG	0,1	3	2004	5	15	5,7	Sand(70%)-Mud(30%)
UK	UK-	MIL	Milford Haven	DG	0,1	3	2004	5	15	4,6	Mud(78,8%)
UK	UK-	TRB	Torbay	DG	0,1	3	2004	5	15	13,7	Muddy sand
UK	UK-	KIL	Kibrannan Sound	DG	0,1	1	2004	10	10	50	soft muds
UK	UK-	GRK	Garroch Head	VV	0,1	10	1979-1998	1	180	69-180	Silt/Clay
E	Spain	SSO	San Sebastian-Pasaia	BC	0,186	9	2000-2004	(combined)	45	33-61	Sand(90%)-Mud(10%)
N	Norway	STA	Stavanger(S5A)	VV	0,1	1	1995	4	4	93	Mud(83%)
N	Norway	TRO	Trondheimsfjord (RAH1)	VV	0,1	1	2001	4	4	50	Mud(88%)
N	Norway	UTN	Utnes (U10)	VV	0,1	1	2001	4	4	38	Sand(89%)
ROI	R. of	GRE	Greatmans Bay	DG	0,1	1	2003	2	2	40,1	Muddy sand
ROI	R. of	KEN	Kenmmare River	DG	0,1	3	2003	4	12	45,9	Muddy sand

A.4.2.1 General multivariate analyses

For the purpose of the multivariate analyses, the common dataset is fourth root transformed to reduce the effect of very abundant species on the overall pattern. Beside this, the rare species (in less than 1% of the samples and with a maximum of 3 individuals) were excluded from these analyses to reduce the effect of rare species on the overall pattern. This lead to a reduced dataset with 576 taxa. The similarity between samples is determined by the Bray-Curtis similarity. The sample groups were determined based on a cluster analyses, with cut-off level at certain similarity level. Multidimensional scaling (MDS) is used to visualize the cluster groups. The analyses were executed in PRIMER6.

The first analyses revealed no obvious rarities, but only some outlier samples. Those samples were excluded for all further analyses.

- The samples of station 3 in the Voordelta (the Netherlands) show an inconsistent pattern (two of them show the lowest similarity in comparison with the rest (outliers); the other two were classified in different cluster groups, depending on the analyses. This rare pattern indicates a problem at this location.
- Station Marebass from Belgium was also directly classified separately from the rest. Also the HA99-93 sample from the Belgian dataset classified different from the related samples and can be considered as outlier.

The general multivariate analyses show the following patterns (Figure 6; Table 16):

- All data clearly grouped per Member State and even data region (North Sea, , when the cluster analyses were sliced at a similarity level of 11. Even if when slicing it further at similarity 15, the grouped data were further split per Member State .
- The North Sea area forms one cluster of samples (cluster h in Table 16), with the samples of Belgium, the Netherlands, Germany and Denmark. The Liverpool Bay samples shows a high similarity with those North Sea samples. Another large cluster group contains a part of the UK data (Garroch Head), the Spanish and Norwegian data. The other Member States (France, Portugal, Ireland) datasets form separate clusters (Table 16; **Error! Reference source not found.**).
- The data of most Member States clustered more or less together in the MDS plot, except the Portuguese data (cluster G), which were more scattered.
- A few samples of the Garroch Head dataset (cluster C) were also split from the others and were very similar. This because those samples contain very high densities of only one species (*Mediomastus fragilis*).

Table 16. Number of samples of each Member State in each cluster group (slice at similarity level 11).

slice11	B	D	DK(NS)	E	Fr	N	NL	PT	ROI	UK	MS/regio
a					34						Fr (QUIW)
b										55	UK(Har, Kill, MILI, TRB)
c										6	UK(GRK)
d				45		12				174	Spain, UK(GRK), N
e									14		ROI
f					15						Fr (MORWI)
g								11			PT
h	130	64	18				56	1		15	UK(liv), NL, DK(NS), D, B

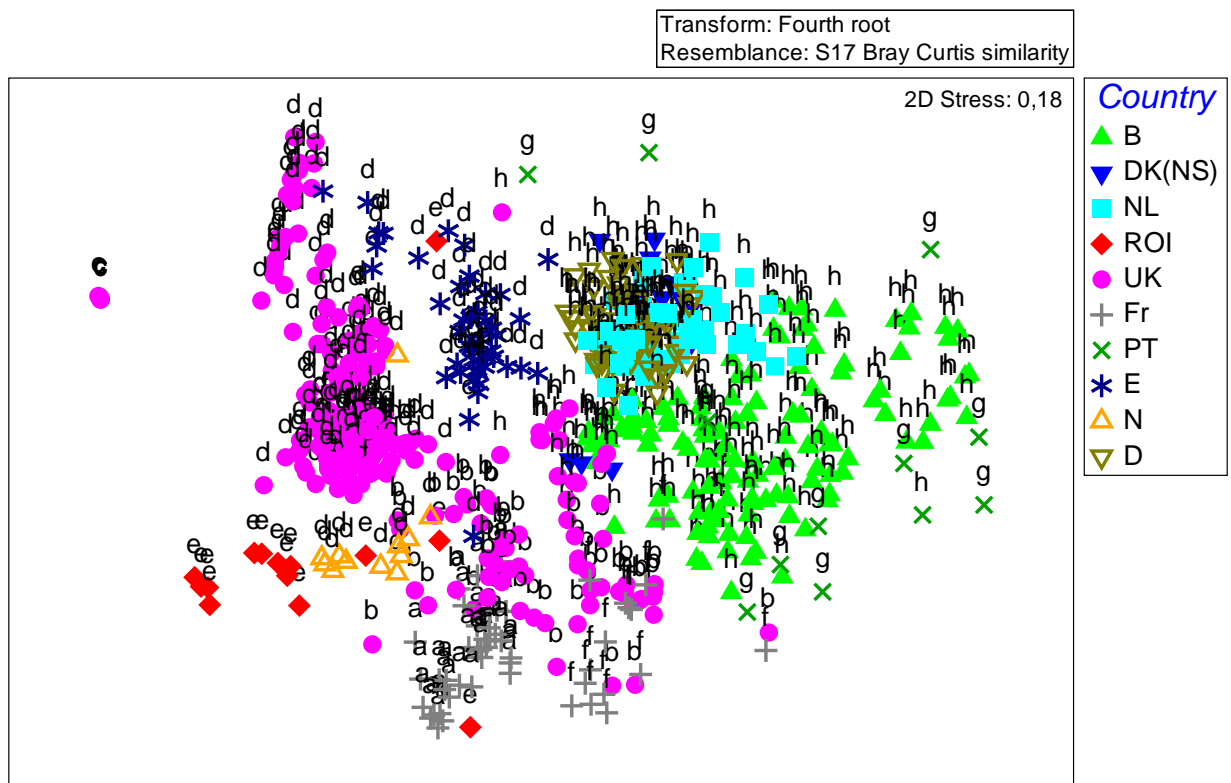


Figure 6. MDS plot of the intercalibration data with indication of the Member States (colored symbols) and the cluster groups (slice at similarity 11)

It can be concluded that the different datasets show a low similarity with each other, because they are clearly split as separate identities at low similarity level. There is no clear grouping of the data in relation to a South-North gradient within the NEA-GIG region. The data seemed to be grouped in a group with the North Sea related datasets and Portugal; a group with the datasets from shallow areas in the UK and France and a group with samples from less shallow areas (>30m depth) of UK, Spain, Norway and Ireland (Table 15). As the analyses show, every region has its own benthic species composition, with commonalities over the NEA-GIG region. The main difference in species composition between the NEA-GIG samples seems in first instance to be related to depth, which can be used as a factor to

delineate sub-regions in the intercalibration. The delineation of sub-regions based on biogeographical reasons (North-South) seems not to be appropriate.

A.4.2.2 Multivariate analyses of the benthic univariate parameters

Species composition on its own is not a parameter that is included in the benthic assessment algorithms. The algorithms are constructed from diversity and species tolerance/sensitivity classification metrics. In these analysis, it is investigated if those parameters are different among the Member States' datasets.

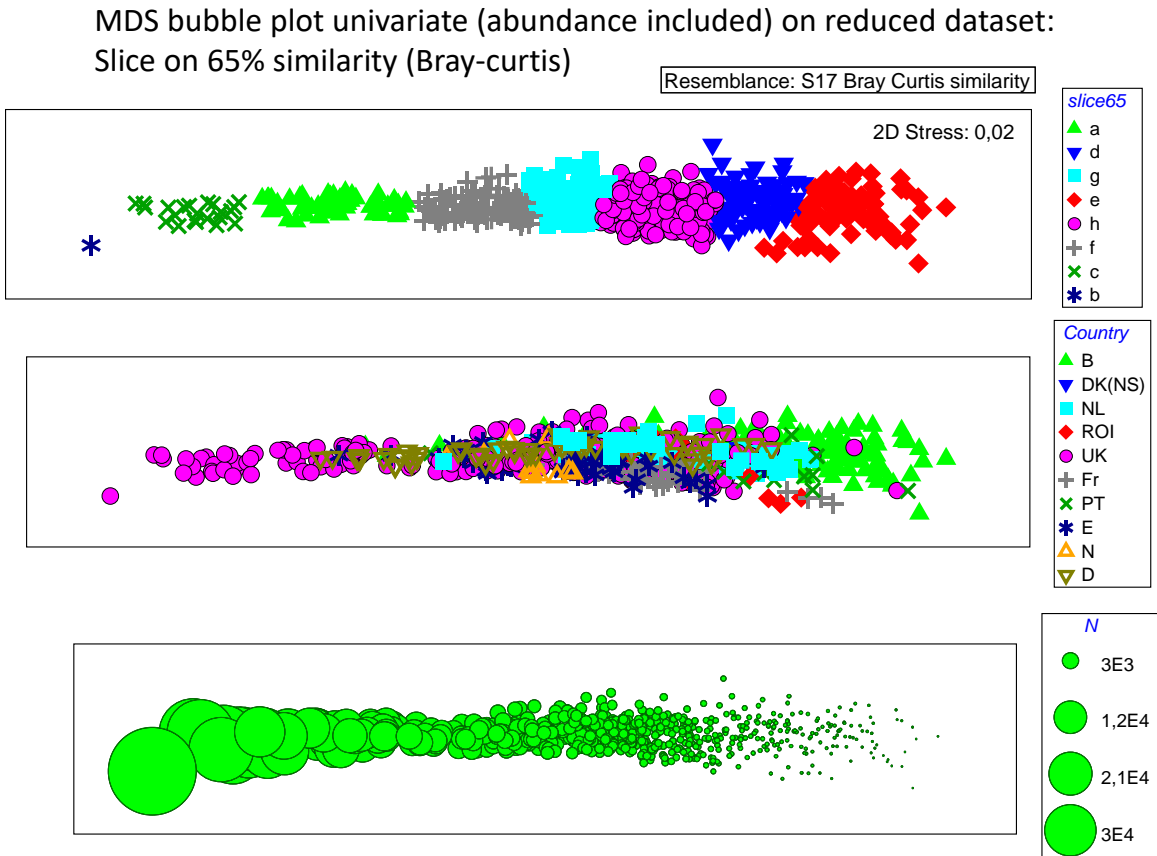


Figure 7. MDS plot of the univariate variables (inclusive abundance), with indication of the cluster groups at slice 65 (upper figure) and the behaviour of the dataset of the different Member States (center figure).

The figure below shows the pattern of abundance in the dataset.

The aim of this analysis, is to confirm if it is necessary to define sub-regions for the intercalibration by testing if there are differences between the samples in their univariate parameters/metrics (e.g., Shannon diversity [\log_{base2}], Margalef, Simpson, number of species, SN [$\ln(S)/\ln(\ln(N))$], abundance, AMBI). These are all the parameters by which the benthic assessment approaches are constructed.

The multivariate pattern is firstly strongly influenced by the parameter abundance. There is no obvious difference between countries and the samples are spread along the univariate gradients. It seems that many of the samples of the Belgian dataset are characterized by low abundances as compared to the other datasets. When abundance is excluded from the analyses, because it is in most approaches only relatively taken into account, the multivariate pattern is different. The gradient is dominated by number of species, and the deviation (at lower number of species) at one end is related to the difference in AMBI (very high values in the upwards gradient) (**Error! Reference source not found.8**). These analyses in the univariate parameters shows that there is a gradient within the dataset based on the univariate parameters from samples with a higher diversity to samples characterized by low diversity (Table 17). The data of the Member States seems to be spread over this gradient. This pattern in univariate parameters seems to correspond with a possible pressure gradient on the benthic data, which cannot be quantified (due to the lack of pressure data). The upper gradient shows the gradient in benthic characteristics, related to the disposal pressure (Garroch head, Spain), whereas the lower diversity gradient can be related to physical pressures (natural, anthropogenic).

MDS bubble plot univariate on reduced dataset: Slice on 75% similarity (Bray-curtis)

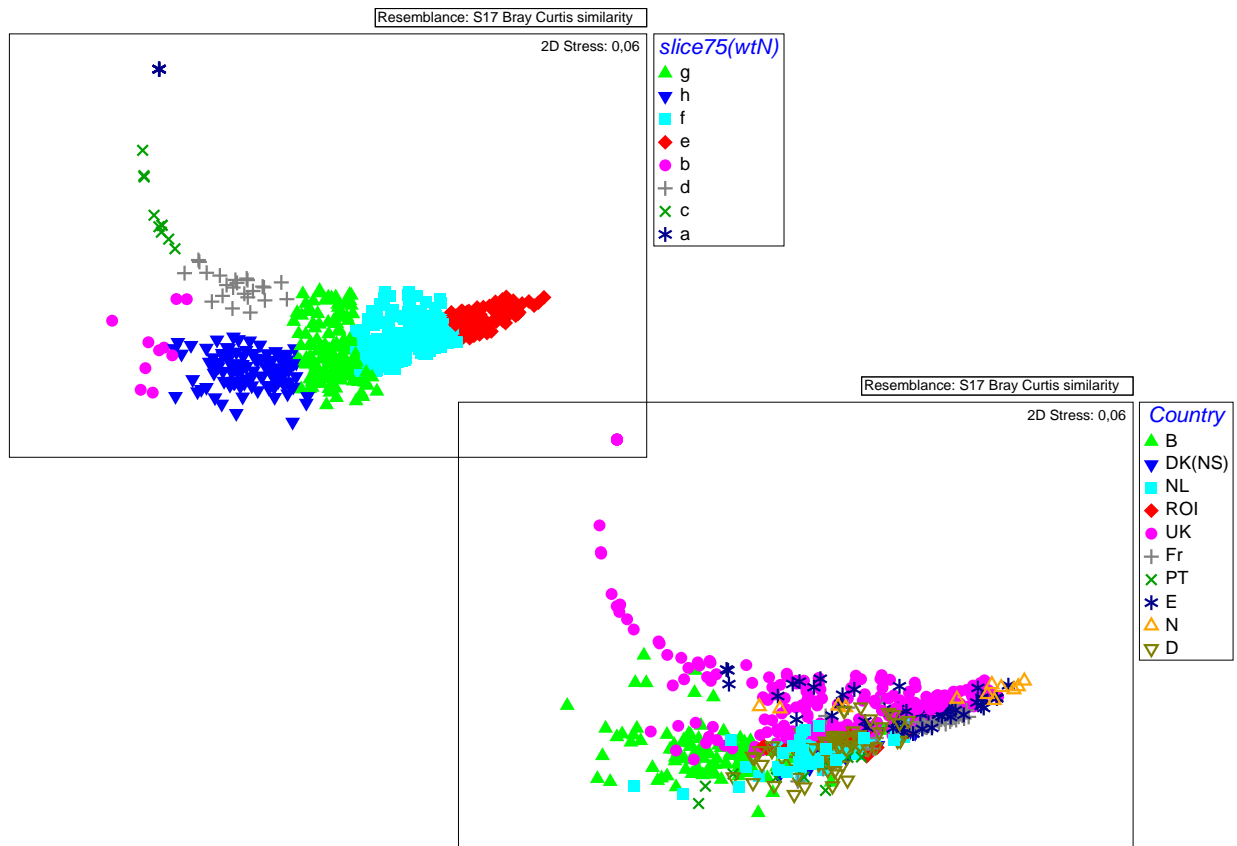


Figure 8. MDS plot of the univariate parameters (exclusive abundance and indication of the cluster groups (slice 75) (upper figure) and the behavior of the datasets of the different Member states (lower figure).

Table 17. Average values of the benthic parameters for each cluster group and their standard deviation.

Group	S		d		H'(log2)		1-Lambda'		AMBI		SN	
a	1,000	0,000	0,000	0,000	0,000	0,000	0,000	0,000	3,000	0,000	0,000	0,000
b	3,100	0,568	0,783	0,161	1,183	0,256	0,860	0,041	1,476	1,066	1,201	0,245
c	3,100	0,876	0,270	0,097	0,532	0,420	0,609	0,112	6,000	0,000	0,537	0,131
d	8,043	1,894	1,027	0,319	0,895	0,493	0,858	0,051	5,441	0,753	1,076	0,167
e	48,933	8,160	8,211	1,355	4,168	0,658	0,992	0,002	2,238	0,556	2,203	0,126
f	27,567	5,816	4,817	1,263	3,150	0,948	0,982	0,010	2,214	1,198	1,939	0,227
g	15,513	2,466	2,810	0,673	2,222	0,904	0,963	0,021	2,221	1,793	1,693	0,272
h	7,423	2,168	1,842	0,426	2,032	0,578	0,951	0,023	1,313	0,574	1,667	0,365

Overall conclusions:

All data are suited for the analysis, except the few outline samples discriminated. Based on the multi-variate analyses on the species-abundance data, we could discriminate the datasets from the different Member States, where the North Sea datasets show most similarity. The samples taken in less shallow regions (>30m) seem to be different regarding species composition compared to the samples taken in the more shallow regions. When this pattern is analyzed based on the metrics of the benthic assessment approaches, all datasets of the Member States are clustered together, but along a gradient. Therefore, no sub-regions based on biogeographical reasons can be discriminated. Only the factor depth seems to delimit two different type of habitats within the common dataset and can be considered as a relevant factor to distinguish between both dataset parts in the intercalibration. The review panel and JRC advise to distinguish this as two sub-types within the common dataset for the comparability analysis.

A4.3 Common benchmark

An alternative procedure for the selection of benchmark sites need to be used in this intercalibration, because we cannot fulfill the guidance principle using this common dataset: "The benchmarking process must use harmonized criteria independent of national classifications (i.e., countries cannot simply nominate the sites they classify as high status as being their benchmark sites without further checking)." The following approaches could be used for benchmarking, but does not make it within the NEA-GIG NEA1/26 intercalibration exercise:

- The absence of qualitative or quantitative pressure data (and it was not the task to collect this, which is an impossible exercise),
- no reference sites for each Member State /region (this approach was tried by Angel with sites from Spain and Norway),
- indirect pressure quantification not appropriate (e.g., LUSI index), due to the selection of data away from point sources (rivers, harbors, etc.) and the majority of the data is time series data from one location.

- An approach that estimates the benthic conditions under least disturbed circumstances could be the selection of samples with the highest diversity characteristics (response variables), which show a theoretical relation with changes in the abiotic environment due to pressures (see Annex 2). This procedure to determine the benchmark samples out of the common dataset is not accepted by JRC. The main reasons argued are, as stated in the IC Guidance, selection of benchmark sites should be done by screening for sites meeting abiotic criteria that represent a similar low level of impairment. The option proposed by the BQE lead for selecting benchmark sites is not acceptable because is based on the diversity, a biotic parameter included in most of the methods to be intercalibrated, and therefore the method values are influenced by this parameter. Moreover, in basis on the Pearson & Rosenberg model, diversity is a critical parameter, as it does not show a monotonic trend along both spatial and temporal gradients of pollution (Subida et al, 2013). When moving away from the source of pollution, the peak of opportunists is often followed by a maximum value in diversity, which then stabilizes at a slightly lower level. This means that, in a gradient of pollution, the highest values for the diversity index may be recorded when the number of species is still low and the community is still in an early stage of recovery (Pearson & Rosenberg, 1978). So, a diversity parameter, in some situations, could indicate high values in moderately disturbed areas.

A review panel argued that from a scientific perspective, the approach is not convincing and that the group should collect pressure data to do the benchmark standardization properly. JRC remained to the review panel the necessity to provide solutions in basis on the available data set. In this sense, JRC proposed to select benchmark sites in basis on the expert judgment.

Based on the knowledge of the coastal areas and the stations included in the dataset, they could indicate the stations that were under minor pressures (or with more distance from the focus of main disturbances) based in the following abiotic criteria:

- no harbours
- no beach regeneration
- no urban sewages
- no industrial sewages
- no fish farms
- no thermal industries
- no influence of agriculture activities
- >3 Km as a distance to the closer city with more than 1000 inhabitants

Therefore, the Member States indicate, the stations with minor pressures. For Spain (Basque country) and Norway, the benchmark sites selected during phase II were used: In the case of Norway because they have reference sites, and I this case of Spain (Basque Country) because they already selected in the previous phase less disturbed sites.

The review panel accepted this proposal.

A.4.3.1 Benchmark standardization

The principal aim of benchmarking in intercalibration is to identify and remove differences among national assessment methods that are not caused by anthropogenic pressure but rather by systematic discrepancies (due to different methodology, biogeography, typology etc.) (Annex V, IC Guidance).

Benchmark standardization will correct for differences in median EQR values between the Member States' benchmark sites obtained by certain assessment approaches. Those median values will be corrected by the benchmark standardization procedure; this correction will be more obvious for cases where the medians are significantly different.

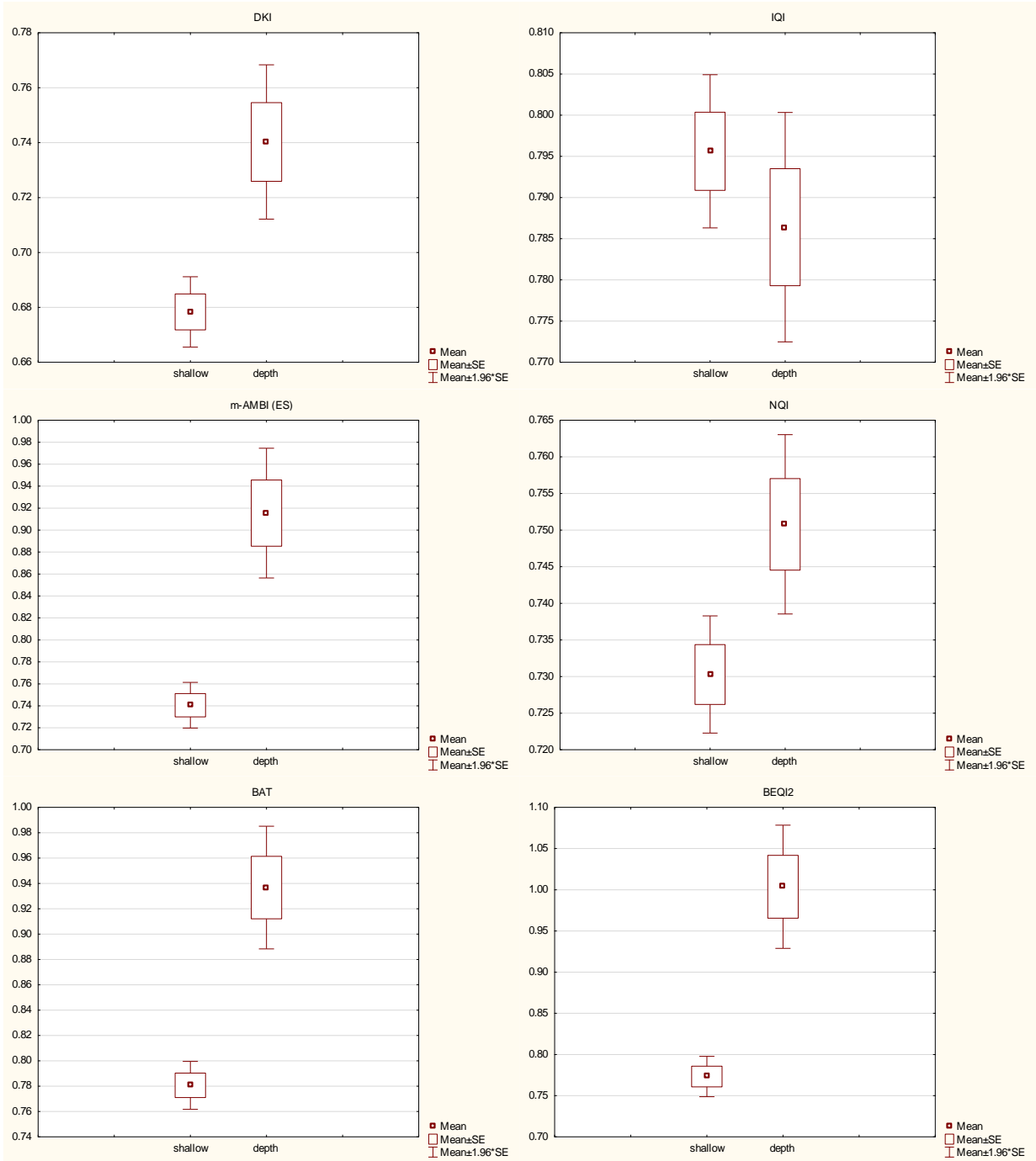
We tested whether benchmark standardization was necessary. Student's *sT* was used to compare the benchmark sites values for the two subtypes (shallow/depth) and the national methods.

There were statistical difference ($P < 0.05$) between both subtypes for all the methods, except for the IQI (UK/ROI method) (Table 18; **Error! Reference source not found.9**). Because of this, benchmark standardization was applied using the Excel sheet for option 3. The correlation between the average value of all national EQRs per survey in the full dataset was significantly correlated ($P < 0.01$) with its standard deviation, thus national EQRs converge towards the bad end of the quality gradient, and therefore, division was used for the standardization.

Benchmark samples were more than three national methods show EQR values less than good status (in accordance to the national boundaries) were excluded. This criteria was used in the previous phase by several MED GIG BQE groups. This were 8 samples of the Belgian dataset (station HA99-117; HA99-77, HA00-1; HAA00-11; HA00-21; HA00-3; HA00-4; HA00-5) and 3 samples of the German dataset (VORWI0700B [replica E]; VORWI0897B; VORWI0897B).

Table 18. Student's *sT* – P values

Method/Member State	P values
DKI (DK)	0.000046
IQI (UK/ROI)	0.28
m-AMBI(ES)	9.409E-07
NQI (NO)	0.0072
BAT (PT)	1.717E-07
BEQI2 (NL)	3.696E-07
m-AMBI(DE)	1.230E-06
m-AMBI(FR)	5.806E-07
BOA2A (ES)	7.884E-09



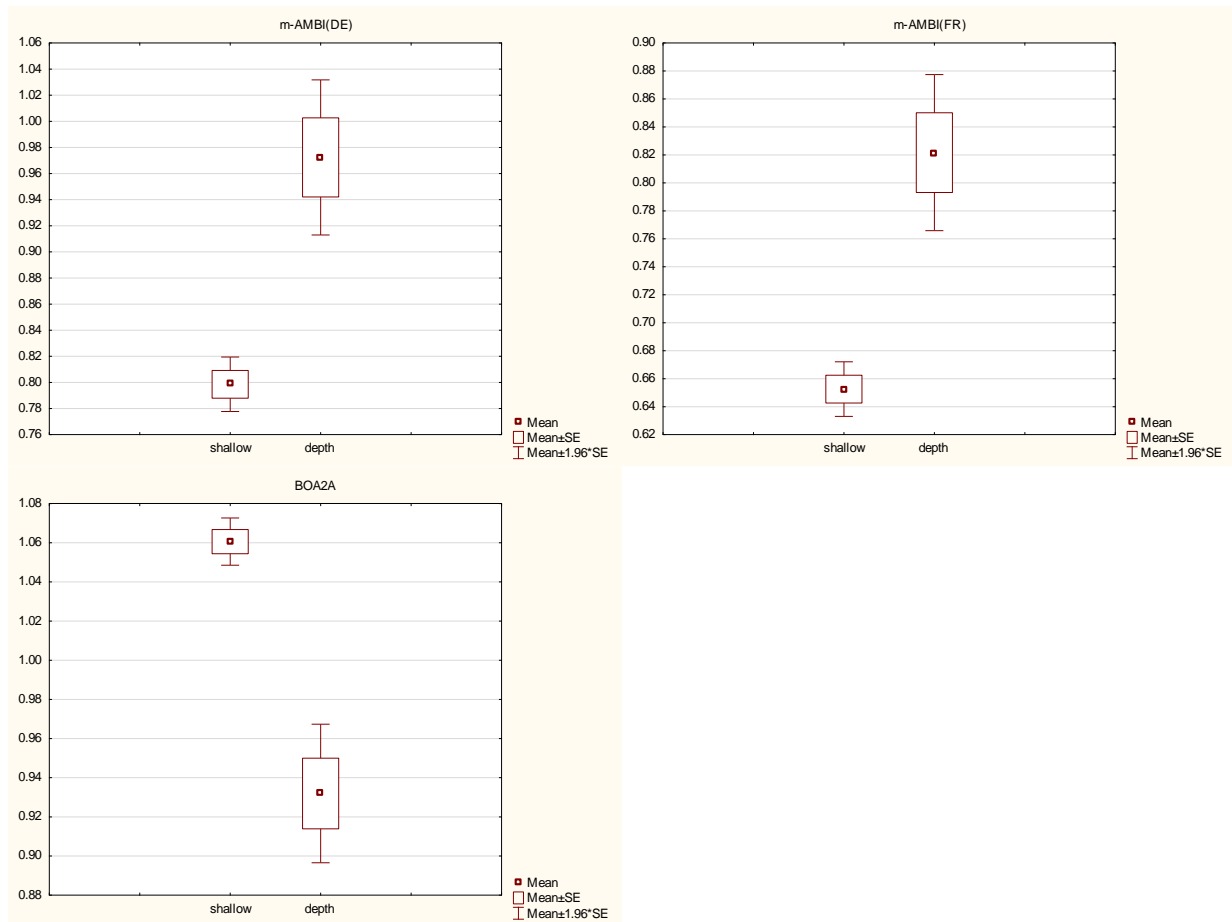


Figure 9. Box-whisker plot median, percentile values and no outlier range of the EQR values at the Member states benchmark sites with the national methods for the two subtypes (shallowness and depth).

A5. Comparison of methods and boundaries

A.5.1 Intercalibration option and common metrics

Option 3a. Intercalibration can be performed based on commonly assessed sites and whether the ecological quality gradient is sufficiently covered. More than three methods are used for this exercise. Following the advice of JRC and the review panel following intercalibration aspects need to be taken into account:

- The benchmark sites selected by the experts and following the review panel recommendations
- As benchmark standardization procedure, the division options is the appropriate one
- Two sub-types, based on depth, need to be distinguished.
- Due to the fact that the BO2A method does not meet the criteria in the previous comparisons, this method can be excluded in the final calculations.

The intercalibration excel sheet IC_Opt3_Div_v1.24.xlsx is used for executing the comparisons.

Because the BEQI assessment approach does not allow the calculation of EQR values on samples level (see 2.1 methods and 3.3 assessment concepts), a separate intercalibration on higher level (set of grouped samples) is executed. This separate intercalibration to analyze if the BEQI assessment approach meets the intercalibration criteria compared to the other assessment approaches. This separate comparability check on higher level implies that there no boundary adjustment could be suggested for the other assessment approaches based on those outcomes.

An intercalibration on sample level and higher level (to include the BEQI approach) was executed, with the benchmark samples selected based on expert judgment.

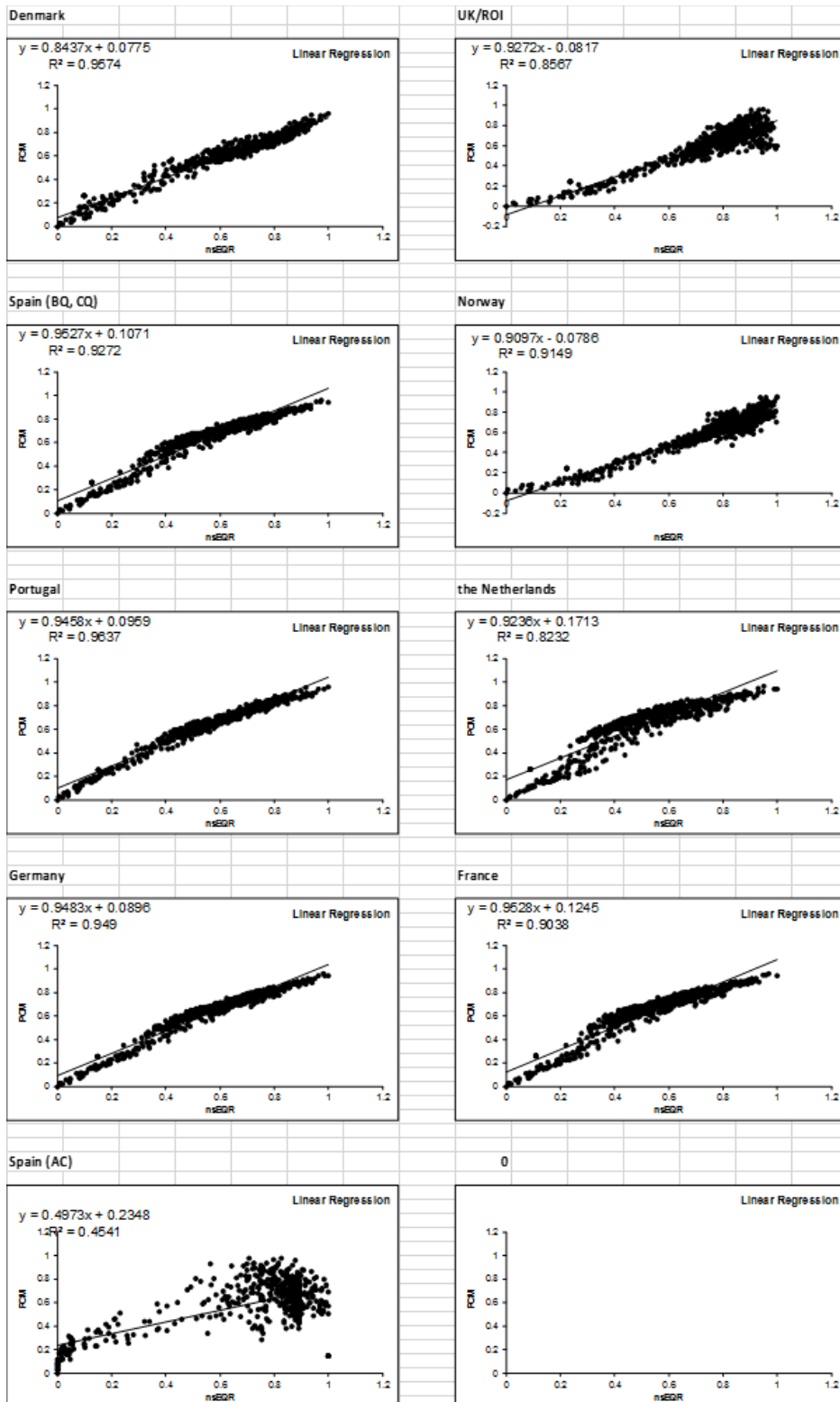
History

A set of comparisons between the benthic assessment approaches are executed during this third intercalibration phase. To keep record of it and to allow for checking which options were tested, this information is included in annex 3 of this report. This were all intermediate comparability analyses to explore the intercalibration and to guide towards the selection of the comparison most in line with the intercalibration guidelines and acceptable for JRC and the review panel.

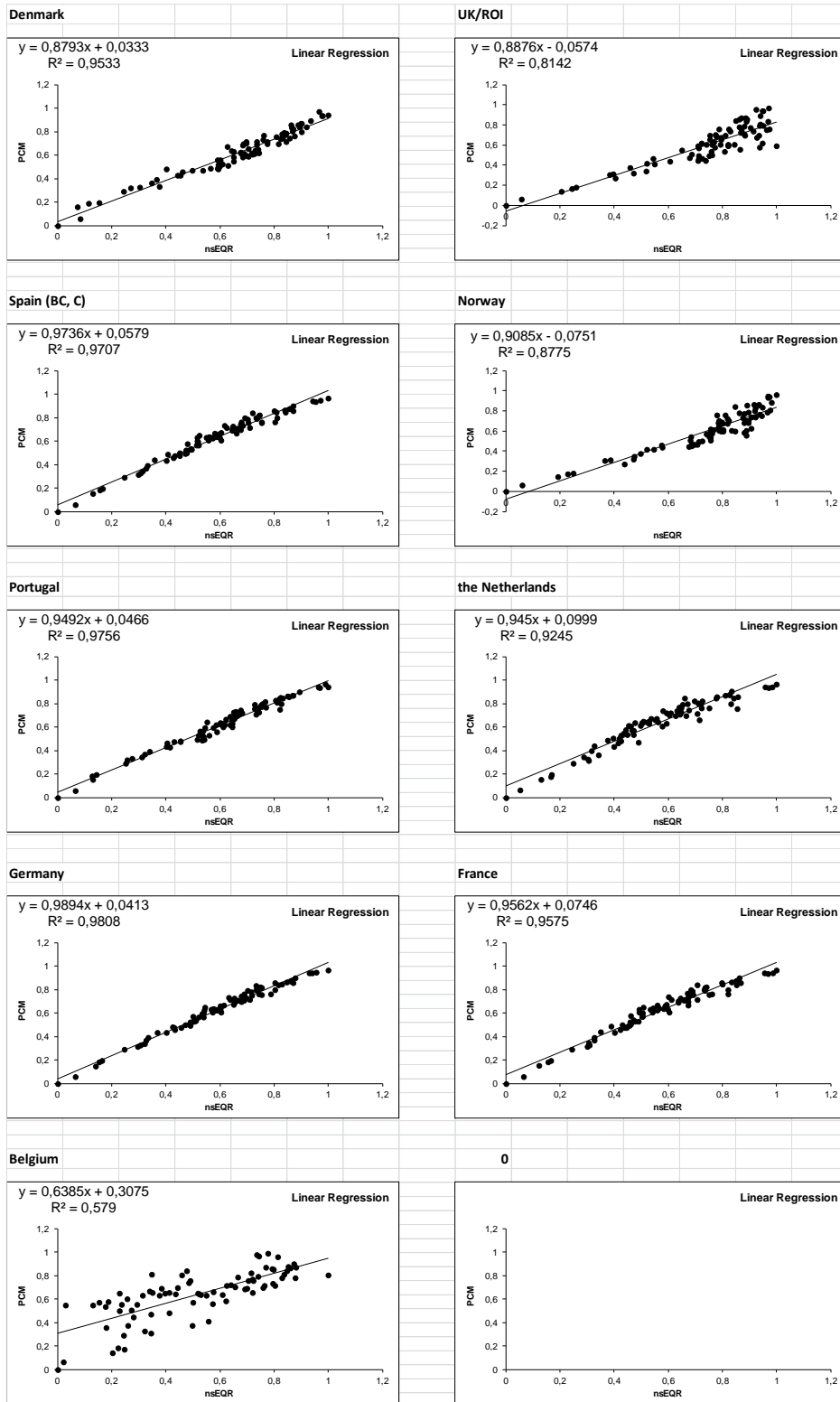
Different outcomes were obtained, based on the different options of benchmarking (biotic or expert judgment), standardization (subtraction or division), inclusion of methods (with or without BO2A), sub-regions (yes or no) and level of comparison (sample or higher). The use of these different options in the comparison lead to difference in the comparability criteria results and the need for boundary adjustments (or not). But the options selected for the final comparability analyses, seems to be the most appropriate regarding the intercalibration guidelines.

A.5.2 Results of the regression comparison

A.5.2.1 Sample level comparison



A.5.2.2. Higher level comparison (+ BEQI, Belgium)



Summary

The correlation between the metrics is determined in the intercalibration excel sheet. For all the intercalibration comparisons, the benthic assessment approaches fulfill the criteria ($R^2 < \frac{1}{2} \max R^2$) of the regression comparison (Table 19). The BO2A of Spain (Andalusia) shows the lowest correlation with the pseudo-common metric. For the IQI and the NQI, the samples were less equally spread over the linear regression line (dominance in upper part) in comparison to the other approaches, as was the case in the analyses on the theoretical behavior of the benthic assessment approaches.

Table 19. Summary of the correlation coefficient (R^2) of each approach with the common metric for the different intercalibration comparisons. Values outside the criteria were put in red.

Method	Sample level comparison	Higher level comparison
	Sub-region	Sub-region
Denmark	0.957	0.9533
UK/ROI	0.854	0.8142
Spain (BC, CR)	0.927	0.9707
Norway	0.914	0.8875
Portugal	0.963	0.9756
The Netherlands	0.823	0.9245
Germany	0.949	0.9808
France	0.903	0.9575
Spain (AC)	0.452	/
Belgium	/	0.579

The Spanish method (Andalusia region) had to be excluded from the comparability analysis due to its low correlation with the PCM ($r=0.452$).

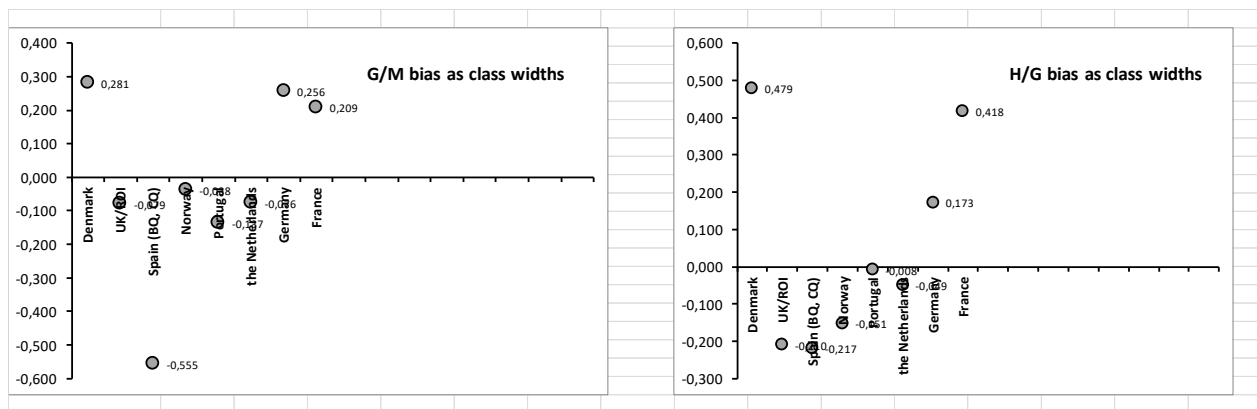
A.5.3 Comparability criteria

A.5.3.1 Sample level comparison

Table 20. Summary of the boundary bias and class differences analyses following the division benchmark standardization, with discrimination of the sub-regions.

	Denmark	UK/ROI	Spain (BQ, CQ)	Norway	Portugal	the Netherlands	Germany	France
Max	1,000	1,000	1,292	1,000	1,130	1,270	1,189	1,027
H/G	0,800	0,750	0,770	0,720	0,790	0,780	0,850	0,770
G/M	0,600	0,640	0,530	0,630	0,580	0,580	0,700	0,530
M/P	0,400	0,440	0,380	0,400	0,440	0,380	0,400	0,380
P/B	0,200	0,240	0,200	0,200	0,270	0,180	0,200	0,200
H/G bias_CW	0,479	-0,210	-0,217	-0,151	-0,008	-0,049	0,173	0,418
G/M bias_CW	0,281	-0,079	-0,555	-0,038	-0,137	-0,076	0,256	0,209

	Denmark	UK/ROI	Spain (BQ, CQ)	Norway	Portugal	the Netherlands	Germany	France
Count	4445	4445	4445	4445	4445	4445	4445	4445
Absolute Class Difference	0,4189	0,3735	0,2650	0,3582	0,2731	0,3028	0,3024	0,3042



For certain national methods do not comply with the comparability criteria. Boundary bias is exceeded by the methods of

- Denmark – Boundaries HG and GM too stringent
- Germany - Boundaries GM too stringent
- France - Boundaries too stringent
- Spain (BQ,Cantabrian)- Bundaries GM too relaxed

The average absolute class difference after boundary harmonization meets the comparability criteria for all national methods.

Spain is requested to adjust the boundaries to allow for completing the intercalibration exercise by raising its Good/moderate boundary to a value of 0.63.

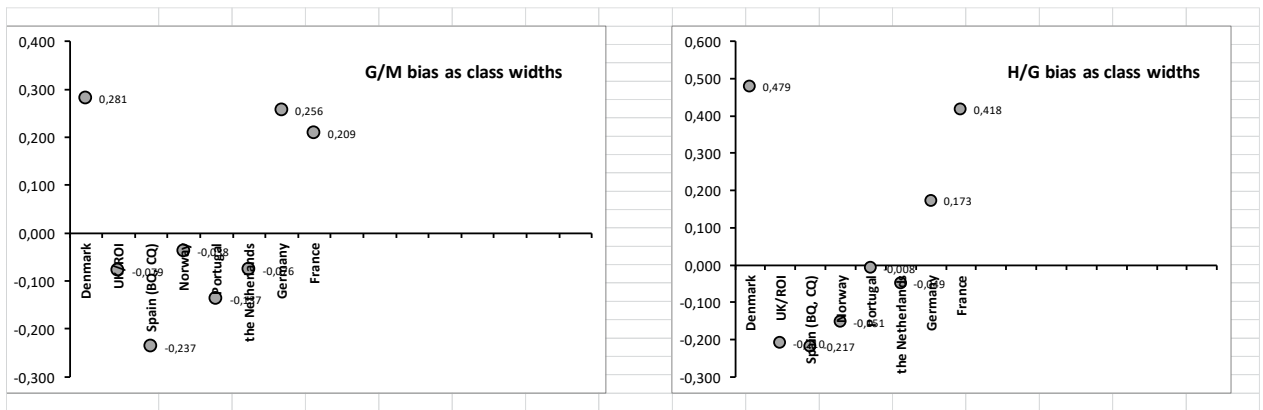
Germany, Denmark and France are not obliged to lower the boundaries that have been identified as being too stringent. The intercalibration criteria values after boundary harmonization are given in Table 21.

Table 21. Summary of the boundary bias and class differences analyses following the division benchmark standardization, with discrimination of the sub-regions, after harmonization of the boundaries.

	Denmark	UK/ROI	Spain (BQ, CQ)	Norway	Portugal	the Netherlands	Germany	France
Max	1,000	1,000	1,292	1,000	1,130	1,270	1,189	1,027
H/G	0,800	0,750	0,770	0,720	0,790	0,780	0,850	0,770
G/M	0,600	0,640	0,630	0,630	0,580	0,580	0,700	0,530
M/P	0,400	0,440	0,380	0,400	0,440	0,380	0,400	0,380
P/B	0,200	0,240	0,200	0,200	0,270	0,180	0,200	0,200

H/G bias_CW	0,479	-0,210	-0,217	-0,151	-0,008	-0,049	0,173	0,418
G/M bias_CW	0,281	-0,079	-0,237	-0,038	-0,137	-0,076	0,256	0,209

	Denmark	UK/ROI	Spain (BQ, CQ)	Norway	Portugal	the Netherlands	Germany	France
Count	4445	4445	4445	4445	4445	4445	4445	4445
Absolute Class Difference	0,4072	0,3874	0,2713	0,3748	0,2947	0,3087	0,2772	0,2893



A.5.3.2 Higher level comparison (BEQI, Belgium)

This higher level comparison is to test the comparability of the BEQI method with the other assessment approaches. Not meeting certain comparability criteria by these other assessment approaches has no consequence for the boundary harmonization (at sample level). The BEQI EQR values are determined on a set of stations (instead of one station).

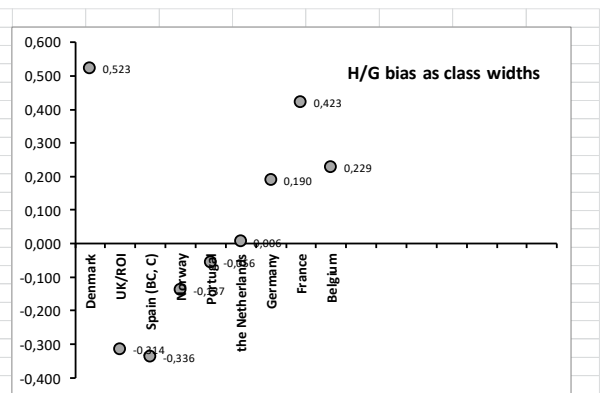
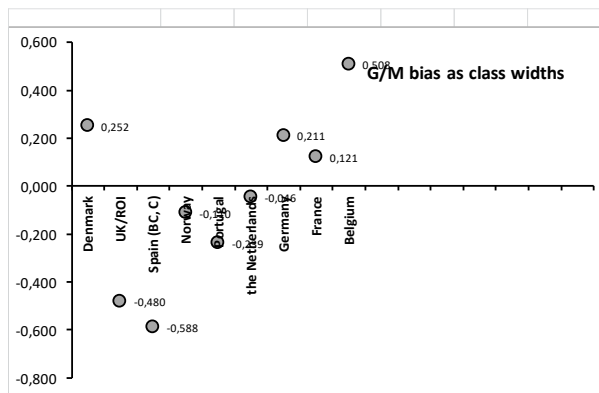
The boundary bias (<0.25) in this analysis is too high for the good/moderate and high/good boundary for the m-AMBI (BC, C) and IQI (Table 22). The DKI and BEQI (Belgium) are more stringent for the good/moderate boundary. The French and Danish approach is also more stringent for the high/good boundary. The class difference (<0.5 class) is below the criteria level for all benthic assessment approaches. The BEQI assessment approach meets the comparability criteria in comparison with the other approaches. Further boundary adjustment cannot be suggested, as this is a comparability check on higher level than sample level; in most assessment approaches, their boundaries were based on a sample level evaluation. Besides this, the BEQI is comparable with all methods applied in sub-region A (very shallow) type - all Belgian coastal waters belong to sub-region A.

Table 22. Summary of the boundary bias and class differences analyses following the division benchmark standardization, with discrimination of the sub-regions.

	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	e Netherland	Germany	France	Belgium
Max	1,000	1,000	1,229	1,000	1,016	1,049	1,040	1,000	1,000
H/G	0,800	0,750	0,770	0,720	0,790	0,780	0,850	0,770	0,800
G/M	0,600	0,640	0,630	0,630	0,580	0,580	0,700	0,530	0,600
M/P	0,400	0,440	0,380	0,400	0,440	0,380	0,400	0,380	0,400
P/B	0,200	0,240	0,200	0,200	0,270	0,180	0,200	0,200	0,200

H/G bias_CW	0,523	-0,314	-0,336	-0,137	-0,056	0,006	0,190	0,423	0,229
G/M bias_CW	0,252	-0,480	-0,588	-0,110	-0,239	-0,046	0,211	0,121	0,508

	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	e Netherland	Germany	France	Belgium
Count	648	648	648	648	648	648	648	648	648
Absolute Class Difference	0,4136	0,5278	0,3194	0,4228	0,3843	0,3611	0,3256	0,3210	0,4799



A.6 Final results to be included in the EC

Table with EQRs

After the boundary harmonization, the final boundaries for the benthic assessment approaches for coastal waters in the Northeast Atlantic (Common Type NEA 1/26) are given in Table 23. These results will be included in the Part I of the EC Decision. For the moment, only the BO2A and the RAT approaches does not meet the comparability criteria and their boundaries (Table 24) will be included in Part 2 the EC Decision.

Table 23. Boundary values of the different benthic assessment approaches after intercalibration. The boundaries in red are those changed after boundary harmonization. Results included in the Part I of the EC Decision.

Country	Benthic assessment approach	Ecological quality ratios			
		High-good boundary	Good-moderate boundary	Moderate-poor boundary	Poor-bad boundary
Denmark	DKI	0.8	0.6	0.4	0.2
France	m-AMBI	0.77	0.53	0.38	0.2
Germany	m-AMBI	0.85	0.70	0.4	0.2
Netherlands	BEQI2	0.78	0.58	0.38	0.18
Norway	NQI	0.72	0.63	0.4	0.2
Portugal	BAT	0.79	0.58	0.44	0.27
Spain (Basque Country and Cantabria)	m-AMBI	0.77	0.63	0.38	0.2
United Kingdom / Ireland	IQI	0.75	0.64	0.44	0.24
Belgium	BEQI	0.8	0.6	0.4	0.2

Table 244. Boundary values of the BO2A and RAT assessment methods. These methods have not been intercalibrated due to justified reason. Boundaries will be included in the Part II of the EC Decision.

Country	Benthic assessment approach	Ecological quality ratios			
		High-good boundary	Good-moderate boundary	Moderate-poor boundary	Poor-bad boundary
Portugal	RAT	0.8	0.6		
Spain (Andalusia)	BO2A	0.83	0.6	0.4	0.2

Correspondence common types versus national types

The common type (NEA1-26) is recognized as type in every Member State and is related to the national types.

A.7 Ecological characteristics

A.7.1 Description of reference or alternative benchmark communities

The description of the benthic community characteristics at reference or alternative benchmark is summarized in Table 25. This information is generated from the WISER database.

A.7.2 Description of good status communities

The description of the benthic community characteristics at good status is summarized in Table 25. This information is generated from the WISER database.

Table 25. Overview of the description by the Member States of the macro-invertebrate reference community and good status community

Member State	Description of reference community	Description of good status community
Belgium	The reference benthic characteristics of each habitat were defined on the randomization of a reference dataset, reflecting the spatial and temporal variability expected in that habitat, based on existing data and knowledge.	Is not defined textually.
Germany	Benthic communities, species numbers, diversity typically for the habitat (sediment, salinity, exposure)- low number of opportunistic species.	High portion of sensitive taxa, complex communities, low number of opportunists, high species number and high diversity assemblages.
Denmark	High diversity (H and richness). Dominance of sensitive species <i>sensu</i> Borja et al. 2000.	High diversity (H and richness). Dominance of sensitive species <i>sensu</i> Borja et al. (2000).
France	High diversity (H and richness). Dominance of pollution sensitive taxa <i>sensu</i> Borja et al., 2000.	Richness and diversity are slightly reduced in comparison to values under reference conditions, while variables according to habitat (community abundance as assessed by AMBI) are slightly unbalanced: sensitive taxa (EG I) abundance may range from high sub-dominant to absent; indifferent taxa (EG II) are of low sub-dominant abundance; tolerant taxa (EG III) of dominant abundance; abundance of opportunistic (EG IV) and indicator taxa (EG V) may range from negligible or low to comparable abundance with indifferent taxa (EG II).
Netherlands	level 3: reference community description is specific for each individual water body. Reference conditions based on historical data from 1970's. Furthermore a general description is given (in Dutch) in: STOWA (2009)	n.a.

Member State	Description of reference community	Description of good status community
	Referenties en maatlatten voor natuurlijke watertypen. report 2007-32	
Norway	n.a.	n.a.
Portugal-BAT method	Reference condition macrobenthic communities are dominated by pollution sensitive taxa (AMBI Ecological Group (EG) I taxa), have low relative abundance of indifferent (EG II) and tolerant (EG III) taxa and negligible relative abundance of opportunist (EG IV) and pollution indicator (EG V) taxa. High numbers of taxa with an even abundance distribution throughout the community is also indicative of reference conditions.	Community species richness (Margalef) and equitability (Shannon-Wiener) values are slightly reduced in comparison to values under reference conditions. While variable according to habitat, community composition (as assessed by AMBI) is slightly unbalanced. Community composition still dominated by EG I and II taxa. Slight reduction of sensitive taxa (EG I), and slight increase on tolerant taxa (EG III).
Portugal- RAT method	Benthic macroinvertebrate communities are characterized by the presence of species from EG I and II, such as <i>Acanthochitona</i> spp., <i>Chthamalus montagui</i> , <i>Dynamene bidentata</i> , <i>Melarhaphe neritoides</i> , <i>Patella depressa</i> , <i>Psammobiidae</i> , <i>Rissoa parva</i> and <i>Sabellaria alveolata</i>	Slight modifications on benthic macroinvertebrate communities are characterized by an increase on the abundance of tolerant species (EG III; e.g. <i>Mytilus galloprovincialis</i>) and opportunistic species (EG IV-V) (e.g. <i>Boccardia polybranchia</i> , <i>Polycirrus</i> sp.). Sensitive species, such as <i>Dynamene bidentata</i> and <i>Melarhaphe neritoides</i> , decrease in abundance.
Spain (Basque Country, Cantabria region)	See: Borja, A., F. Aguirrezabalaga, J. Martinez, J.C. Sola, L. Garciaarberas & J.M. Gorostiaga, 2003. Benthic communities, biogeography and resources management. In: Borja, A. & M. Collins, (Ed.). <i>Ocenaography</i>	Borja, A., A.B. Josefson, A. Miles, I. Muxika, F. Olsgard, G. Phillips, J.G. Rodríguez & B. Rygg, 2007. An approach to the intercalibration of benthic ecological status assessment in the North Atlantic ecoregion, according to the European Water

Member State	Description of reference community	Description of good status community
	and Marine Environment of the Basque Country, Elsevier Oceanography Series n. 70: 27-50.	Framework Directive. Marine Pollution Bulletin 55: 42-52.
Spain (Andalusia)	n.a.	n.a.
United Kingdom/Ireland	Reference condition macrobenthic communities are dominated by pollution sensitive taxa (AMBI Ecological Group (EG) I taxa), have low relative abundance of indifferent (EG II) and tolerant (EG III) taxa and negligible relative abundance of opportunist (EG IV) and pollution indicator (EG V) taxa. High numbers of taxa with an even abundance distribution throughout the community is also indicative of reference conditions.	Taxa number and Simpsons evenness are slightly reduced in comparison to values under reference conditions, while variables according to habitat (community abundance as assessed by AMBI) are slightly unbalanced: sensitive taxa (EG I) abundance may range from high sub-dominant to absent; indifferent taxa (EG II) are of low sub-dominant abundance; tolerant taxa (EG III) of dominant abundance; abundance of opportunistic (EG IV) and indicator taxa (EG V) may range from negligible or low to comparable abundance with indifferent taxa (EG II).

PART B- Common type NEA 3/4

B.2 Description of national assessment methods

A benthic assessment approach consists of an indicator algorithm, boundary settings and a reference setting approach. Two benthic assessment approaches need to be intercalibrated in this case. The Netherlands used the BEQI2 method to evaluate the ecological status in type 3/4; whereas Germany selected the m-AMBI method.

B.2.1 Methods and required BQE parameters

The current intercalibration exercise is based on the latest versions of the multi-metric indicator algorithms (Table 26). The BEQI2 consist of the parameters species richness, Shannon wiener and AMBI and were equally weighted in the EQR determination (Van Loon et al., 2015). The m-AMBI takes into account the same parameters, but the EQR is determined based on a factor analysis (Borja et al., 2004; Muxika et al., 2007). The EQR values determined for the samples within the common dataset are re-calculated based on those algorithms. The benthic parameters (species richness, Shannon diversity and AMBI) for the multi-metric or multivariate analyses are derived from the AMBI tool.

The WFD requires the inclusion of certain metrics within the national assessment method for benthic invertebrates, which are summarized for each Member State in Table 37. Both assessment methods contain the required parameters.

Table 26. Overview of the algorithms of the two assessment methods. H': Shannon wiener diversity; S: Number of species; AZTI: Marine Biotic Index.

MULTIMETRIC		
BEQI2 (The Netherlands)	$\text{EQR (ecotope)} = \frac{1}{3} * [\text{S}_{\text{ass}} / \text{S}_{\text{ref}}] + \frac{1}{3} * [\text{H}'_{\text{ass}} / \text{H}'_{\text{ref}}]^1 + \frac{1}{3} * [(6 - \text{AMBI}_{\text{ass}}) / (6 - \text{AMBI}_{\text{ref}})]$	Van Loon et al., 2015
MULTIVARIATE		
M-AMBI (Germany)	Factor analysis: S, AMBI, Shannon diversity index ¹	(Borja et al., 2004 and Muxika et al., 2007) http://ambi.azti.es

¹Shannon diversity: log base 2.

Table 27. Overview of the metrics included in the national assessment methods

Member state	Full BQE method	Taxonomic composition	Abundance	Disturbance sensitive taxa	Diversity	Bio-mass	Taxa indicative of pollution	Combination rule of metrics
Netherlands	Yes	Not strictly – only as groups (5) of different sensitivity	As relative abundance of different sensitivity groups and proportional abundance in Shannon Wiener index	5 sensitivity classes (AMBI)	Yes, number of species and Shannon Wiener index	No	Group of opportunistic species	Average of 3 univariately normalized indicator EQR scores
Germany	Yes	Not strictly – only as groups (5) of different sensitivity	As relative abundance of different sensitivity groups and proportional abundance in Shannon Wiener index	5 sensitivity classes (AMBI)	Yes, number of species and Shannon Wiener index	No	Group of opportunistic species	Factorial analyses, calculating vectorial distances to reference conditions

B.2.2 Sampling and data processing

The benthic sampling procedure for the WFD Monitoring within the Netherlands and Germany for type NEA 3/4 is slightly different, especially regarding the sampling design.

The benthic sampling in the intertidal habitats in Germany are done by cores (different sizes possible) at certain locations. At each location 10 replicate samples were taken. In the Netherlands transect sampling is applied. Each transect is composed of 10 (Balgzand) or 20 (Piet Scheveplaat) stations. At each station, 2 (Piet Scheveplaat from 2009 onwards), 3 (Piet Scheveplaat before 2009) or 5 (Balgzand) replicate small core samples have been sampled and combined. The sample area of the cores and the number of cores combined per station show some changes during the years, which is document in several monitoring reports of NIOZ and Koeman and Bijkerk, the external benthos laboratories.

The processing of the samples is similar, with identification and counting of the individuals to species level. The taxonomy in both countries is standardized regarding WORMS. The level of the species determination and truncation rules are country specific and applied on the entire data set.

B.2.3 National reference conditions

The determination of the reference conditions is a complicated subject (Van Hoey et al., 2010; Birk et al., 2013). The ecological status in the WFD has to be measured as a deviation from a reference condition. These reference conditions need to correspond to largely undisturbed (=‘near-pristine’) conditions (no or minor impact from human activities). Indeed, the lack of appropriate reference sites or robust historical datasets is one of the major problems addressed in the intercalibration exercises and in setting the good ecological status boundaries (Borja et al., 2007; 2009). Scientists are faced with virtual lack of undisturbed sites along the European coasts and estuaries, and historical data are not easily accessible (Borja et al., 2004). Reference settings will need to be based on clear stressor-response relationships, a knowledge of the ‘naturalness’ of the system; and expert judgment may also have a role to play (Van Hoey et al., 2010). As summarized in Table 4, both countries used the best available information (e.g. areas with least disturbed conditions) and their expert judgment to delineate appropriate reference values for their metrics. For most methods, the principle is to use highest indicator value which is not an outlier. For this reason, high percentile values (99 to 95p) (for AMBI low percentile values; 1 to 5 p) are mainly used (Van Loon et al., 2015).

The reference values used to calculate the EQR values for each sample within a habitat (also referred to as ecotopes in the BEQI2 MMI) in the common dataset are listed in 29. Those values were applied per benthic assessment approach on the common dataset.

Table 28. Overview of the methodologies used to derive the reference conditions for the national assessment methods included in the Ic exercise

Member State	Type and period of reference conditions	Number of reference sites	Location of reference sites	Reference criteria used for selection of reference sites
Germany	Expert knowledge, Historical data, Least Disturbed Conditions; reference time: 1959 up to now. Habitat-specific. The highest values from the reference data sets were selected as reference values for AMBI, Diversity and richness. As reference value for the bad conditions 0 is used for Richness and Diversity, 6 for AMBI.	Not true reference sites, but least disturbed sites, 6 sites for subtidal, 9 sites for littoral stations (two in the common intercalibration dataset).	different sites Wadden Sea of Lower Saxony	The communities at the sites had to correspond with description of the reference community description referring to a certain habitat. This approach is based on the hypothesis that most undisturbed areas are still found in small patches and will be represented by the best sites in the data set of the corresponding habitat.
Netherlands	(a) Historical data for 1991-2006; (b) Estimation of reference values: AMBI(ref): the 1 percentile value; S(ref) and H'(ref): 99 percentile of S and H' for dataset 1992-2006 (15 years). The principle is to use highest indicator value which is not an outlier. (c) theoretical bad values: S(bad) = 0; H'(bad) = 0; AMBI(bad) = 6. (c)	Not true reference sites, but least disturbed sites can be selected if necessary, primarily in the intertidal area Piet Scheveplaat, where the fishery is minimal.	The Piet Scheveplaat in the Wadden Sea is a reference site for intertidal habitat.	Not applicable because marine waters in The Netherlands are always subject to at least some level of anthropogenic impact. However, least disturbed samples from distinct sampling locations can be selected based on expert judgment using information on pressures at the sampling locations.

¹Changed compared to the WISER input, based on Van Hoey et al., 2014 report.

Table 29. Overview of the reference values for benthic characteristics used in the intercalibration exercise.

Intertidal	Habitat	Sampled surface (m ²)	Sampling device	Species richness	Shannon (H' log2)	AMBI
Germany	Sand	0.2	plastic tubes	20	3.24	0.02
Germany	Muddy Sand	0.2	plastic tubes	21	3.11	1.61
Germany	mud	0.04	plastic tubes	20	2.9	2
Netherlands	muddy sand	0.1m ²	Manual cores (0,008m ²)	29	3.6	0.54

Subtidal	Habitat	Sampled surface (m ²)	Sampling device	Species richness	Shannon (H' log2)	AMBI
Germany	Subtidal high dynamic (sand)	0.9	Van Veen	36	3.61	0.36
Germany	Subtidal low dynamic (muddy sand to sand)	0.9	Van Veen	30	3.77	0.05
Netherlands	Subtidal	0.12 (2 boxcores of 0.06 m ² pooled)	Boxcorer	23	3.5	0.54

Two questions arose from analyzing the table 29:

- 1) The species richness between the muddy intertidal and other intertidal habitats in Germany, is not that different, despite the difference in sampling surface (0.04 compared to 0.2 respectively).

This estimation of the reference values is appropriate for this moment, because no differences in the number of species could be detected if the sampled area was enlarged. Therefore, the reference values for the intertidal mud for an area of 0.181m² can be considered as the same as for an area of 0.04m².

- 2) There is a difference between the reference values for the intertidal habitats of Germany and the intertidal habitat of the Netherlands. The values in the Netherlands were higher than in Germany, despite the lower sampling surface.

This difference in reference values, especially for species richness can be attributed to the following facts:

- The sampling design, which is point sampling (10 samples) in Germany and transect sampling (3*20 samples) per location in the Netherlands.
- The species richness in the Netherlands is also estimated based on pooling and aggregating samples over a wider spatial range (more than one location). This leads to relatively higher reference values for S (see Van Loon et al. 2015). In Germany it is location specific.
- And also some difference in the taxonomical truncation rules between the countries.

There is a big difference in total sampled area per country in the common dataset, which result in a different amount of species encountered in the data. For the intertidal muddy sand habitat, Germany founds 85 species (19 rare species), whereas the Netherlands 143 (40 rare species). This differences in species pool for both datasets, resulted from difference in total sampled area and sampling strategy, reasons for difference in reference values.

B.2.4 National boundary setting

The boundary setting procedure for both countries is summarized in 31. The boundary values used in the intercalibration for Germany and the Netherlands for type NEA3/4 were summarized in 30.

Table 30. The boundary values for the different assessment approaches as used in the Ic exercise

	High/Good	Good/Moderate	Moderate/Poor	Poor/Bad
Germany	0.85	0.70	0.40	0.20
Netherlands	0.80	0.60	0.40	0.20

Table 31. Explanations for national boundary setting of the national methods included in the Ic exercise

Member State	Type of boundary setting	Specific approach for H/G boundary	Specific approach for G/M boundary	BSP: method tested against pressure
Germany	Boundaries taken over from the intercalibration exercise (Borja et al., 2007 ¹). Calibrated against pre-classified sampling sites. The boundary setting procedure is in line with the WFD's normative definitions.			The boundaries were additionally adjusted by the assessment of expert judgment (Heyer 2007). The m-AMBI relates to pressures of sediment enrichment, eutrophication and hazardous substances (Muxika et al. 2007).
Netherlands		The Good/Moderate boundary of 0.60 is primarily derived from the initial G/M boundary for sheltered coastal waters (Van Hoey et al., 2015), which was estimated using expert judgment and set at 0.60 (see. Van Loon et al. 2015, paragraph 2.7. for more information).		

B.2.4 Results of WFD compliance checking

Table 32. WFD Compliance checking criteria

Compliance criteria	Compliance checking conclusions
1. Ecological status is classified by one of five classes (high, good, moderate, poor and bad).	Yes, for both benthic assessment approaches
2. High, good and moderate ecological status are set in line with the WFD's normative definitions (Boundary setting procedure)	Yes, for both benthic assessment approaches
3. All relevant parameters indicative of the biological quality element are covered (see Table 1 in the IC Guidance). A combination rule to combine parameter assessment into BQE assessment has to be defined. If parameters are missing, Member States need to demonstrate that the method is sufficiently indicative of the status of the QE as a whole.	The two Member States included the relevant parameters (see Table 3), A combination rule to combine parameter assessment is defined by both benthic assessment approaches.
4. Assessment is adapted to intercalibration common types that are defined in line with the typological requirements of the WFD Annex II and approved by WG ECOSTAT	Yes, for both Member States
5. The water body is assessed against type-specific near-natural reference conditions	No (see Table 4). Alternative benchmark conditions (based on a "least disturbed condition" criteria) had to be defined due to the absence of near-natural reference conditions in the intercalibrated type.
6. Assessment results are expressed as EQRs	Yes, for both benthic assessment approaches
7. Sampling procedure allows for representative information about water body quality/ecological status in space and time	In most cases, the monitoring is considered as representative by the Member State itself. This aspect is not confirmed by specific, standardized analyses to test their representativeness. Sampling procedures are outlined in general, but not linked with the running WFD monitoring programs.
8. All data relevant for assessing the biological parameters specified in the WFD's normative definitions are covered by the sampling procedure	Yes, for both benthic assessment approaches. The sampling procedure defined by each Member State allows the collection of species-abundance data, which is necessary to calculate all metrics of the different benthic assessment approaches.
9. Selected taxonomic level achieves adequate confidence and precision in classification	Yes, for both benthic assessment approaches, with some difference in taxonomic detail per Member State, but sufficient comparability. The taxonomic discrimination rules are country species and applied to each member states dataset.

There can be concluded that all compliance criteria were met for both benthic assessment approaches.

B.3 Feasibility checking

B.3.1 Typology

In the NE Atlantic, seven basic intercalibration types have been agreed upon. In this report the type NEA3/4 is taken into account (see outline of characteristics in 33).

Table 33. NEA GIG Intercalibration Type NEA 3/4

New Type ID	Name	Salinity [PSU]	Tidal range (m)	Depth (m)	Current velocity (knots) [m/s]	Exposure	Mixing	Residence time
CW – NEA3/4	Polyhaline, exposed or moderately exposed (Wadden Sea type)	Polyhaline (18 - 30)	Mesotidal (1 - 5)	Shallow (< 30)	Medium (0,51- 1,54m/s)	Exposed or moderately exposed	Fully mixed	Days

This type is only discriminated in the Netherlands and Germany.

B.3.2 Pressures addressed

The BEQI2 and m-AMBI assessment approach are well tested against a pressure gradient. This pressure-response relation of both approaches are published in literature (Borja et al., 2009; Van Loon et al., 2015) and intercalibration report (NEA-GIG coastal waters, Van Hoey et al., 2015). Both methods are sensitive to various types of pressures, as eutrophication, oxygen depletion (see Dutch example), physical disturbance (see German sand extraction example) and increased suspended matter (see Dutch example).

Dutch example (Van Loon et al., 2015):

The sensitivity of the BEQI2 for human and natural induced stressors was explored by regression analysis of regional BEQI2 and time-series of measurements of dissolved oxygen in the Westerschelde mesohaline-intertidal ecotope and of the suspended matter concentration in the Dollard mesohaline-intertidal ecotope (Figure 10). The BEQI2 shows a positive, significant correlation with oxygen concentration, meaning that an increase in oxygen concentration leads to a higher BEQI2 EQR. Beside it, the BEQI2 shows a negative, significant correlation with suspended matter, meaning that a higher SPM concentration leads to a lower BEQI2 EQR.

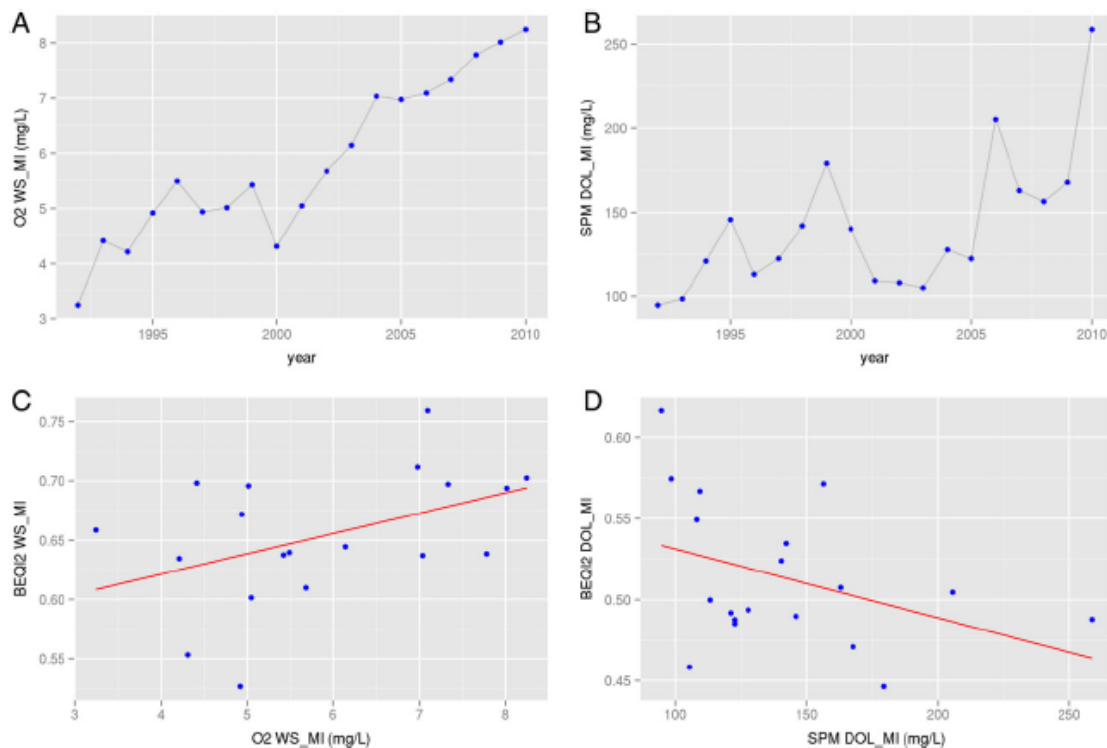


Figure 10. A-B Time trends of the state of the parameters oxygen and suspended matter. C-D state impact correlations for oxygen concentration and suspended matter with BEQI EQRs. waterbody ecotopes Westerschelde mesohaline-intertidal (WS_MI) and Dollard mesohaline-intertidal (DOI_MI), respectively.

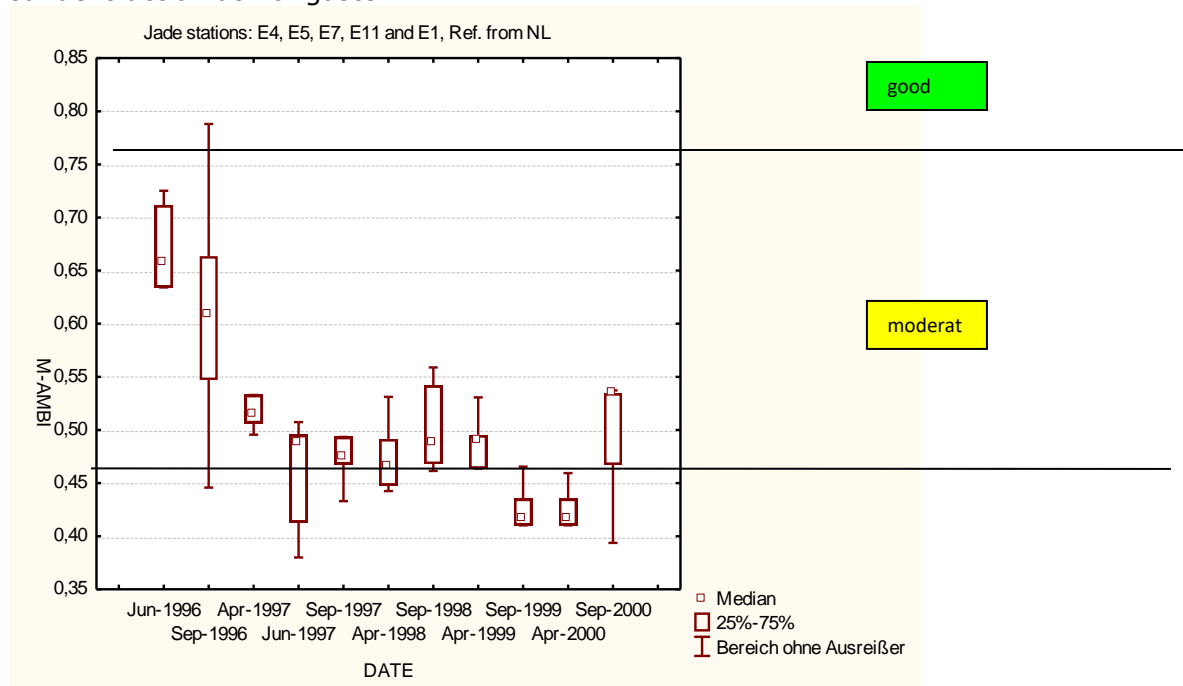
German example:

In the Dangaster Außentief (German Wadden Sea) in July 1996 huge sand extraction (1,2 million m³ sand) took place. Before (June 1996) and after sand extraction the macrozoobenthos was investigated at several stations (Fischer et al. 2004) twice or thrice a year (April, June and September) until June 2000. With the data of five (E4, E5, E7, E11 and E17) out of these stations the M-AMBI values were calculated (Figure 2). The chosen stations laid to the south and in a distance between 50 m to 300 m from of the sand extraction area.

The M-AMBIs were calculated with the NL reference values given by (van Hoey et al. 2007) (AMBI 0.6, diversity 2.35 and richness 24). It is a static and correlative comparison, as no specific pressure linked variable (as organic matter content, sediment re-suspension or suspended matter), is available.

The ecological status decreased from a 'good' ('II') to a 'moderate' ('III') (Figure 11). In September 2000 the M-AMBI increased again.

Figure 11. M-AMBi values at each sampling data in the BACI design monitoring for sandextraction at Dangaster



B.3.3 Assessment concept

Do all national methods follow a similar assessment concept?

The two benthic assessment approaches for type NEA3/4 are very similar. They consist both of the same metrics (parameters) and differ only in their EQR calculation algorithm. The BEQI2 has a fixed formula and a priori pooling of the samples, whereas the m-AMBI is based on a factor analysis.

The main difference in assessment concept between the Netherlands and Germany is situated in how the raw data is pooled for determining the EQR values per habitat type. The BEQI2 assessment approach executed a randomisation procedure, which pool the small core samples obtained within a single habitat-year at random to 0.1m² (sample pool size) and repeat this 10 times to calculate per habitat the average BEQI2 score. This lead to an EQR value per year for each habitat within a waterbody. The Germany assessment approach pool the core samples per station a priori to the calculation of the EQR values for that station by the m-AMBI. The number of samples can vary between station and habitat type. If more stations are available per habitat type/waterbody, those EQR values need to be 'averaged' to come to an EQR value per habitat within a waterbody. For both assessment methods, the reference values were in accordance with the pooling principle and obtained sample pool sizes.

Due to this situation, we have different levels (habitat versus location) and sampling areas between both assessment approaches to calculate the EQR values. Therefore, this difference in concept is harmonized for intercalibration purpose. It is clear that it is not appropriate to calculate the EQR values on sample level (core or grab), due to the fact that both countries

do it on a higher level (standardised sample pool surface). Therefore, we decided to work with a 'common' fixed sample size of 0.81m² for the intercalibration, which is the standard for the German assessment approach, but not in correspondence with the Dutch assessment way. For harmonization purpose, the data of the Netherlands is split in separate location assessments instead of an entire habitat assessment. This is feasible and acceptable and the relation between both approaches should be more or less the same, regardless the level of pooling.

	BEQI2	m-AMBI
Dutch dataset	A priori pooling of the subsamples to corresponding sample pool size of the Dutch reference values.	A priori pooling of the subsamples to corresponding sample size of the German reference values. By this the German reference values can be used for the assessment of the Dutch data.
German dataset	BEQI2 calculated on the a priori pooled German subsamples. The BEQI reference values can be used, despite their is a slight difference in total sample surface.	A priori pooled subsamples (10) to corresponding surface per location, as the German assessment method is.

In this case, we have compared 143 (German dataset) and 180 (Dutch dataset) sample assessments, which should give enough values to test the comparability criteria (Table 34). This create an unequal balance in data between both countries, but this has no influence on the comparison results. If the data of the years 2000 and 2001 in the Dutch dataset were not considered, the same results were obtained regarding the boundary adjustment (from 0.6 to 0.611).

Is the Intercalibration feasible in terms of **assessment concepts**?

Yes, despite some small difference in the way the EQR calculation occur for both benthic indicator approaches.

B.4 Collection of intercalibration dataset and benchmarking

B.4.1 Dataset description

At the start of the project, we had an expert meeting where we discussed the data availability and appropriateness. First, we decided to use autumn data only, to exclude seasonal variation. Second, we decided to focus on intertidal habitats, because most appropriate intercalibration data could be derived for it. This in the light of selecting benchmark samples. For the subtidal habitats, no appropriate pressure data was available, neither sites could be selected as benchmark sites by expert judgment. For the intertidal habitats, sites for both countries with similar level of eutrophication and negligible fishery pressure could be selected. Finally, the benthic data from the muddy sand habitat in the intertidal was selected because the Dutch monitoring focused on this habitat type and also a lot of German sites belong to this habitat type (Table 34). The similarity in the samples of the Netherlands and German for the intertidal habitats is investigated in "Multivariate analyses" section and is very good.

Therefore, due the availability of benchmark sites for the intertidal muddy sand in both countries and a large amount of data, the comparability of the assessment approaches is tested on this data set.

Table 34. Overview of the available data and its metadata information.

Dataset	Station	program	#assessments	Time period	Grouping of subsamples	Total surface	Waterbody type	Habitat/ecotoop	Benchmark
GE1	AuWe_MZB_3	NLWKN	8	2007-2014	10*0,0181	0,181	N4_4900_01	intertidal sand	no
GE1	Nney_MZB_1	NLWKN	8	2007-2014	10*0,0181	0,181	N4_3100_01	intertidal sand	no
GE1	Nney_MZB_2	NLWKN	8	2007-2014	10*0,0181	0,181	N4_3100_01	intertidal sand	yes
GE1	Nney_MZB_3	NLWKN	8	2007-2014	10*0,0181	0,181	N4_3100_01	intertidal muddy sand	no
GE1	Nney_MZB_5	NLWKN	8	2007-2014	10*0,0181	0,181	N4_3100_01	intertidal mud	no
GE1	Nney_MZB_6	NLWKN	8	2007-2014	10*0,0181	0,181	N4_3100_01	intertidal mud	no
GE1	Nney_MZB_7	NLWKN	8	2007-2014	10*0,0181	0,181	N4_3100_01	intertidal mud	no
GE1	Nney_MZB_8	NLWKN	8	2007-2014	10*0,0181	0,181	N4_3100_01	intertidal muddy sand	yes
GE1	WuKu_MZB_6	NLWKN	8	2007-2014	10*0,0181	0,181	N4_4900_02	intertidal muddy sand	no
GE1	WuKu_MZB_10	NLWKN	1	2007	10*0,0181	0,181	N4_5900_01	intertidal muddy sand	no
GE2	HH T1	HH	14	2000-2013	75*0,00166	0,1245		intertidal muddy sand	no
GE2	HH T2	HH	14	2000-2013	75*0,00166	0,1245		intertidal muddy sand	no
GE2	HH T3	HH	14	2000-2013	75*0,00166	0,1245		intertidal muddy sand	no
GE2	HH T4	HH	14	2000-2013	75*0,00166	0,1245		intertidal muddy sand	no
GE2	HH T5	HH	14	2000-2013	75*0,00166	0,1245		intertidal muddy sand	no
NL1	Balgzand-Raai J_A	Balgzand	15	2000-2014	substaal 1-12 (12*0,0157)	0,1884	Waddensea	intertidal mud-muddy	no
NL1	Balgzand-Raai J_B	Balgzand	15	2000-2014	substaal 13-24	0,1884	Waddensea	intertidal muddy sand	no
NL1	Balgzand-Raai B_A	Balgzand	15	2000-2014	substaal 1-12	0,1884	Waddensea	intertidal muddy sand	no
NL1	Balgzand-Raai B_B	Balgzand	15	2000-2014	substaal 13-24	0,1884	Waddensea	intertidal muddy sand	no
NL1	Balgzand-Raai C_A	Balgzand	15	2000-2014	substaal 1-12	0,1884	Waddensea	intertidal muddy sand	no
NL1	Balgzand-Raai C_B	Balgzand	15	2000-2014	substaal 13-24	0,1884	Waddensea	intertidal muddy sand	no
NL2	Piet Scheveplaat - Raai 600_A	Piet Schev	15	2000-2014	substaal 1-10 (10*0,0157)	0,157	Waddensea	intertidal muddy sand	yes
NL2	Piet Scheveplaat - Raai 600_B	Piet Schev	15	2000-2014	substaal 11-20	0,157	Waddensea	intertidal muddy sand	yes
NL2	Piet Scheveplaat - Raai 601_A	Piet Schev	15	2000-2014	substaal 1-10	0,157	Waddensea	intertidal muddy sand	yes
NL2	Piet Scheveplaat - Raai 601_B	Piet Schev	15	2000-2014	substaal 11-20	0,157	Waddensea	intertidal muddy sand	yes
NL2	Piet Scheveplaat - Raai 602_A	Piet Schev	15	2000-2014	substaal 1-10	0,157	Waddensea	intertidal muddy sand	yes
NL2	Piet Scheveplaat - Raai 602_B	Piet Schev	15	2000-2014	substaal 11-20	0,157	Waddensea	intertidal muddy sand	yes

B.4.2 Data acceptance criteria

The Netherlands and Germany have delivered data for the intercalibration exercise.

To explore the common intercalibration dataset for benthic macro-invertebrates, we performed some standard multivariate analyses. This to evaluate the following aspects:

- to check for outliers (samples very different from the rest and showing a problem)
- If there were regional or sub-regional differences between the samples and habitats
- If different benthic communities could be detected, which can be related to different physical habitats (sedimentology).
- If there is any pattern in the data that justifies the delineation of sub-types for benchmarking, even the fact that we already select common types.

B.4.3 General multivariate analyses

For the purpose of the multivariate analyses, the common dataset is fourth root transformed to reduce the effect of very abundant species on the overall pattern. Beside this, the rare species (with less than 3 individuals) were excluded from these analyses to reduce the effect of rare species on the overall pattern. The similarity between samples is determined by the Bray-Curtis similarity. The sample groups were determined based on a cluster analyses, with cut-off level at certain similarity level (31). Multidimensional scaling (MDS) is used to visualize the cluster groups (**Error! Reference source not found.12**). The sample groups discriminated from the cluster analyses were compared with the habitat type considered by the experts (**Error! Reference source not found.13**). The analyses were executed in PRIMER6.

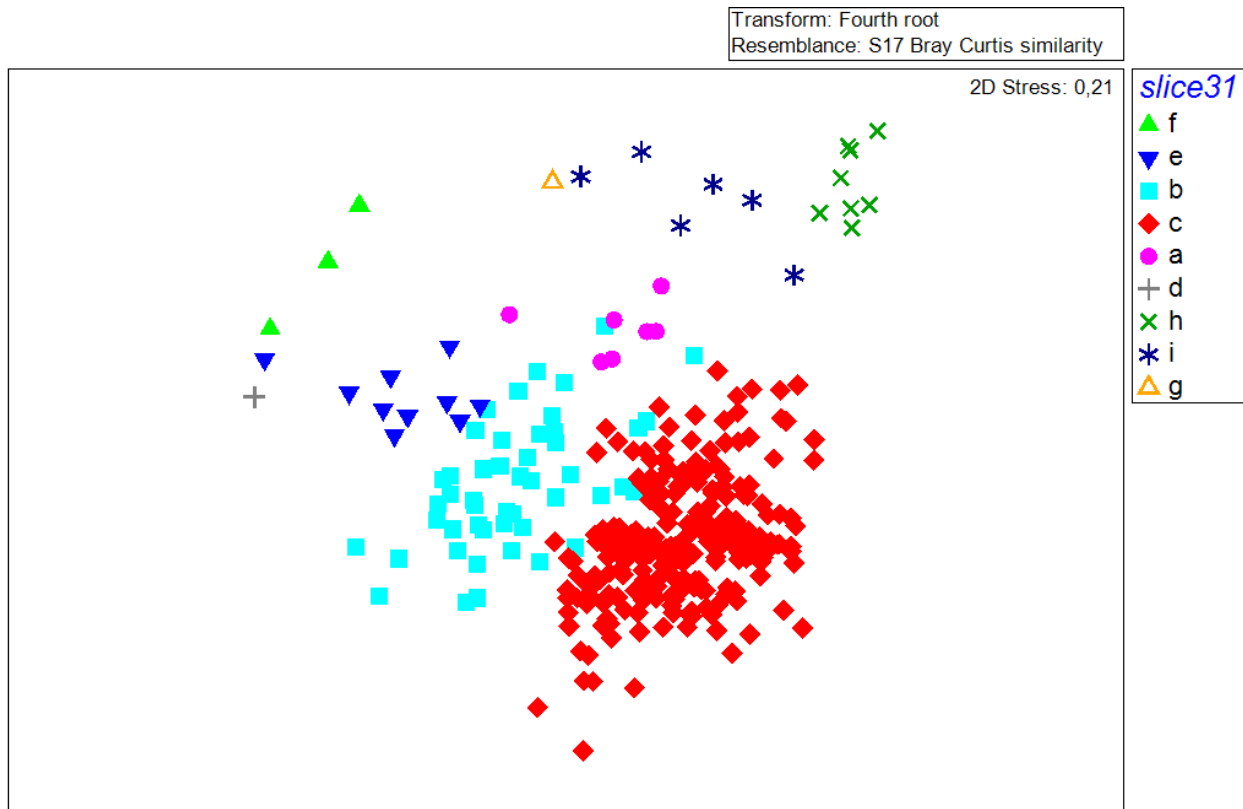


Figure 12. MDS of the cluster groups 9 slice 31 Bray Curtis similarity), which result in 9n groups are coded alphabetically (a-i)

Some explanation on the cluster groups:

- No outlier samples present in the common dataset (no very different sample from the rest).
- The subtidal habitats clearly separated from the intertidal habitats, both in the cluster groups (a, e, f) as by the habitat groups (subtidal mud and fine sand). Those were not further considered for the intercalibration, because the focus is on the intertidal habitats.
- The intertidal mud habitat (Germany) clearly clustered separately from the others, in cluster i, g and h (location dependent). This means, that this habitat type could be a separated sub-type for the Wadden sea. Due to the absence of Dutch data for this type, this is not further considered.
- The samples, considered located in an intertidal sand habitat, could not be discriminated from the intertidal muddy sand habitat in the cluster analyses (belong to cluster b and c). This can mean that the location considered as intertidal sand, should not be a separate subtype for this intercalibration.
- The majority of the samples in the common dataset were from the intertidal muddy sand habitat and clustered together in two main clusters (b and c).
 - o Cluster b contains the samples of Balgzand 'raai' ZDJ en AuWe-MZB3 and are slightly different from the other intertidal muddy sand locations.
 - o Cluster c contains the majority of the samples and are reflecting the species composition of an intertidal muddy sand habitat in the Wadden Sea area. This cluster clearly groups the samples of this habitat type of both countries.

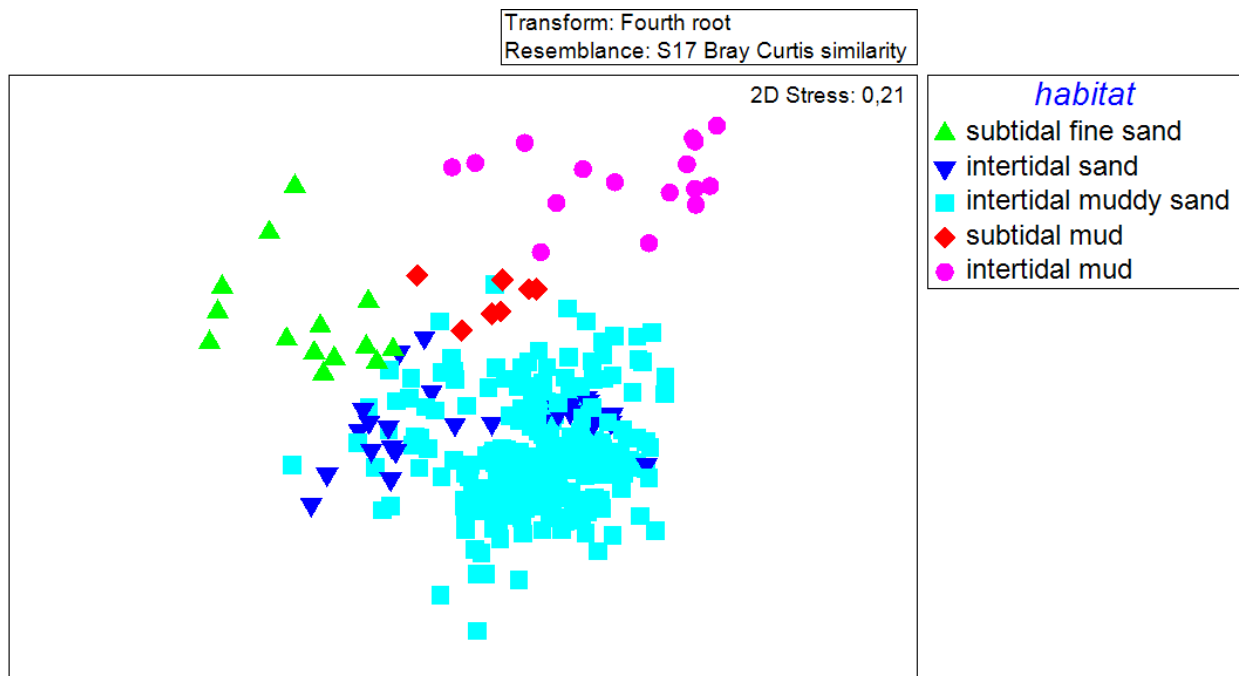


Figure 13. MDS with indication of the habitat types

We can have concluded, based on the species composition, that the benthic fauna in the Wadden Sea area is similar between Germany and the Netherlands. There is no geographical difference in species composition and main characteristics within the intertidal muddy sand habitat type. This analyses also shows that it is relevant to consider the habitats separately, as sub-types if necessary. This means that it is preferred that the reference conditions are habitat specific, as Germany does. Only, the difference in community characteristics between intertidal sand and muddy sand is not obvious, due to the position of the intertidal sand samples in the MDS.

For the intercalibration exercise, we can clearly use the samples of the intertidal muddy sand habitat of both countries to test the comparability between both benthic assessment approaches.

B.4.4 Common benchmark

Both countries have select a benchmark site that is subjected to a similar level of eutrophication but consider the lowest influence of fishery. Details on the level of eutrophication and fishery for the German locations are given in the table in annex 4. Both pressures are the main driver for changes in the benthic system within the Wadden Sea area.

For the Netherlands this is the Piet Scheveplaat for the intertidal habitat and for Germany that is the Nney_MZ8 site for the intertidal muddy sand habitat.

B.4.5 Benchmark standardization

The principal aim of benchmarking in intercalibration is to identify and remove differences among national assessment methods that are not caused by anthropogenic pressure but

rather by systematic discrepancies (due to different methodology, biogeography, typology etc.; see remarks in section on reference settings) (Annex V, IC Guidance).

Benchmark standardization will correct for differences in median EQR values between the Member States' benchmark sites obtained by certain assessment approaches. Those median values will be corrected by the benchmark standardization procedure; this correction will be more obvious for cases where the medians are significantly different.

We tested whether benchmark standardization was necessary. Student's *sT* was used to compare the benchmark sites values for the two national methods.

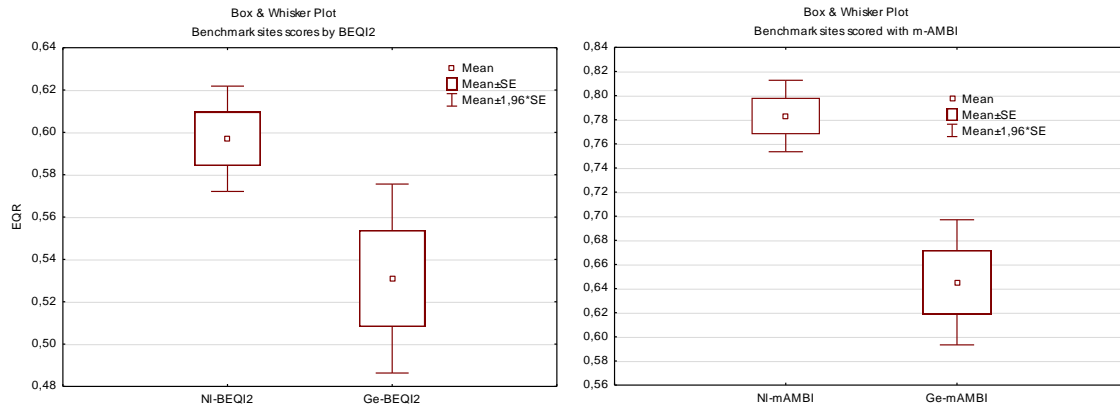


Figure 14. Box-whisker plot of the assessment of the Dutch and German benchmark sites with each benthic assessment approach.

Figure 14. Box-whisker plot of the assessment of the Dutch and German benchmark sites with each benthic assessment approach.

The benchmark sites of both countries were not significantly different from each other for the BEQI2 ($p = 0,155$) (left box whisker plot) (14), despite the difference in the box plot. The benchmark sites of both countries were significant different with the m-AMBI approach ($p = 0.0135$) (right box-whisker plot) (Figure 14). This indicated that benchmark standardization is necessary.

The correlation between the average value of all national EQRs per survey in the full dataset was not significantly correlated with its standard deviation, therefore national EQRs does not converge towards the bad end of the quality gradient, and therefore, subtraction was used for the standardization.

B.5 Comparison of methods and boundaries

B.5.1 Intercalibration option and common metrics

Option 3a. Intercalibration can be performed based on commonly assessed sites and whether the ecological quality gradient is sufficiently covered. Only two methods are involved in the intercalibration, which involve that there is a direct comparison (pseudo-metric=other method).

B.5.2 Results of the regression comparison

The regression comparison shows that both methods correlated very well ($R^2 = 0.9103$).

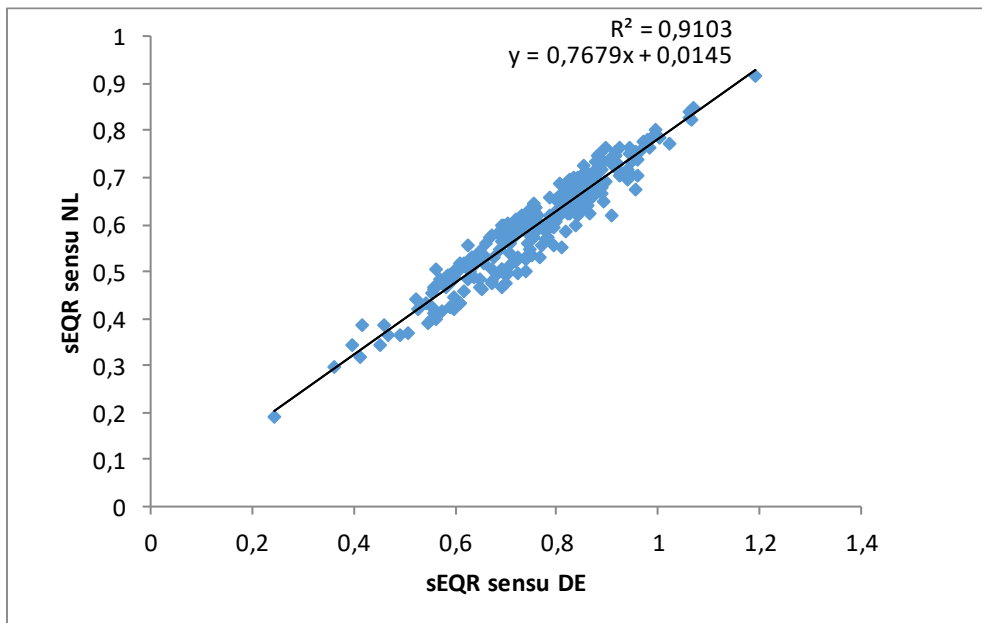


Figure 15. Scatter plot of EQR values of Germany and Netherlands, with linear regression line.

Figure 15. Scatter plot of EQR values of Germany and Netherlands, with linear regression line.

B.5.3 Comparability criteria

The boundary bias criteria are above 0.25 for the H/G boundary of the m-AMBI and G/M boundary of the BEQI2. The H/G boundary of the m-AMBI is slightly above the criteria, but a change is not suggested by the excel sheet. The G/M boundary of the BEQI2 need to be slightly increased to meet the boundary bias criteria by 0.11 to 0.611.

Table 35. Boundary bias values for the High/Good and Good/ Moderate boundaries for the German and Dutch benthic assessment methods.

Boundary	A Germany	A on scale of B	B Netherlands	B on scale of A	A average bias	B average bias	A excess as classes	A harmonised boundary	B excess as classes	B harmonised boundary
MP	0,400	0,415	0,400	0,400						
GM	0,700	0,715	0,600	0,582	0,194	-0,306		no change	0	0,611
HG	0,850	0,865	0,800	0,764	0,252	-0,200	0,002	0,850		no change

The average absolute class difference for the five classes between both methods is 0.35 (<0.5). If the poor and bad classes are not taken into account, the average absolute class difference is 0.39 (<0.5).

These results seem to be logically, because the boundaries for Germany are higher than for the Netherlands, but for the reference values it is the reverse. This lead to the fact that both benthic assessment approaches are comparable.

B.6 Final results to be included in the EC

Table with EQRs

A boundary adjustment for the G/M boundary by the Netherlands is needed. They accepted to increase the boundary to 0,61. The final boundaries for the benthic assessment approaches (BEQI2 and m-AMBI) for the Wadden Sea in the North-east Atlantic are given in the table 36.

Table 36. Boundary values of the different benthic assessment approaches after intercalibration. The boundaries in red are those changed after boundaries harmonization

Country	Benthic assessment approach	Ecological quality ratios			
		High-good boundary	Good-moderate boundary	Moderate-poor boundary	Poor-bad boundary
Germany	m-AMBI	0.85	0.70	0.4	0.2
Netherlands	BEQI2	0.80	0.61	0.4	0.2

Correspondence common types versus national types

The common type (NEA3/4) is recognized as type in every Member State and is related to the national types.

Gaps of the current intercalibration

Not all habitat types within the Wadden Sea could be considered, due to the absence of a comparable dataset for those habitats between both countries, especially in the light of discriminating appropriate benchmark sites for those habitats.

B.7 Ecological characteristics

B.7.1 Description of reference or alternative benchmark communities

The description of the benthic community characteristics at reference or alternative benchmark is summarized in 37. This information is generated from the WISER database. Only for France, Norway and Spain (Andalusia) this information is not available.

B.7.2 Description of good status communities

The description of the benthic community characteristics at good status is summarized in Table 2537. This information is generated from the WISER database.

Table 37. Overview of the description by the member states of the macroinvertebrate reference community and good status community

Member State	Description of reference community	Description of good status community
Germany	Benthic communities, species numbers, diversity typically for the habitat (sediment, salinity, exposure)- low number of opportunistic species.	High portion of sensitive taxa, complex communities, low number of opportunists, high species number and high diversity assemblages.
Netherlands		

PART C- Common type NEA 7

C.2 Description of national assessment methods

Table 38. Overview of the national assessment methods.

Member State	Method	Included in this IC exercise?
Norway	Norwegian Quality Index (NQIvI)	Yes
United Kingdom	Infaunal Quality Index (IQIvIV)	Yes

NQIvI (Rygg 2006): The NQIvI is a multimetric index composed of the following metrics:

- (i) AZTI Marine Biotic Index (AMBI) (sensitivity component)
- (ii) SN (number of taxa (S) and abundance (N)) (diversity factor)
- (iii) a correction factor for down-weighting artificially high index values of small samples (few individuals (N/N+5)).

The index is a weighted algorithm (50 % AMBI and 50 % species/abundance) formulated as follows:

$$NQIvI = \left(0.5 * \left(1 - \frac{AMBI}{7}\right) + \left(0.5 \frac{SN}{2.7} * \frac{N}{N+5}\right)\right)$$

The class boundaries are: High/Good = 0.72, Good/Moderate = 0.63.

IQIvIV (Phillips et al. 2014, UKTAG 2014): The IQIvIV is a multimetric index composed of three individual metrics:

- (i) AZTI Marine Biotic Index (AMBI) (sensitivity component)
- (ii) Simpson's Evenness (1-') (diversity factor)
- (iii) number of taxa (S).

Infaunal Quality Index (IQIvIV): The individual metrics have been weighted and combined within the IQIvIV in order to best describe the changes in the benthic invertebrate community in response to anthropogenic pressures. The IQIvIV is formulated as follows:

$$IQI_{v,IV} = \left(\left(0.38 \times \left(\frac{1 - (AMBI/7)}{1 - (AMBI_{Ref}/7)} \right) \right) + \left(0.08 \times \left(\frac{1 - \lambda'}{1 - \lambda'_{Ref}} \right) \right) + \left(0.54 \times \left(\frac{S}{S_{Ref}} \right)^{0.1} \right) - 0.4 \right) / 0.6$$

The four class boundaries are: High/Good = 0.75, Good/Moderate = 0.64, Moderate/Poor = 0.44, Poor/Bad = 0.24.

To calculate the IQivIV the following information is required:

- (i) Abundance of benthic invertebrates (identified to lowest taxonomic level)
- (ii) Characterisation of the habitat sampled (salinity and substratum)
- (iii) Sampling methodology (e.g. sample method area and gear used)
- (iv) Processing methodology (e.g. sieve mesh).

Reference condition metrics are specific for the habitat sampled and sample method used.

C.2.1 Methods and required BQE parameters

Both National methods include the aspects of the benthic invertebrate community that must be included in the ecological status assessment of a water body as defined in Annex V (1.2) of the WFD.

Table 39. Overview of the metrics included in the national assessment methods.

Member State	Full BQE met	Composition	Abundance	Disturbance sensitive taxa	Diversity	Taxa indicative of pollution	Combination rule of metrics
Norway	Yes	Yes – expressed as groups (5) of different sensitivity	Yes – species abundance as correction factor (Ntot/Ntot+5) and relative abundance of different sensitivity groups	Yes – 5 sensitivity classes (AMBI)	Yes, number of species	Yes - Specific opportunistic species	Weighted algorithm. See National description.
United Kingdom	Yes	Yes – expressed as groups (5) of different sensitivity	Yes – expressed as relative abundance of different sensitivity groups and proportional abundance in Simpson index	Yes – 5 sensitivity classes (AMBI)	Yes, number of taxa and Simpson index	Yes - Specific opportunistic species	Weighted algorithm. See National description.

C.2.2 Sampling and data processing

Table 40. Overview of the sampling and data processing of the national assessment methods.

	Norway	United Kingdom
Sampling/survey device	0.1m ² grab, processed using a 1mm sieve	0.1m ² grab, processed using a 1mm sieve
How many sampling/survey occasions (in time) are required to allow for ecological quality classification of survey site or area?	One	One
Sampling/survey months	Recommended sampling period is spring, but classification is possible using data collected throughout the year. July and August should be avoided if possible to avoid large numbers of juveniles.	Recommended sampling period is February to June, inclusive but classification is possible using data collected throughout the year as long as the potential impact of seasonal bias on the classification is considered.
Which method is used to select the sampling /survey site or area?	Sites must be representative of the water body and are selected by expert judgement	Single samples taken from stations spread across suitable habitats within a water body.
How many spatial replicates per sampling/ survey occasion are required to allow for ecological quality classification of sampling/ survey site or area?	Minimum 3 grab replicates per site. Number of sites within each water body vary.	Number of samples required is dependent on the level of inherent variability in the biological community being sampled and associated environmental conditions (UKTAG 2014). Number of sites can vary between water bodies.
Total sampled area or volume, or total surveyed area, or total sampling duration on which ecological quality classification of sampling/survey site or area is based	Minimum sampling area of 0.3 m ² per site. Number of sites within each water body varies.	Water body, single sampling occasion
Short description of field sampling/survey procedure and processing (sub-sampling)	Sampling follows NS-ISO 16665:2013 (2013). Water quality - Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna.	Sampling follows BS-ISO 16665:2013 (2013). Water quality - Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna

C.2.3 National reference conditions

Table 41. Overview of the methodologies used to derive reference conditions for the national assessment methods.

Member State	Approach to setting reference conditions	Number of reference sites	Location of reference sites	Reference criteria used for selection of reference or benchmark sites
Norway	Expert judgement, recent data from least impacted sites	n.a.	Outer coast of Skagerrak, southern Norway.	Reference sites were selected by the following criteria: Deeper than 5m, limited fresh water influence (> 1km from nearest estuary) and of sufficient distance (based on expert judgment) from any known pollution sources, such as large cities or industrial activity.
United Kingdom	Suitable reference conditions are derived based on physiochemical conditions and sampling methodologies using data from undisturbed sites or sites with minor disturbance, combined with expert judgement and models to accommodate changes in habitat	No specific reference sites but data from over 1000 sampled data points contribute to expert judgement assessment and models	No specific reference sites but data from multiple locations from UK coastal and transitional waters	All samples used if of sufficient data quality with matched physicochemical data

C.2.4 National boundary settings

Table 42. Explanations for national boundary setting of the national methods.

Member State	Type of boundary setting: Expert judgment – statistical – ecological discontinuity – or mixed for different boundaries?	Specific approach for HG boundary	Specific approach for GM boundary	BSP: method tested against pressure
Norway	National boundaries (Molvær et al., 1993) adjusted following the NEAGIG Phase 1 intercalibration exercise (Borja et al., 2007)	<p>The HG boundary was set to ensure the benthic communities at good and high status respectively displayed the following characteristics:</p> <p>High ecological status:</p> <p>Diversity and abundance of invertebrates within the range normally associated with pristine conditions. All taxa that are sensitive to disturbance and associated with pristine conditions are present.</p> <p>Good ecological status:</p> <p>Diversity and abundance of invertebrates just outside the range normally associated with type-specific conditions. Most sensitive taxa of the type specific communities are present.</p>	<p>The GM boundary was set to ensure the benthic communities at moderate and good status respectively displayed the following characteristics:</p> <p>Moderate ecological status:</p> <p>Diversity and abundance of invertebrates moderately outside the range normally associated with type-specific conditions. Taxa that indicate disturbance are present. Many of the sensitive species from type specific communities are absent.</p>	
United Kingdom	Boundaries established from the NEAGIG Phase 1 intercalibration exercise (Borja et al., 2007). Full explanation in Phillips et al. 2014.	The HG boundary was set to ensure the benthic communities at good and high status respectively displayed the following characteristics as assessed by AMBI:	The GM boundary was set to ensure the benthic communities at moderate and good status respectively displayed the following	AMBI ecological group proportions were established for samples over a sewage sludge disposal pressure gradient. Initially, equidistant class boundaries were set and each AMBI EG proportion was

Member State	Type of boundary setting: Expert judgment – statistical – ecological discontinuity – or mixed for different boundaries?	Specific approach for HG boundary	Specific approach for GM boundary	BSP: method tested against pressure
		<p>High ecological status:</p> <ul style="list-style-type: none"> • sensitive taxa (EGI) of dominant abundance • indifferent and tolerant taxa (EGII and EGIII) absent or of sub-dominant abundance • opportunistic taxa (EGIV) absent or of negligible abundance • indicator taxa (EGV) absent or of negligible abundance <p>Good ecological status:</p> <ul style="list-style-type: none"> • sensitive taxa (EGI) abundance may range from high sub-dominant to absent • indifferent taxa (EGII) of low sub-dominant abundance • tolerant taxa (EGIII) of dominant abundance • opportunistic taxa (EGIV) and indicator taxa (EGV) abundance may range from negligible or low to equi-abundance with indifferent taxa 	<p>characteristics as assessed by AMBI:</p> <p>Moderate ecological status:</p> <ul style="list-style-type: none"> • sensitive taxa (EGI) of negligible abundance or absent • indifferent taxa (EGII) of low sub-dominant abundance • tolerant taxa (EGIII), opportunistic taxa (EGIV) and indicator taxa (EGV) co- dominate the abundance <p>Good ecological status:</p> <ul style="list-style-type: none"> • sensitive taxa (EGI) abundance may range from high sub-dominant to absent • indifferent taxa (EGII) of low sub-dominant abundance • tolerant taxa (EGIII) of dominant abundance • opportunistic taxa (EGIV) and indicator taxa (EGV) abundance may range from negligible or low to equi-abundance with indifferent taxa 	<p>calculated for i) the overall status and ii) the lower and upper quartiles of the data in each status. Where the AMBI EG proportions did not conform to those interpreted from the WFD Normative Definitions, the status boundary was adjusted towards the quartile that gave a more accurate representation. Boundaries were further optimised during Intercalibration Phase I.</p>

C.2.5 Results of WFD compliance checking

Table 43. List of the WFD compliance criteria and the WFD compliance checking process and results of the national methods included in the IC exercise.

Compliance criteria	Compliance checking conclusions
1. Ecological status is classified by one of five classes (high, good, moderate, poor and bad).	Yes
2. High, good and moderate ecological status are set in line with the WFD's normative definitions (Boundary setting procedure)	Yes
3. All relevant parameters indicative of the biological quality element are covered (see Table 1 in the IC Guidance)?	Yes
4. Assessment is adapted to intercalibration common types that are defined in line with the typological requirements of the Annex II WFD and approved by WG ECOSTAT?	Yes
5. The water body is assessed against type-specific near-natural reference conditions?	Yes – reference conditions are adapted for specific habitats and sample collection and processing method.
6. Assessment results are expressed as EQRs?	Yes
7. Sampling procedure allows for representative information about water body quality/ecological status in space and time?	Yes
8. All data relevant for assessing the biological parameters specified in the WFD's normative definitions are covered by the sampling procedure?	Yes
9. Selected taxonomic level achieves adequate confidence and precision in classification?	Yes

Conclusion on compliance checking: Both National methods meet the compliance criteria.

C.3 Feasibility checking

C.3.1 Typology

The common intercalibration water body type, NEA7, shared between Norway and the United Kingdom is described below:

Common IC type	Type characteristics	MS sharing IC common type
NEA7 - Deep, fjordic type	Fully saline (>30), mesotidal (1-5m), deep (>30m), sheltered, low current velocity (< 1knot)	Norway, United Kingdom (Scotland)

What is the outcome of the feasibility evaluation in terms of typology? Are all assessment methods appropriate for the intercalibration water body types, or subtypes?

Method	Appropriate for IC type	Remarks
NQIV	Yes	Soft sediment benthic infauna assessment
IQIV	Yes	Soft sediment benthic infauna assessment
<p>Conclusion</p> <p>The Intercalibration is feasible in terms of typology. Both classification schemes intercalibrated relate only to the soft sediment infauna component.</p>		

C.3.2 Pressures addressed

Table 44. Pressures addresses by the national methods included in the Ic exercise and overview of the relationships between national methods and the pressures.

Member State	Method tested	Pressure	Pressure indicators	Amount of data	Strength of relationship
Norway	NQIV	Mine waste	Titanium Dioxide (TiO ₂ %)	n.a.	R2 = 0.8168
		Oxygen deficiency (organic enrichment)	Oxygen (O ₂ (ml/l))	n.a.	R2 = 0.6955
		Urban pollution (industry, boat)	Distance from Oslo harbor (m)	n.a.	R2 = 0.3884

Member State	Method tested	Pressure	Pressure indicators	Amount of data	Strength of relationship
		traffic, road traffic, run off, waste water)			
		Industry	Nickel (Ni ppm)	n.a.	R2 = 0.6498
United Kingdom	IQiVIV	Sewage Sludge disposal (organic enrichment and metals)	Contaminant concentration	169 samples	R2 = 0.674 (p<0.001)
		Mine waste (particulates and metals)	Contaminant concentration, sediment loading	212 samples	R2 = 0.455 (p<0.001)
		Aquaculture (organic enrichment and biocides)	Distance from pressure	326 samples	Varies between sites, average R2 = 0.57

Plots showing the relationship between the Norwegian method and pressures:

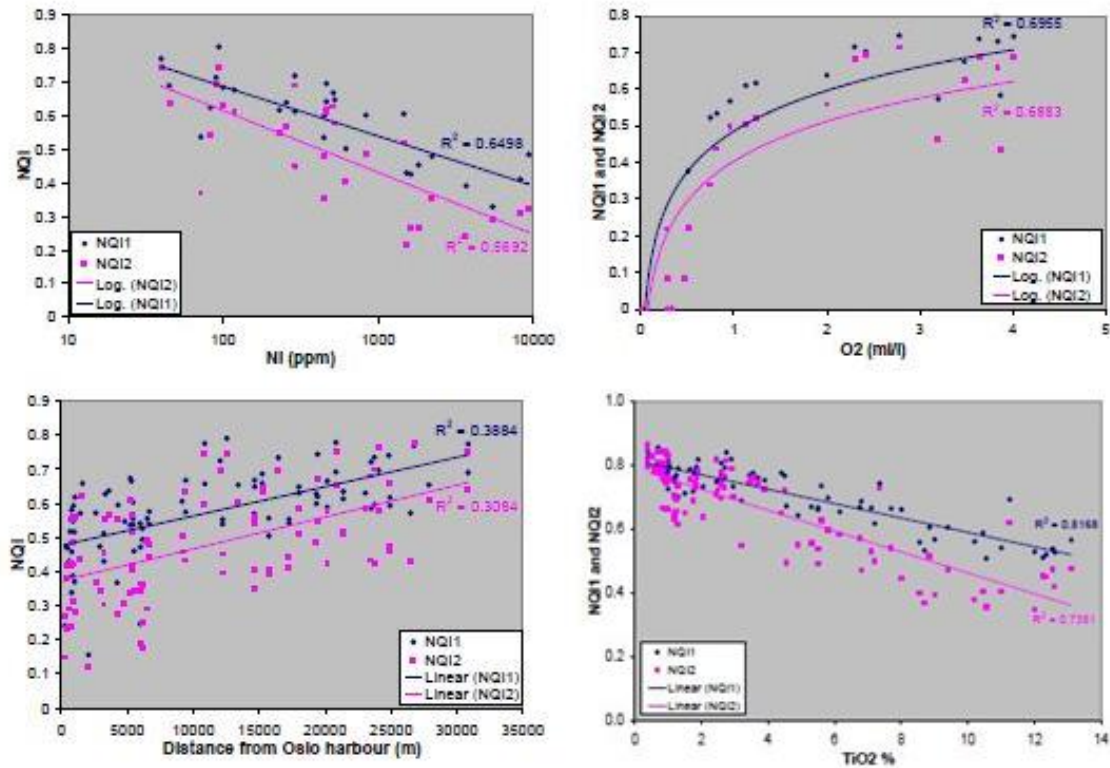


Figure 16. Correlation between Norwegian method and pressures

Plots showing the relationships between the Infaunal Quality Index (IQiVIV) and pressures:

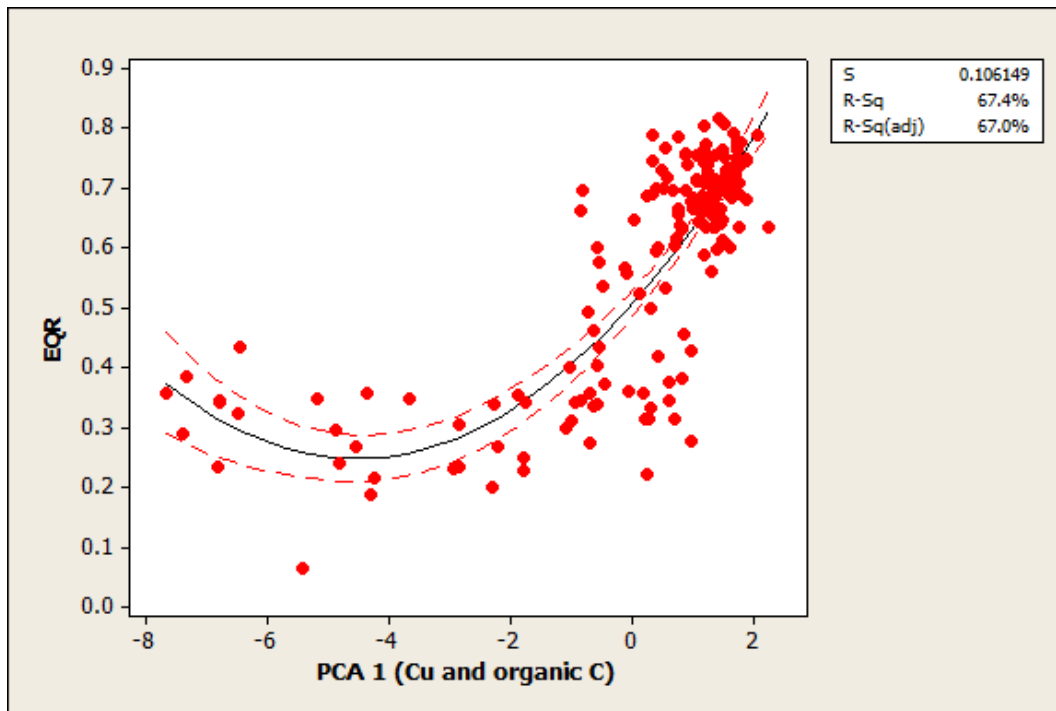


Figure 17. Correlation between EQR (Infaunal Quality index) and principal component (PCA1) of Cu and organic carbon data (sewage sludge disposal pressure).

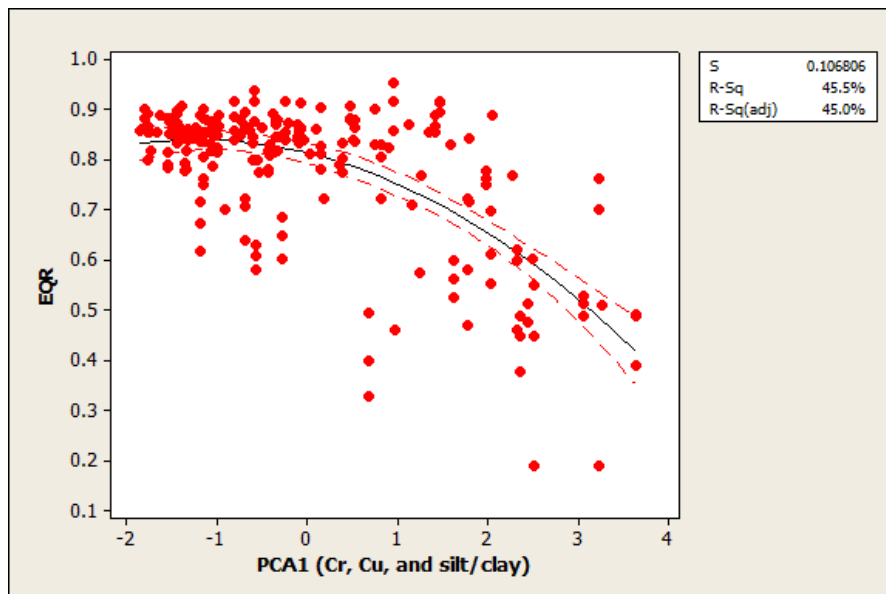


Figure 18. Correlation between EQR (Infaunal Quality index) and principal Component assessment (PCA1) of Cu, Cr and silt/clay data (Mine waste pressure).

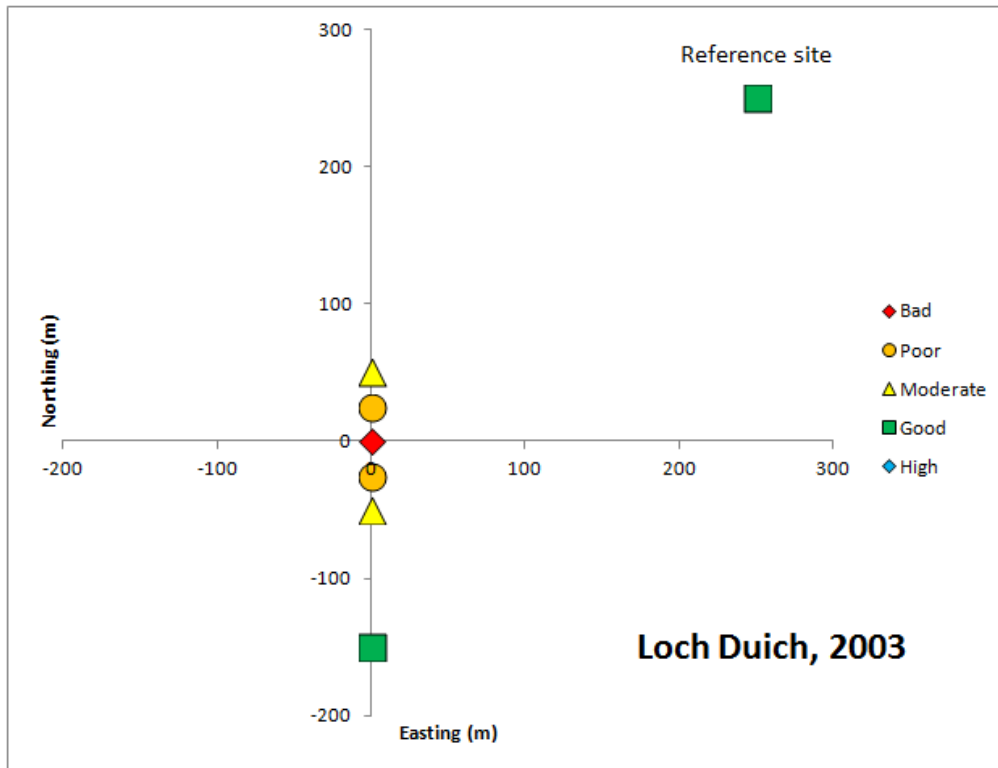


Figure 19. Ecological status as assessed by the Infaunal Quality IQivIV at distance from fish farm pressure (Loch Duich, 2003)

Figure 19. Ecological status as assessed by the IQivIV at distance from fish farm pressure (Loch Duich, 2003).

Conclusion: Both assessment methods have been demonstrated to have a measurable response to pressure.

The NQivI correlates to several different pressures, including oxygen deficiency, industrial pollution and mine waste (Rygg 2011).

The IQivIV has been demonstrated to correlate to a selection of different pressures, including organic enrichment, metal contamination and sediment loading. Included within this is the response of the IQivIV to fish farming, which is an important pressure within type NEA7 water bodies.

C.3.3 Assessment concept

The benthic assessment approaches used by Norway and the United Kingdom follow a similar assessment concept.

Method	Assessment concept	Remarks
IQiVIV / NQiVI	<p>These approaches consist of different parameters (AMBI, number of taxa, Shannon wiener, Simpson, or abundance) and a different algorithm (factorial or simple algorithm).</p> <p>The assessment is performed on sample level.</p>	The simple algorithm differences is based on a different weighing of the parameters or using it as a correction factor (e.g. abundance)

C.4. Collection of intercalibration dataset and benchmarking

C.4.1 Dataset description

The NEA7 benthic dataset contains 426 samples with standardised sampling methodology (0.1 m² grab, processed using 1 mm sieve). These data were originally collated in the NEAGIG Intercalibration Phase I.

The dataset comprises data provided by the Norsk Institutt for Vannforskning (NIVA, 100 samples, including reference sites) and the Scottish Environment Protection Agency (SEPA, 326 samples from aquaculture, including impact and reference sites).

The EQRs in the analysis have been calculated using data truncated according to the 2008 UK data treatment rules (UK truncation rules were also applied to the IC dataset for all MS calculations in Phase I). Details of the 2008 UK data treatment rules are available in Phillips et al (2014).

Table 45. Overview of the number of sites/samples/data values.

Member State	Number of sites or samples or data values		
	Biological data	Physico-chemical data	Pressure data
Norway	100	100	31 sites described as non-reference
United Kingdom	326	326	Distance from pressure source provided with all samples

Table 46. Overview of the data acceptance criteria used for the data quality control

Data acceptance criteria	Data acceptance checking
Data requirements (obligatory and optional)	Sample level, quantitative, benthic invertebrate data, Definition of the habitat sampled (sediment parameters from particle size analysis or qualitative description)
The sampling and analytical methodology	0.1 m ² grab data, processed using a 1 mm sieve
Level of taxonomic precision required and taxa lists with codes	Lowest taxonomic level. Taxon lists closely aligned with the Ulster Museum and Marine Conservation Society Marine Species Directory and AMBI score lists (www.azti.es).
The minimum number of sites / samples per intercalibration type	Yes – exceeds minimum number of data records of 20-25 per Member State as recommended in Intercalibration guidance version 5 (September 2010).
Sufficient covering of all relevant quality classes per type	Yes – pressure gradient data included

C.4.2 Common benchmarking or reference conditions

For the Intercalibration of the common type, NEA7, common reference conditions were defined.

Reference conditions

For the Norwegian NEA7 data, reference sites were identified in accordance with expert judgement considering distance from pressure sources and the physical characteristics of the sites.

The United Kingdom NEA7 data was based on surveys monitoring the effects of fish farms. Reference sites were defined for these surveys as being beyond the influence of the fish farms and other anthropogenic pressures. (Samples where the percentage of the silt/clay fraction was >90% were excluded from the reference set on the basis that the poor circulation of dissolved oxygen through the sediments were likely to be impacting the biology, resulting in the samples being non-representative of reference conditions in relation to the rest of the data.)

Reference sites

The number of reference sites for Norway and the United Kingdom were 31 and 104 respectively. For both Member States, this exceeds the recommended minimum requirements of 20-25 discrete data points classified by each Member State as described in the Intercalibration guidance version 5 (September 2010) so is considered sufficient for the process.

Benchmark standardization

To account for potential biogeographical differences between Norway and the United Kingdom, data from each Member State was assigned a different subtype. Reference sites were present in each subtype of the common dataset.

C.5 Comparison of methods and boundaries

C.5.1 Intercalibration option and common metrics

For the Intercalibration of the common type, NEA7, IC option 3 was used.

IC Option 3 - Similar data acquisition, but different numerical evaluation (BQE sampling and data processing generally similar, so that all national assessment methods can reasonably be applied to the data of other countries)

As specified by Intercalibration Guidance version 5.0 (September 2010), Option 3a (direct comparison with regression) was used for the Intercalibration between Norway and the United Kingdom as the approach was i) based on commonly assessed sites, ii) inclusive of data from across a pressure gradient and iii) based on <3 methods.

C.5.2 Results of regression comparison

The correlation coefficient (r) and the probability (p) for the correlation of the methods (only two methods included in this common type) are shown below.

Member State/Method	R	P
Norwegian NQIvI vs. United Kingdom IQIvIV	0.992 (Pearson correlation)	<0.001

- the relationship is highly significant $p \leq 0.001$
- assumptions of normally distributed error and variance (homoscedasticity) of model residuals are met
- both methods adequately represent the other method ($r^2 > 0.5$)
- the slope of the regression lies between 0.5 and 1.5.

C.5.3 Comparability criteria

The comparison of national boundaries using comparability criteria (see Annex V of IC guidance) is summarised below.

Boundary bias

Boundary bias between the High/Good and Good/Moderate boundaries are provided below:

	NQIVI	IQIVIV
H/G bias	0.081	0.227
G/M bias	0.199	0.099

- In each case, bias is within the acceptable limits of between -0.25 and 0.25.
- Assessing class agreement (absolute average class difference)
- Average class difference between the NQIVI and IQIVIV is 0.185. This is below the required threshold of 1 and is therefore acceptable.
- Kappa agreement:
- The Kappa agreement coefficient between the NQIVI and IQIVIV is 0.921. This is above the required threshold of 0.4 and is therefore acceptable.

C.6 Final results to be included in the EC

Table with EQRs

Table 47. Overview of the Ic results for the national methods included in the Ic exercise. The results are included in the Part I of the Annex of the EC Decision.

Member state	National classification system intercalibrated	Ecological Quality Ratios	
		High-Good boundary	Good-Moderate boundary
Norway	Norwegian Quality Index (NQIVI)	0.72	0.63
United Kingdom	Infaunal Quality Index (IQIVIV)	0.75	0.64

Correspondence common types versus national types

Common boundaries will be applied within the national systems of Norway and the United Kingdom as presented in the above table.

In the UK, common European type NEA7 equates to UK coastal water type 11 (CW 11). These boundaries will be utilised in all coastal water types, with the specific reference conditions for the samples defined by habitat and sampling method.

C.7 Ecological characteristics

C.7.1 Description of reference or alternative benchmark communities

Reference condition macrobenthic communities are dominated by pollution sensitive taxa (e.g. AMBI Ecological Group (EG) I taxa), have low relative abundance of indifferent (EG II) and tolerant (EG III) taxa and negligible relative abundance of opportunist (EG IV) and pollution indicator (EG V) taxa. High numbers of taxa with an even abundance distribution throughout the community are also indicative of reference conditions. Communities are also characterized by relatively high species numbers and evenness.

C.7.2 Description of good status communities

At good ecological status, taxa number and Simpsons evenness are slightly reduced in comparison to values under reference conditions, whilst variables according to habitat (community abundance as assessed by AMBI) are slightly unbalanced: sensitive taxa (EG I) abundance may range from high sub-dominant to absent; indifferent taxa (EG II) are of low sub-dominant abundance; tolerant taxa (EG III) of dominant abundance; abundance of opportunistic (EG IV) and indicator taxa (EG V) may range from negligible or low to comparable abundance with indifferent taxa (EG II).

Under borderline conditions, taxa number and Simpson's evenness are expected to be slightly to moderately reduced in comparison to reference conditions. In terms of community abundance as assessed by AMBI; sensitive taxa (EGI) abundance is expected to be between high sub-dominant to absent; indifferent taxa (EGII) are expected to be of low sub-dominant in abundance; tolerant taxa (EGIII) are expected to be between dominant and co-dominant in abundance; opportunistic taxa (EGIV) are expected to be between negligible to co-dominant in abundance and; indicator taxa (EGV) are expected to be between negligible to co-dominant in abundance.

PART D-Type NEA 5

The type NEA 5 covers the small (18.5 km²) water body which represents the euhaline rocky coastal water around Helgoland. The salinity is >30 PSU. Due to its unique hydromorphological characteristics the type NEA 5 is not part of a common intercalibration type and has not been part of the intercalibration process.

D.2 Description of national assessment methods

Method: Marine Biotic Index Tool (MarBIT) adapted to NEA 5 conditions based on three different Habitats. The index uses the metrics abundance, species richness, the proportion of sensitive and the proportion of tolerant taxa to calculate a quality status for each metric. Based on autecological species data and historical references, different lists of taxa serve as references for each differentiated habitat sampled.

The data are processed using different methods like taxonomic spread, log-normal abundance distribution etc. The results are then normalized to calculate the EQR range. The median of all four metric EQRs serves as the final status assessment for each habitat. Three habitats are differentiated in NEA 5 Helgoland. The final quality assessment either uses the EQRs separately as calculated for each of the three habitats or one EQR combined from the 3 sub-EQRs by averaging.

D.2.1 Methods and required BQE parameters

Table 48. Overview of the metrics included in the national assessment methods.

Member State	Full BQE met	Composition	Abundance	Disturbance sensitive taxa	Taxa indicative of pollution	Combination rule of metrics
Germany	Yes	Taxonomic spread index TSI based on reference taxa list for each area	Correlation with reference log-normal abundance distribution	Fraction of taxa sensitive to disturbance in relation to reference taxa list for each area	Fraction of taxa tolerant to disturbance in relation to reference taxa list for each area	Weighted algorithm. See National description.

D.2.2 Sampling and data processing

Table 49. Overview of the sampling and data processing of the national assessment methods.

Sampling/survey device	<p>Intertidal: 50 x 50 cm frame with subdivisions of 5x5 cm</p> <p>Laminaria-holdfasts: manual sampling by divers of holdfast in a bag with ambient water retaining mobile fauna</p> <p>Tiefe Rinne: Dredging, subsample of 3 replicates of 2 L volume</p>
How many sampling/survey occasions (in time) are required to allow for ecological quality classification of survey site or area?	1 survey per year
Sampling/survey months	Summer: June-September
Which method is used to select the sampling /survey site or area?	<p>Intertidal: measuring percentage cover</p> <p>Laminaria-holdfasts: collecting 10 holdfasts at two different sites;</p> <p>sampling all mobile species through successive sieves down to 300µm mesh-size; recording sessile species directly on Laminaria-holdfasts</p> <p>Tiefe Rinne: dredging at 5 different transects for 2 min each; sampling all mobile species through successive sieves down to 300µm mesh-size;</p> <p>recording sessile species directly on substratum (mainly shells and fewrocks); all samples fixed in 4 % Formalin/seawater or 70 % alcohol,</p> <p>taken to the lab for species identification and counting</p>
Identification level	Whenever possible down to species level according to available and most recent identification references. All macroscopic species identified, according to international nomenclature and national quality guidelines.
Data processing	All data are listed in spread-sheets showing abundance either as number of individuals per sample (most mobile fauna), relative abundance based on frequency per unit substratum (most sessile fauna), or percentage cover (fauna of intertidal habitat). These data are exported into the MarBIT to calculate the different metrics.

D.2.3 National reference conditions

Reference conditions were derived/modelled from collected and analysed autecological data of potentially occurring species and the corresponding abiotic conditions in the water body. Together with the analysis of historical samples, this resulted in species reference lists valid for the water body (= assessment unit). The only waterbody of Helgoland represents the water type NEA 5.

D.2.4 National boundary settings

Table 50. Explanations for national boundary setting of the national methods.

Member State	Type of boundary setting: Expert judgment – statistical – ecological discontinuity – or mixed for different boundaries?	Specific approach for HG boundary	Specific approach for GM boundary	BSP: method tested against pressure
	<p>Mixed boundary setting</p> <p>The Moderate/Poor and Poor/Bad boundaries were derived from the normative definitions and translated into ecologically sensible values for each of the 4 metrics in the MarBIT.</p> <p>The boundary index values were divided into five groups reflecting quality classes – High, Good, Moderate, Poor and Bad - applying the natural breaks method (Jenks and Caspall, 1970), included in the ArcGIS software</p>	<p>The HG boundary was in general set approximately halfway from the GM boundary and up to the maximum index value. If possible, the normative definitions were applied, taking into account the variability of the metrics at reference conditions.</p>	<p>The GM boundary was set where a statistical significance occurs with respect to the change of the 4 metrics in the MarBIT from the reference value (derived individually and separately for each of the 4 metrics).</p>	<p>No, not Helgoland</p>

D.2.5 Results of WFD compliance checking

Table 51. List of the WFD compliance criteria and the WFD compliance checking process and results of the national methods included in the Ic exercise.

Compliance criteria	Compliance checking conclusions
1. Ecological status is classified by one of five classes (high, good, moderate, poor and bad).	Yes
2. High, good and moderate ecological status are set in line with the WFD's normative definitions (Boundary setting procedure)	Yes
3. All relevant parameters indicative of the biological quality element are covered (see Table 1 in the IC Guidance)?	Yes
4. Assessment is adapted to intercalibration common types that are defined in line with the typological requirements of the Annex II WFD and approved by WG ECOSTAT?	No, NEA 5 is a type not shared by MS. Only Germany.
5. The water body is assessed against type-specific near-natural reference conditions?	No, no reference sites available
6. Assessment results are expressed as EQRs?	Yes
7. Sampling procedure allows for representative information about water body quality/ecological status in space and time?	Yes
8. All data relevant for assessing the biological parameters specified in the WFD's normative definitions are covered by the sampling procedure?	Yes
9. Selected taxonomic level achieves adequate confidence and precision in classification?	Yes

Conclusion on compliance checking: Both National methods meet the compliance criteria.

D.3 Feasibility checking

D.3.1 Typology

Due to its unique hydromorphological characteristics the type NEA 5 is not part of a common intercalibration type and has not been part of the intercalibration process. **Therefore IC is not feasible**

D.3.2 Pressures addressed

D.3.2 The index addresses eutrophication and/or general degradation as the main pressures similar to other methods.

D.4 Ecological characteristics

D.4.1 Description of reference or alternative benchmark communities

No or very scarce anthropogenic pressures. There is a diverse community of mobile and sessile species with high species richness. Species richness is similar to that of the historical reference. Tolerant species at low abundance, whereas many sensitive taxa are present.

D.4.2 Description of good status communities

Anthropogenic pressures are low. There is a slight deviation in species abundance and richness from reference sites and with lower levels of species richness. Tolerant species show increased abundance and sensitive taxa are well presented but less abundant.

Conclusions

Coastal water bodies has been classified into different types. The IC exercise has been successfully completed for all these common types.

The benthic assessment approaches of all Member States meet all WFD compliance criteria. Only, the benthic assessment approach of the Andalusia region (Spain) does not meet the requirements of compliance criteria N°3, due to the lack of a diversity parameter within their approach (scientific justification available and accepted by review panel).

All methods described can show in one or another way, a certain response to certain pressures.

IC was feasible for all Member States, excepting for BO2A and RAT methods (included in the common type NEA 1/36) and the MarBit method (in the common type NEA 5)

References

- Birk, S., Strackbein, J. & Hering, D., 2010. WISER methods database. Version: March 2011. Available at <http://www.wiser.eu/results/method-database/>.
- Birk, S., Willby, N.J., Kelly, M.G., Bonne, W., Borja, A., Poikane, S., van de Bund, W., 2013. Intercalibrating classifications of ecological status: Europe's quest for common management objectives for aquatic ecosystems. *Science of the total Environment* 454-455, 490-499.
- Borja, A., Franco, F., Valencia, V., Bald, J., Muxika, I., Belzunce, M.J., et al., 2004. Implementation of the European Water Framework Directive from the Basque country (northern Spain): a methodological approach. *Marine Pollution Bulletin* 48 (3-4), 209-218.
- Borja, A., Josefson, A.B., Miles, A., Muxika, I., Olsgard, F., Phillips, G., Rodriguez, J.G., Rygg, B., 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., Muxika, I., Rodriguez, J.G., 2009. Paradigmatic responses of marine benthic communities to different anthropogenic pressures, using M-AMBI, within the European Water Framework Directive. *Marine Ecology – An Evolutionary Perspective* 30, 214-227
- Borja et al., A., A. Miles, A. Occhipinti-Ambrogi, T. Berg, 2009. Current status of macroinvertebrate methods used for assessing the quality of European marine waters: implementing the Water Framework Directive. *Hydrobiologia*, 633: 181-196.
- Davies, Susan P., 2012. Peer review of the intercalibration exercise phase II: European water framework directive.
- Fitch, J.E., Cooper, K.M., Crowe, T.P., Hall-Spencer, J.M., Philips, G., 2014. Response of multi-metric indices to anthropogenic pressures in distinct marine habitats: The need for recalibration to allow wider applicability. *Marine Pollution Bulletin* dx.doi.org/10.1016/j.marpolbul.2014.07.056.
- ISO 16665:2013 (2013). Water quality - Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna, International Organization for Standardization, Geneva, Switzerland.
- Josefson, A.B., Blomqvist, M., Hansen, J.L.S., Rosenberg, R., Rygg, B., 2009. Assessment of marine benthic quality change in gradients of disturbance: comparison of different Scandinavian multi-metric indices. *Marine Pollution Bulletin* 58, 1263-1277
- Marques J.C., F. Salas, J. Patrício, H. Teixeira & J.M. Neto. 2009. *Ecological Indicators for Coastal and Estuarine Environmental Assessment. A user guide*. Ed. 00, ISBN: 978-1-84564-209-9. UK: WIT Press.
- Molvær, J., Knutzen, J., Magnusson, J., Rygg, B., Skei, J., & Sørensen, J. (1993). *Classification of environmental quality in fjords and coastal waters. A guide*. ISBN: 82-7655-267-2.

Phillips, G.R., Anwar, A., Brooks, L., Martina, L.J., Miles, A.C., Prior, A., 2014. Infaunal Quality Index: Water Framework Directive Classification Scheme for Marine Benthic Invertebrates Environment Agency (UK) R&D Technical Report, No SC080016.”

Phillips, G. R., Miles, A. C., Prior, A., Martina, L. J., Brooks, L., & Anwar, A. (2014). Infaunal Quality Index: WFD classification scheme for marine benthic invertebrates. R&D Technical Report. Bristol: Environment Agency. ISBN: 978-1-84911-319-9.

Teixeira, H., Neto, J.M., Patrício, J., Veríssimo, H., Pinto, R., Salas, F. & Marques, J.C. 2009. Quality assessment of benthic macroinvertebrates under the scope of WFD using BAT, the Benthic Assessment Tool. *Marine Pollution Bulletin*, 58: 1477-1486. (doi:10.1016/j.marpolbul.2009.06.006).

Van Hoey, Gert; Borja, Angel; Birchenough, Silvana; Degraer, Steven; Fleischer, Dirk; Kerckhof, Francis; Magni, Paolo; Buhl-Mortensen, Lene; Muxika, Iñigo; Reiss, Henning; Schröder, Alexander; Zettler, Michael, 2010. The use of benthic indicators in Europe: from the Water Framework Directive to the Marine Strategy Framework Directive. *Marine Pollution Bulletin* 60: 2187-2196

Van Hoey, Gert; David Cabana Permuy; Magda Vincx ; Kris Hostens, 2013. An Ecological Quality Status assessment procedure for soft-sediment benthic habitats: Weighing alternative approaches. *Ecological Indicators* 25, 266-278

Van Hoey, G., Vanaverbeke, J., Degraer, S., 2014. Study related to the realization of the Water Framework Directive intercalibration for the Belgian Coastal waters, to design the descriptive elements 1 and 6 of the Marine Strategy Framework Directive and the nature objectives of the Habitat Directive for invertebrate bottom fauna of soft substrates. ILVO-mededeling 170.

van Loon, W.M.G.M., Boon A.R., Giitenberger, A., Walvoort, D.J.J., Lavaleye, M., Duineveld, G.C.A. and Verschoor A.J., 2015. Application of the Benthic Ecosystem Quality Index 2 to benthos in Dutch transitional and coastal waters. *Journal of Sea Research* 103, 1-13

Rygg, B. (2006). Developing indices for quality-status classification of marine soft-bottom fauna in Norway. NIVA report 5208. ISBN: 82-577-4927-3.

Rygg, B. (2011). Uttesting av indekser på marin bløtbunnsfauna. NIVA report 6255. ISBN: 978-82-5990-2.

UKTAG (2014) UKTAG Transitional and Coastal Water Assessment Methods, Benthic invertebrate fauna, invertebrates in soft sediments, Infaunal Quality Index. ISBN:978-1-906934-34-7.

Annex 1. Common type NEA 1/26: Sampling and data processing information

	Denmark	Belgium	United Kingdom / Ireland	Germany
Sampling guideline	<p>Holme, N.A. & A.D. McIntyre, 1984. Methods for the study of marine benthos. IBP Handbook 16, Blackwell, Oxford.</p>	<p>ISO standard (ISO 16665:2005(E)) "Water quality – Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna".</p>	<p>ISO standard (ISO 16665:2005(E)) "Water quality – Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna".</p>	<p>Muster-Standardarbeitsanweisung für Laboratorien des Bund/Länder-Messprogramms Prüfverfahren-SOP: Makrozoobenthos-Untersuchungen in marinen Sedimenten (Weichboden)</p>
Sampling description	<p>Three to six Van Veen are taken (blindly) at a site or area using ships. Alternatively 40 Haps are taken, one at each geographical position, mostly regularly spaced within an area. For the case of point sites, 5-10 Haps are taken blindly at each site and sampling occasion.</p>	<p>Habitat approach, the main habitat types within a water body were sampled in such way to get a confident ecological quality classification (enough samples, spatially and eventually temporal distributed within a habitat). The samples were taken randomly within the habitat area.</p>	<p>Sampling design variable according to UK and Ireland monitoring authority. Samples taken from soft bottom habitats, either i) spread as single samples or ii) taken as replicates at one or more stations. Surveys are undertaken either i) annually or ii) once in a reporting cycle according to monitoring authority. Biological samples require an associated sediment field sample for particle size analysis and supporting depth and salinity information.</p>	<p>5 to 20 sediment samples are taken from 1 ecotope. Each sample is sieved separately (1mm, 0,5mm mud) and residue is stored and transferred to the laboratory. Benthic species are separated and identified to the lowest taxonomic level.</p>

	Spain (Basque country, Cantabria region)	Netherlands	Portugal
Sampling guideline	ISO standard (ISO 16665:2005(E)) “Water quality – Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna”.	STOWA, 2009. Instructie; Richtlijn Monitoring Oppervlaktewater en Protocol Toetsen en Beoordelen (28 april 2009); STOWA, NN. Quality Handbook Hydrobiology (in prep).	ISO standard (ISO 16665:2005(E)) “Water quality – Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna”.
Sampling description	2-6 sampling locations are visited per water body once a year in winter. At each location 3 van Veen grab replicates are taken (0.1 square-metres each), and sieved on board by 1 mm mesh.	Normally sediment cores are collected at sampling stations with a device like the Reineck Box corer operated from a ship for subtidal stations. The sediment is washed through a 1mm mesh. Specimens are sorted from the residue, identified to the species level, counted and weighed. Biomass is most accurately measured by the difference between dry weight and ash weight, the ash free dry weight AFDW.	Biological samples are collected from soft bottom habitats, by using a 0.1 m ² sampling area Van Veen Grab (or equivalent). Sampling stations are placed at representative sites of water bodies, and in sufficient number to cover natural variations, according to monitoring authority. A minimum of 3 replicates per sampling station are collected. Biological samples require an associated sediment field sample for particle size and organic matter content analysis, and supporting depth, salinity, temperature, and chemical parameters information (bottom water).

	Spain (Andalusia)	France	Norway
Sampling guideline	ISO standard (ISO 16665:2005(E)) “Water quality – Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna”.	ISO standard (ISO 16665 :2005(E)) “ Water Quality – Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna “	ISO standard (ISO 16665:2005)) Water quality – “Guidelines for quantitative sampling and sample processing of marine soft-bottom macrofauna”.
Sampling description	Overall, one sampling station was defined for each water body, provided it was considered representative of the whole water body. Soft-bottom sampling is carried out, in broad daylight, with the vessels owned by the Regional Agency of Environment (Regional Government of Andalucía), except in shallower areas where it may be carried out by direct sampling or with small auxiliary vessels. A sample corresponds to the average of 3 sampling units. The sediment collected in each sampling unit is posteriorly sieved through a 0.5 mm mesh.	Above all, a monitoring location is defined on the basis of its representativeness across the whole WB. In order to consider the intra-stational variability, it was decided that each location will be studied in 3 points (3 replicates per point), bringing to 9 the number of replicates for each monitoring locations. In subtidal areas, the sampling (one replicate) is carried out by the mean of a grab (area=0.1 m ²) and sieved on board by a 1mm mesh. In intertidal areas, the sampling (one replicate) is carried out by the mean of a hand corer (area = .029 m ²) and sieved by a 1mm mesh. Biological samples require an associated sediment field sample (each of the 3 points constituting the monitoring location), for analysis of particle size and organic matter.	Samples collected by using a 0.1m ² van Veen grab, and sieved on board by 1 mm mesh. 4 replicates per station. An associated sediment field sample taken for grain size and TOC.

	Denmark	Belgium	United Kingdom / Ireland	Germany
Method to select the survey site	Expert knowledge, Random sampling/surveying	Stratified Random sampling/surveying	Stratified Random sampling/surveying	Expert knowledge, Random sampling/surveying
Sampling Device	Corer, Grab	Grab	Corer, Grab	Corer, Grab
Specification of sampling device	0.1 m ² Van Veen Grab, 0.0143 m ² Haps-corer	Van Veen Grab (0.1m ²)	Van Veen Grab (0.1m ²), Day Grab (0.1m ²), Hand Core (0.01m ²)	Van Veen-grab (0.1m ²), corers with 9-15cm diameter
Sampled habitat	Single habitat(s)	All available habitats per site (Multi-habitat)	Single habitat(s)	Single habitat(s)
Specification of sampled habitat	Soft bottom (sand - mud)	soft bottom sediments (muddy sediments [Macoma balthica habitat], fine muddy sand [Abra alba habitat], clean sands [Nephtys cirrosa habitat])	Soft bottom	Soft bottom
Sampled zones in tidal areas	Subtidal zone	Subtidal zone	Both tidal zones	Both tidal zones
Sampling months	April to June	October	February to May (current recommended target months)	May or September/October
Number of sampling occasions in time	One per year	One occasion per year (preferential autumn)	Minimum of one occasion for classification (varies between 1-3 for UK and Ireland monitoring authorities)	One occasion per sampling season

	Spain (Basque country, Cantabria region)	Netherlands	Portugal
Method to select the survey site	Expert knowledge; Fixed sampling stations, representative of the water body	Fixed locations	Expert knowledge
Sampling Device	Grab	Corer	Grab
Specification of sampling device	Van Veen Grab	corer tube; box corer (e.g. Reineck Box corer), flushing sampler (only in saline lakes 0-2 m)	Van Veen Grab (0.1 m ²) or equivalent
Sampled habitat	Single habitat(s)	Single habitat(s)	Single habitat(s)
Specification of sampled habitat	Soft bottom	All present habitats in the water body.	Soft bottom (sandy-mud)
Sampled zones in tidal areas	Both tidal zones	Both tidal zones	Subtidal zone
Sampling months	Winter (Basque country); Summer (Cantabria)	Coastal water types (NEA1, NEA3): March 1st to June 15th	February - March
Number of sampling occasions in time	Once a year	Minimum one survey per year (preferably fall), and scores and classification preferably averaged over three years.	Minimum of one occasion per the chosen sampling season

	Spain (Andalusia)	France	Norway
Method to select the survey site	Expert knowledge	Expert knowledge, Fixed sampling stations representative of the WB	Expert knowledge
Sampling Device	Grab	grab	grab
Specification of sampling device	Van Veen Grab	Van Veen Grab or Day Grab or Smith-McIntyre Grab	Van Veen grab (0.1 m ²)
Sampled habitat	Single habitat(s)	Single habitat(s)	Single habitat(s)
Specification of sampled habitat	Soft bottom	Soft bottom	Soft bottom
Sampled zones in tidal areas	Subtidal zone	Subtidal and intertidal zone	Subtidal zone
Sampling months	Summer: June - August	From February to April	May, August, September
Number of sampling occasions in time	One occasion per sampling season	One occasion per sampling season	one per year

	Denmark	Belgium	United Kingdom / Ireland	Germany
Number of spatial sampling replicates	Six 0.1 m ² Van Veen, or 40 Haps samples	Depends on habitat type samples (18 for <i>Macoma balthica</i> habitat, 20 for <i>Abra alba</i> habitat and 18 for <i>Nephtys cirrosa</i> habitat)	Variable according to habitat, number of years/ stations, methodology and required confidence.	6-10 replicates per ecotope
Total sampled area or duration	0.6 m ²	Depends on habitat type samples (1.8 m ² for <i>Macoma balthica</i> habitat, 2.0 m ² for <i>Abra alba</i> habitat and 1.8 m ² for <i>Nephtys cirrosa</i> habitat)	Variable according to habitats, number of years/ stations, methodology and required confidence.	1 m ² per ecotope, 2-4 ecotopes per water body, average of several years
Minimum size of sampled organisms	1 mm (mesh-size of sieve)	1 mm	1000 µm (Coastal Waters)	1000 µm, 500 µm in mud sediments
Sample treatment	Organisms of the complete sample are identified.	Organisms of the complete sample are identified.	Organisms of the complete sample are identified.	Organisms of the complete sample are identified.
Level of taxonomic identification	Other, Species/species groups	Family, Genus, Other, Species/species groups	Species/species groups	Genus, Species/species groups

	Spain (Basque country, Cantabria region)	Netherlands	Portugal
Number of spatial sampling replicates	3 replicates per station (2-6 stations per water body)		Variable according to habitat, number of years/stations, and required confidence.
Total sampled area or duration	0.3 m ² (each replicate has 0.1 m ²)		Variable according to habitat, number of years/stations, and required confidence.
Minimum size of sampled organisms	1 mm mesh	1 mm	1000 µm for Coastal Waters
Sample treatment	Organisms of the complete sample are identified.	Organisms of the complete sample are identified.	Organisms of the complete sample are identified.
Level of taxonomic identification	Species/species groups	Species/species groups	Other, Species/species groups

	Spain (Andalusia)	France	Norway
Number of spatial sampling replicates	3	3	4
Total sampled area or duration	0.025 m ² (average of 3 spatial replicates)	0,9m ² (3 locations, 3 replicates per location)	0,4m ²
Minimum size of sampled organisms	0.5 mm mesh size	1 mm	1 mm
Sample treatment	Organisms of the complete sample are identified.	Organisms of the complete sample are identified.	Organisms of the complete sample are identified.
Level of taxonomic identification	Family, Other, Species/species groups	Species/species groups	Species/species groups

	Denmark	Belgium	United Kingdom / Ireland	Germany
Specification of level of determination	Species level (or if not possible to determine, genus or family level): Echinodermata, Polychaeta, Crustacea, Mollusca; Higher Group level: Nemertea, Nematoda, Turbellaria	Determination to the lowest level possible. Oligochaeta to level of order. Some Polychaeta to the level of family (Cirratulidae). Taxonomy between assessment and reference data were set consistently.	n.a.	All to species level except some Oligochaeta, Diptera, Priapulida,...
Determination of abundance	Individual counts	Individual counts	Individual counts	Individual counts
Abundance is related to	Area	Area	Area	Area
Unit of the record of abundance	individuals per m ²	Number of individuals per one square-metre	Number of individuals per area of sample	Number of individuals per one m ²
Quantification of biomass	n.a.	Wet weight	n.a.	n.a.
Other biological data	none	none	none	none

	Spain (Basque country, Cantabria region)	Netherlands	Portugal
Specification of level of determination	Some groups can be indentified to higher taxonomical levels.	n.a.	Truncation rules (Borja et al., 2007)
Determination of abundance	Individual counts	Individual counts	Individual counts
Abundance is related to	Area	Area	Area
Unit of the record of abundance	Number of individuals per one m ²	Number of individuals per one m ²	Number of individuals per sampling area
Quantification of biomass	n.a.	n.a.	n.a.
Other biological data	none		none

	Spain (Andalusia)	France	Norway
Specification of level of determination	Plathelminthes, Nemertina and Nematoda to phylum level; oligochaetes to sub-class level; harpacticoid copepods to order level; insects to class level, except chironomids; chironomids to family level; hemichordates to phylum level.	Species level, except for the following groups: <i>Echiura</i> , <i>Hemichordata</i> , <i>Hydrozoa</i> , <i>Insecta</i> , <i>Nemertea</i> , <i>Oligochaeta</i> , <i>Phoronida</i> , <i>Platyhelminthes</i> et <i>Priapulida</i>	Species level or lowest level possible
Determination of abundance	Individual counts	Individual counts	Individual counts
Abundance is related to	Area	Area	Area
Unit of the record of abundance	Number of individuals per one m ²	Number of individuals per 0,1 m ²	Number of individuals per 0,1 m ²
Quantification of biomass	n.a.	n.a.	n.a.
Other biological data	none	none	none

	Denmark	Belgium	United Kingdom / Ireland	Germany
Special cases or additions of sampling	none	none	Presence/absence recorded where taxa are unsuitable for quantification (e.g. colonial taxa). Truncation rules are applied to the data to exclude non-benthic and non-invertebrate fauna from the IQI assessment.	none
Comments on 'data acquisition' part	The DKI is applied on 0.1 m ² samples and therefore Haps samples are pooled to this sample size (6-7 Haps)	none	none	none

	Spain (Basque country, Cantabria region)	Netherlands	Portugal
Special cases or additions of sampling	none		Presence/absence recorded where taxa are unsuit
Comments on 'data acquisition' part	none	The present Dutch surveillance monitoring (BIOMON program) can be split up in 3 areas, with differences in sampling strategy, namely (1) the Delta area, (2) the Dutch coast and (3) the Waddenzee; Eems-Dollard. The macrobenthic fauna monitoring activities are all under the responsibility of one agency (but different offices) could lead to some small taxonomic differences in the methodology. Since these differences also exist in the reference data sets, it is expected that the impact on the EQR-scores are very small.	none

	Spain (Andalusia)	France	Norway
Special cases or additions of sampling	none	none	none
Comments on 'data acquisition' part	none	none	none

Annex 2. Common type NEA 1/26: Alternative benchmark approach (based on biotic variables)

This procedure to determine the benchmark samples out of the common dataset was not accepted by JRC and the review panel. In the authors' point of view, this gives a reliable, objective alternative for the determination of the benchmark samples, which is explained in this annex.

An alternative procedure for the selection of benchmark sites can be used in this intercalibration due to the absence of quantitative and even qualitative pressure data of each sample within the common dataset. The collection of such information on sample level in a standard way is rather impossible (except for some sub-data sets, e.g., the Garroch Head analyses), due to the absence of such information at the Member State level. Alternative pressure quantifications, as general pressure index, distance from the coast, are not appropriate for this NEA-GIG dataset due to the type of data (many samples from the same location), indirect influence of harbors and rivers being rather low for the majority of samples, other pressures being probably more important (local pollution [such as dumping activity at Garroch Head dataset, Basque Country dataset is at a submarine outfall], fishery, and the like). Besides this, the variation in pressure quantification will be low and many samples will be cataloged within the same pressure status, due to the absence of detailed pressure information. Such a general pressure index approach was tested for the intercalibration of transitional waters within the NEA-GIG region in phase II and was inadequate.

For the dataset, where some pressure information was available (see Garroch Head dataset), we could objectively distinguish least disturbed samples (lower copper concentration), and showed that there is some variability in the classification of those samples by the different benthic assessment approaches. Unfortunately this does not meet the set-up of the benchmarking in the intercalibration guidance (benchmark sites in each Member State are necessary).

An approach that estimates the benthic conditions under least disturbed circumstances could be the selection of samples with the highest diversity characteristics. In theory, areas characterized by samples with a high diversity (expressed as any type of diversity indices) are less subjected to pressures on the system than areas characterized by lower diversity (Pearson & Rosenberg, 1978) (**Error! Reference source not found.**). This relationship is not linear, but a clear gradient exists. The multivariate analysis on the common dataset (see higher) show the benthic variability within the data, but also a clear gradient in benthic characteristics (**Error! Reference source not found.**). The gradient within these benthic univariate parameters can be used as a proxy for the pressures on the samples of the NEA-GIG common dataset.

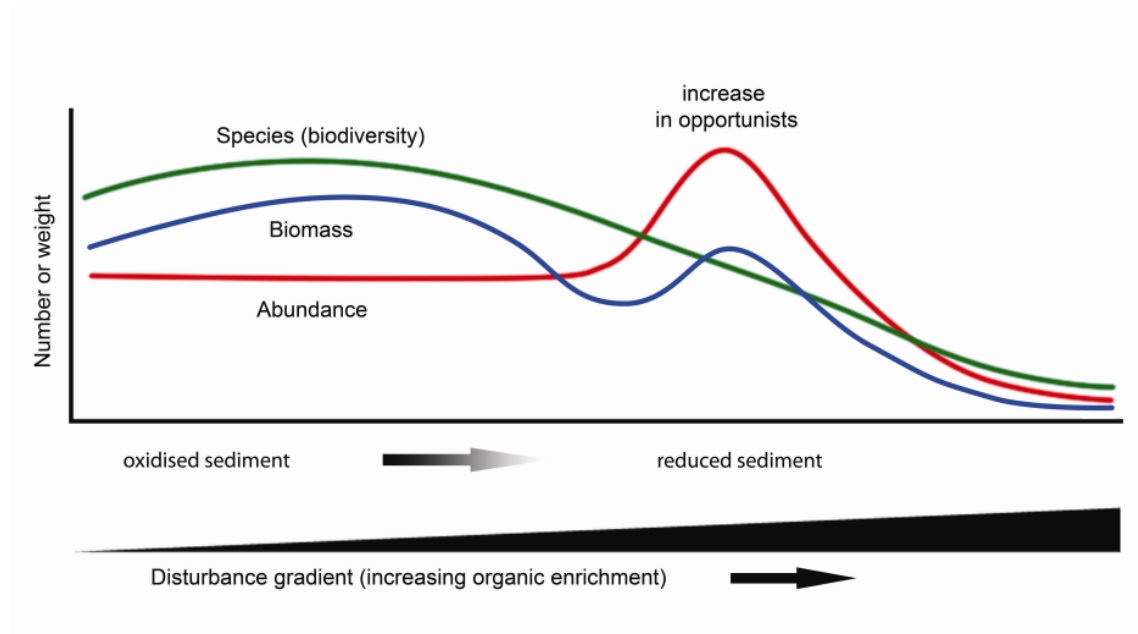


Figure 20. Pearson & Rosenberg relation between the benthic characteristics and a disturbance gradient (organic pollution).

Therefore, the X-axis of **Error! Reference source not found.** can be used as a proxy for the pressure gradient within the NEA-GIG benthic coastal dataset (or the first principle component of the multivariate analysis). Along this gradient, the samples clustered in group E and F can be considered as alternative benchmark sites, because they are characterized by similar diversity characteristics

These diversity characteristics should reflect the status of benthos under least disturbed conditions. The amount of samples in group E and F is high, which allows a good characterization of the natural variability of the benthos within the NEA-GIG region under least disturbed conditions and covered the upper part of the theoretical gradient in benthic characteristics along a disturbance gradient.

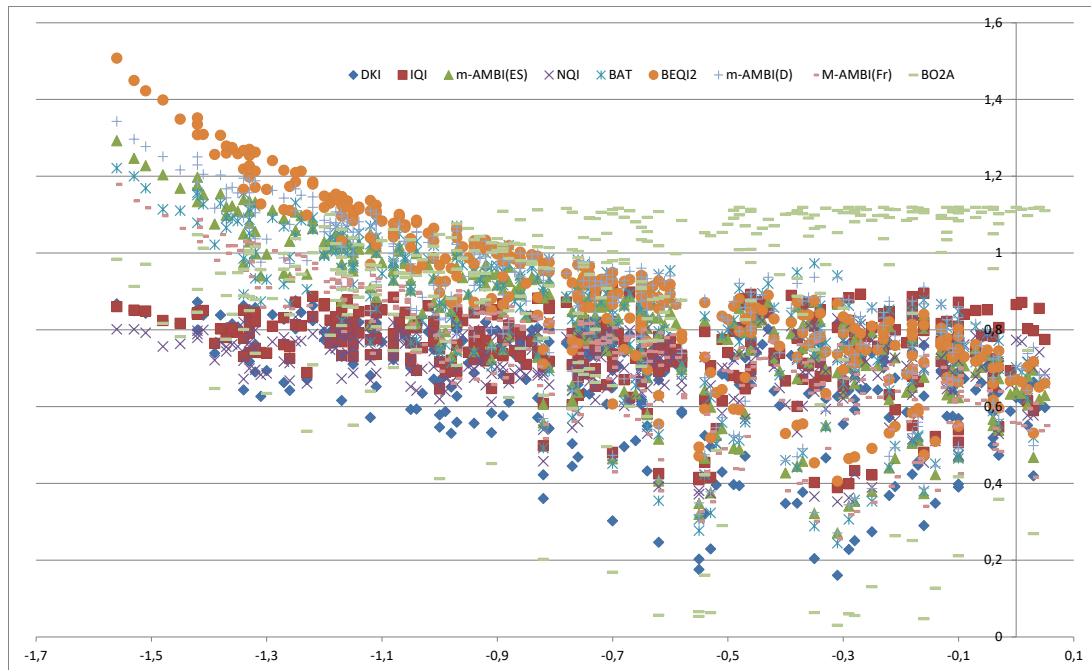


Figure 21. EQR values of the assessment approaches for the benchmark samples

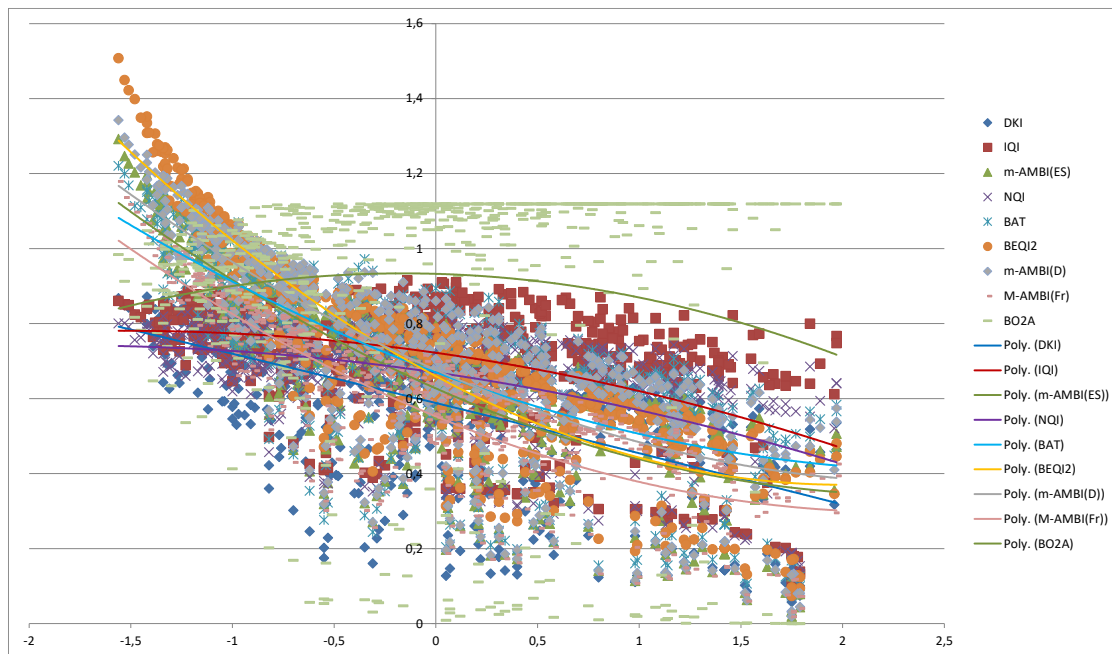


Figure 22. EQR values of the different methods with trend line (2nd order polynomial trend line) along the pressure gradient (X-axis values of MDS).

Figure 21. EQR values of the different methods with trend line (2nd order polynomial trend line) along the pressure gradient (X-axis values of MDS).

The analysis of those benchmark sites (**Error! Reference source not found.**) and g radients (Figure) show that most benthic assessment approaches have a high variability along the gradient, but were more or less in line with each other. The BO2A shows the lowest affinity with this gradient and the highest variability in EQR values for the benchmark sites. The trend line of the BO2A is not in line with the others. Beside this, the M-AMBI approaches, BAT and BEQI2 show the same trend line, whereas the NQI and IQI deviate a little bit from this. They show a more buffered pattern, characterized by less variability at high status, which is related to their algorithm.

Benchmark standardization?

General pattern

In general, significant differences between the different assessment approaches were observed on the benchmark sites within the common dataset (**Error! Reference source not found., Error! Reference source not found.**), except in a few cases (DKI and NQI; IQI and m-AMBI (Fr); BAT and m-AMBI (ES &D); BEQI2 and m-AMBI (D)).

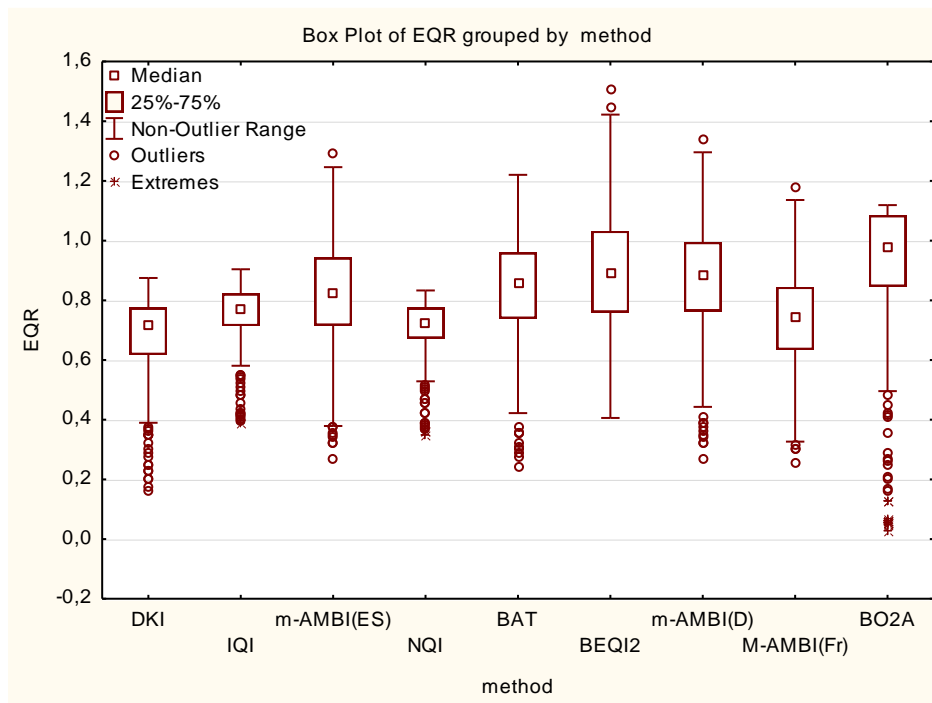


Figure 23. Box-whisker plot of the EQR values at the benchmark sites for the different benthic assessment methods, with indication of the outlier values.

Table 52. Kruskal-Wallis p levels by comparison the EQR values of each approach for the benchmark sites (samples of cluster group E and F)

	DKI	IQI	m-AMBI(E) NQI	BAT	BEQI2	m-AMBI(D)	M-AMBI(Fr)	BO2A
DKI		0,000000	0,000000	1,000000	0,000000	0,000000	0,000304	0,000000
IQI	0,000000		0,000005	0,000059	0,000000	0,000000	1,000000	0,000000
m-AMBI(ES)	0,000000	0,000005		0,000000	1,000000	0,000360	0,021182	0,000000
NQI	1,000000	0,000059	0,000000		0,000000	0,000000	0,013701	0,000000
BAT	0,000000	0,000000	1,000000	0,000000		0,041123	0,829322	0,000000
BEQI2	0,000000	0,000000	0,000360	0,000000	0,041123		1,000000	0,000000
m-AMBI(D)	0,000000	0,000000	0,021182	0,000000	0,829322	1,000000		0,000607
M-AMBI(Fr)	0,000304	1,000000	0,000000	0,013701	0,000000	0,000000		0,000000
BO2A	0,000000	0,000000	0,000000	0,000000	0,000000	0,032067	0,000607	

Benchmark standardization will correct for differences in median EQR values between the Member States benchmark sites obtained by certain assessment approaches. Therefore, we analyze the median EQR values of the Member States (per type [$<30m$ and $>30m$]) benchmark sites for each of the different assessment approaches separately. Those median values will be corrected by the benchmark standardization procedure and this correction will be more obvious for cases where the medians are significantly different.

Benthic assessment approaches at the Member States benchmark sites and comparability results

1) M-AMBI (Germany)

The EQR values at the benchmark sites of Spain, France and Norway are significantly different from the German and UK-type 1 benchmark sites by the m-AMBI (Germany) approach (**Error! Reference source not found., Error! Reference source not found.55**). UK-type 2, the Dutch and UK-type 2 are also significantly different with the French benchmark sites.

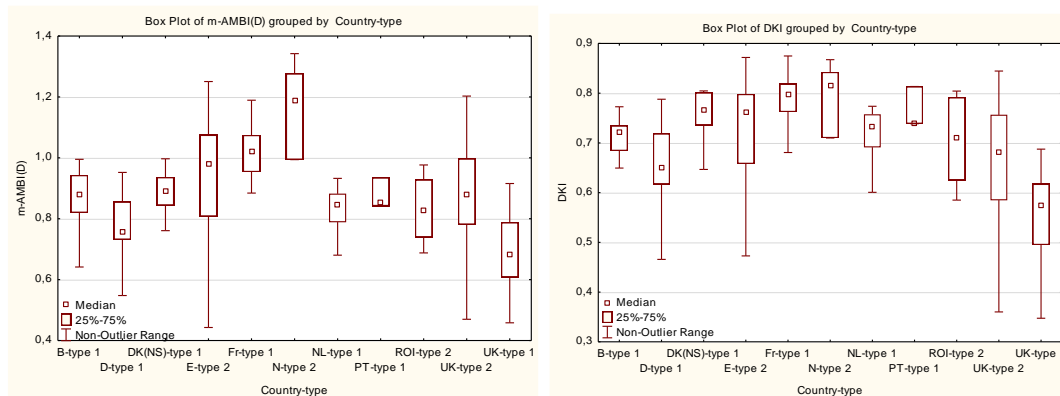


Figure 24. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the m-AMBI (Germany) (left) and the DKI (Denmark) (right).

Table 53. Kruskal Wallis p levels (multiple comparison of mean ranks for all groups) by comparison, the EQR values of each Member state benchmark site with the m-AMBI (Germany) (white fields) and the DKI (Denmark) (grey fields)

DKI	→	B-type 1	D-type 1	DK(NS)-type 1	E-type 2	Fr-type 1	N-type 2	NL-type 1	PT-type 1	ROI-type 2	UK-type 2	UK-type 1
m-AMBI(D)	↓											
B-type 1			1,000	1,000	1,000	0,021	1,000	1,000	1,000	1,000	1,000	0,017
D-type 1		1,000		0,219	0,065	0,000	0,107	1,000	1,000	1,000	1,000	0,692
DK(NS)-type 1		1,000	1,000		1,000	1,000	1,000	1,000	1,000	1,000	0,784	0,000
E-type 2		1,000	0,000	1,000		0,377	1,000	1,000	1,000	1,000	0,219	0,000
Fr-type 1		0,069	0,000	1,000	1,000		1,000	0,173	1,000	1,000	0,000	0,000
N-type 2		0,504	0,000	1,000	1,000	1,000		1,000	1,000	1,000	0,391	0,000
NL-type 1		1,000	1,000	1,000	0,405	0,002	0,064		1,000	1,000	1,000	0,002
PT-type 1		1,000	1,000	1,000	1,000	1,000	1,000	1,000		1,000	1,000	0,207
ROI-type 2		1,000	1,000	1,000	1,000	0,205	0,449	1,000	1,000		1,000	0,156
UK-type 2		1,000	0,013	1,000	1,000	0,000	0,230	1,000	1,000	1,000		0,003
UK-type 1		0,117	1,000	0,350	0,000	0,000	0,000	1,000	1,000	1,000	0,000	

2) DKI

The EQR values at the UK-type1 are significantly different (lower) from most other benchmark sites, except the Portuguese, Irish and German sites (**Error! Reference source not found, Error! Reference source not found.55**). The French and UK-type 1, Belgian and German sites are also significantly different to the DKI benthic assessment approach.

3) M-AMBI of France

The EQR values at the benchmark sites of Spain, France and Norway are significantly different from the German and UK-type 1 benchmark sites by the m-AMBI (France) approach (Figure 25, **Error! Reference source not found.56**). UK-type 2, the Dutch and UK-type 2 are also significantly different from the French benchmark sites. The benchmark sites of the Member States which are significantly different from each other are the same as with the m-AMBI approach of Germany and Spain.

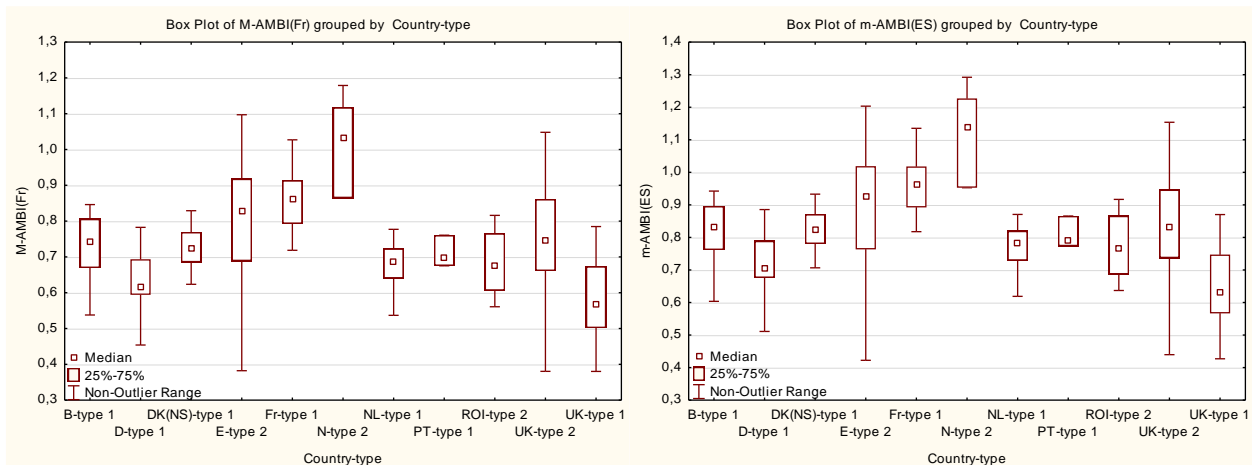


Figure 25. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the m-AMBI (France) (left) and the m-AMBI (Basque Country, Cantabria region) (right).

Table 54. Kruskal-Wallis p levels 9multiple comparison of mean ranks for all groups) by comparison the EQR values of each Member State benchmark sites with the m-AMBI (France) (white fields) and the m-AMBI (Basque Country; Cantabria) (grey fields).

BC/CR	→	B-type 1	D-type 1	DK(NS)-type 1	E-type 2	Fr-type 1	N-type 2	NL-type 1	PT-type 1	ROI-type 2	UK-type 2	UK-type 1
France	↓											
			0,840	1,000	1,000	0,102	0,471	1,000	1,000	1,000	1,000	0,142
		0,687		1,000	0,000	0,000	0,000	1,000	1,000	1,000	0,002	1,000
		1,000	1,000		1,000	1,000	1,000	1,000	1,000	1,000	1,000	0,777
		1,000	0,000	1,000		1,000	1,000	0,138	1,000	1,000	1,000	0,000
		0,128	0,000	0,729	1,000		1,000	0,001	1,000	0,163	0,002	0,000
		0,439	0,000	0,894	1,000	1,000		0,026	1,000	0,292	0,309	0,000
		1,000	1,000	1,000	0,066	0,000	0,014		1,000	1,000	1,000	1,000
		1,000	1,000	1,000	1,000	1,000	1,000	1,000		1,000	1,000	1,000
		1,000	1,000	1,000	1,000	0,141	0,219	1,000	1,000		1,000	1,000
		1,000	0,001	1,000	1,000	0,005	0,384	1,000	1,000	1,000		0,000
		0,190	1,000	1,000	0,000	0,000	0,000	1,000	1,000	1,000	0,000	

4) M-AMBI of Spain (Basque Country, Cantabria region)

The EQR values at the benchmark sites of Spain, France and Norway are significantly different from the German and UK-type 1 benchmark sites by the m-AMBI (Basque Country, Cantabria) approach (Figure 25, **Error! Reference source not found.**). UK-type 2, the Dutch and UK-type 2 are also significantly different from the French benchmark sites. The benchmark sites of the Member States which were significantly different from each other are the same as with the m-AMBI approach of Germany and France.

5) BEQI2 of the Netherlands

The EQR values at the benchmark sites of Spain, France and Norway are significantly different from the German, Dutch and UK-type 1 benchmark sites by the BEQI2 approach (Figure 26, **Error! Reference source not found.57**). The EQR values of the benchmark site of UK-type 2 and UK-type 1 and Germany are also significantly different. The benchmark sites of the Member States which were significantly different from each other are the same as with the m-AMBI approach of Germany, Spain and France.

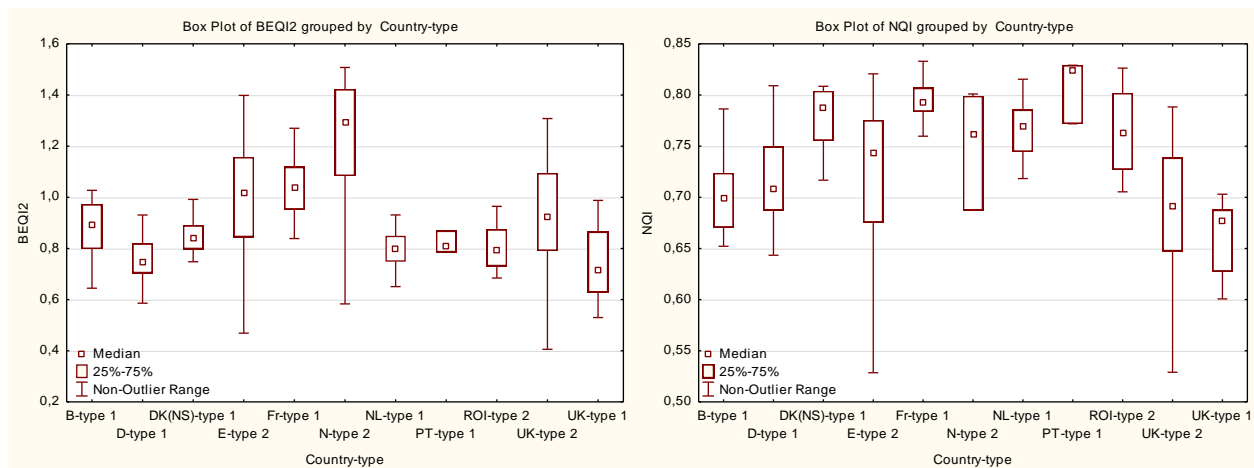


Figure 26. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the BEQI2 (the Netherlands) (left) and the NQI (Norway) (right).

Table 55. Kruskal-Wallis p values (multiple comparison mean ranks for all groups) by comparison the EQR values of each Member State benchmark sites with the BEQI2 (the Netherlands) white fields and the NQI (Norway) (grey fields).

NQI	→	B-type 1	D-type 1	DK(NS)-type 1	E-type 2	Fr-type 1	N-type 2	NL-type 1	PT-type 1	ROI-type 2	UK-type 2	UK-type 1
BEQI2	↓											
B-type 1			1,000	0,047	1,000	0,000	1,000	0,023	0,260	0,531	1,000	1,000
D-type 1	0,609			0,221	1,000	0,000	1,000	0,123	0,727	1,000	1,000	0,198
DK(NS)-type 1	1,000	1,000			1,000	1,000	1,000	1,000	1,000	1,000	0,001	0,000
E-type 2	1,000	0,000	1,000			0,001	1,000	0,916	1,000	1,000	0,199	0,012
Fr-type 1	0,190	0,000	0,449	1,000			1,000	1,000	1,000	1,000	0,000	0,000
N-type 2	0,324	0,000	0,386	1,000	1,000			1,000	1,000	1,000	1,000	0,124
NL-type 1	1,000	1,000	1,000	0,014	0,000	0,004			1,000	1,000	0,000	0,000
PT-type 1	1,000	1,000	1,000	1,000	1,000	1,000	1,000			1,000	0,085	0,012
ROI-type 2	1,000	1,000	1,000	0,572	0,064	0,071	1,000	1,000			0,065	0,005
UK-type 2	1,000	0,000	1,000	1,000	1,000	0,089	0,641	0,166	1,000	1,000		1,000
UK-type 1	0,749	1,000	1,000	1,000	0,000	0,000	0,000	1,000	1,000	1,000	0,000	

6) NQI of Norway

The EQR values at the benchmark sites of UK-type 1 are significantly (lower) different from most other benchmark sites by the NQI, except for the Belgian, German, Norwegian and UK-type 2 benchmark sites (Figure 26, **Error! Reference source not found.57**). The French benchmark sites are also significantly different from many other sites (Belgium, Germany, Spain, UK-type 1 and UK-type 2). There are also significant differences between the Belgian and Danish and Dutch benchmark sites with the NQI approach.

7) BAT of Portugal

The EQR values at the benchmark sites of UK-type 1 are significant different from Spain, France, Norway and UK-type 2 benchmark sites by the BAT benthic assessment approach (**Error! Reference source not found., Error! Reference source not found.58**). The German benchmark sites are significant different with Spain, France, Norway and UK type 2. Significant difference are also observed between the Dutch and French and Belgian and French benchmark sites and the French and the UK-type 2 sites.

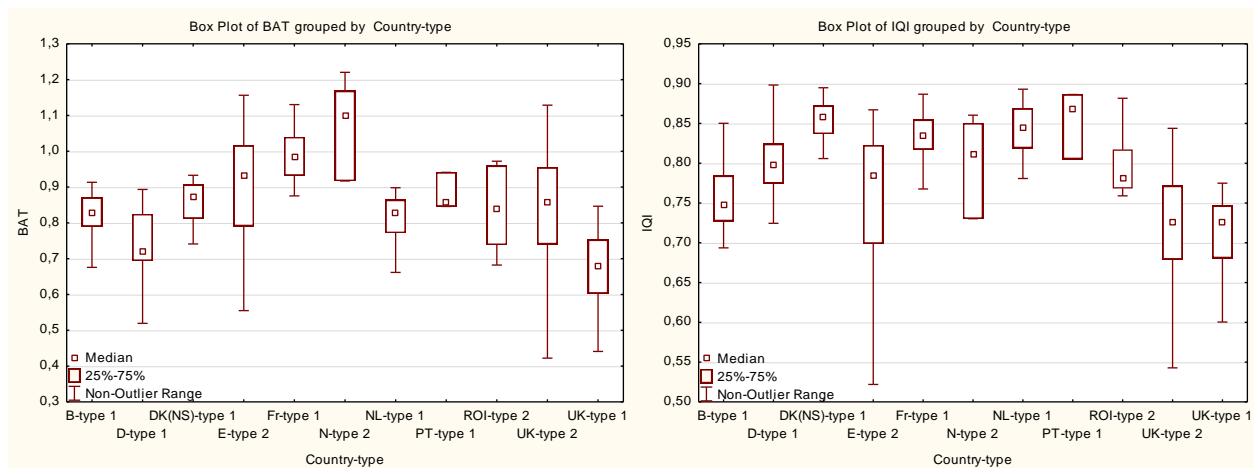


Figure 27. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the BAT (Portugal) (left) and the IQI (United Kingdom, Ireland) (right).

Table 56. Kruskal-Wallis p values (multiple comparison of mean ranks for all groups) by comparison the EQR values of each member state benchmark sites with the BAT (Portugal) (white fields) and the IQI (UK and Ireland) (grey fields).

IQI	→	B-type 1	D-type 1	DK(NS)-type 1	E-type 2	Fr-type 1	N-type 2	NL-type 1	PT-type 1	ROI-type 2	UK-type 2	UK-type 1
BAT	↓											
B-type 1			1,000	0,011	1,000	0,002	1,000	0,001	1,000	1,000	1,000	1,000
D-type 1		1,000		1,000	1,000	1,000	1,000	0,539	1,000	1,000	0,000	0,002
DK(NS)-type 1		1,000	1,000		0,034	1,000	1,000	1,000	1,000	1,000	0,000	0,000
E-type 2		1,000	0,000	1,000		0,004	1,000	0,003	1,000	1,000	0,148	0,394
Fr-type 1		0,004	0,000	1,000	1,000		1,000	1,000	1,000	1,000	0,000	0,000
N-type 2		0,321	0,000	1,000	1,000	1,000		1,000	1,000	1,000	0,622	0,545
NL-type 1		1,000	1,000	1,000	0,791	0,001	0,181		1,000	1,000	0,000	0,000
PT-type 1		1,000	1,000	1,000	1,000	1,000	1,000	1,000		1,000	0,104	0,079
ROI-type 2		1,000	1,000	1,000	1,000	1,000	1,000	1,000	1,000		0,400	0,342
UK-type 2		1,000	0,013	1,000	0,893	0,000	0,307	1,000	1,000	1,000		1,000
UK-type 1		0,291	1,000	0,189	0,000	0,000	0,000	0,590	1,000	0,755	0,000	

8) IQI of UK/Ireland

The classification of the benchmark sites of the different Member States by the IQI leads also to some significant differences (**Error! Reference source not found., Error! Reference source not found.**). The Danish sites are significantly different from the Belgian, Spanish, UK-type 2 and UK-type 1 sites. The French sites are significantly different from the Belgian, Spanish, Dutch, UK-type 2 and UK-type 1 sites. The Dutch sites are also significant different from the Spanish, UK-type 2 and UK-type 1 sites.

9) BO2A of Spain (Andalusia region)

From Andalusia region, no benthic data was included in the common dataset. Therefore, no benchmark sites were delimited for this region of this Member State. The median values of the benchmark sites of the different Member States, evaluated with the BO2A are also different in some cases (**Error! Reference source not found.**27). The sites of UK-type 2, Spain, France and Norway have lower values than the others.

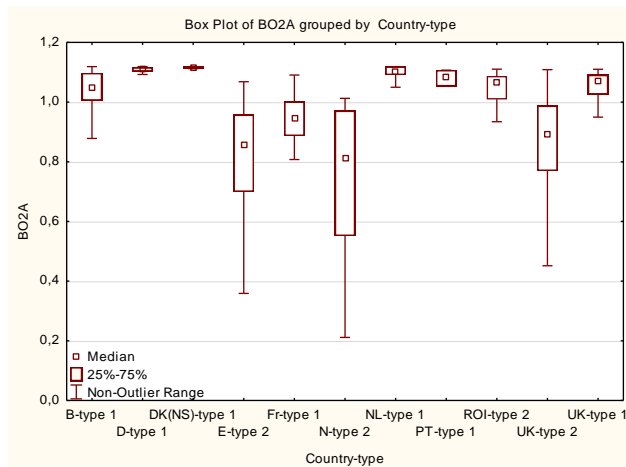


Figure 28. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the BO2A (Spain, Andalusia region).

Results of the regression comparison

For all the intercalibration comparisons, the benthic assessment approaches fulfill the criteria ($R^2 > \frac{1}{2} \max R^2$) of the regression comparison (Table 19), except BO2A. The BO2A of Spain (Andalusia) shows the lowest correlation with the pseudo-common metric. For the IQI and the NQI, the samples were less equally spread over the linear regression line (dominance in upper part) in comparison to the other approaches, as was the case in the analyses on the theoretical behavior of the benthic assessment approaches.

Table 57. Summary of the correlation coefficient (R2) of each approach with the common metric for the different intercalibration comparisons. Values outside the criteria were put in red.

Method	Subtraction standardization		Division standardization	
	No sub-region	Sub-region	No sub-region	Sub-region
Denmark	0.9553	0.9549	0.9536	0.9566
UK/ROI	0.8402	0.8692	0.8267	0.8105
Spain (BC, CR)	0.8864	0.9406	0.8887	0.9261
Norway	0.8965	0.9148	0.8911	0.9141
Portugal	0.9477	0.9671	0.9465	0.9621
The Netherlands	0.7869	0.8762	0.7923	0.8503
Germany	0.9121	0.9569	0.9227	0.9475
France	0.8514	0.9224	0.8542	0.9026
Spain (AC)	0.3599	0.4315	0.3546	0.4508

Comparability criteria

Subtraction benchmark standardization

1) no sub-regions (deep/shallow areas) within NEA 1/26 type.

The boundary bias (<0.25) is too high for the BO2A (**Error! Reference source not found.**). Denmark, France and Germany are more stringent than the other approaches, especially for the good/moderate boundary. The class differences (<0.5 class) is too high for the BO2A and around criteria level for the DKI.

Table 58. Summary of the boundary bias and class differences analyses following the subtraction benchmark standardization, no discrimination of sub regions

Boundary bias	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (AC)
Max	1,000	1,000	1,292	1,000	1,220	1,508	1,342	1,179	1,119
H/G	0,800	0,750	0,770	0,720	0,790	0,780	0,850	0,770	0,830
G/M	0,600	0,640	0,530	0,630	0,580	0,580	0,700	0,530	0,500
M/P	0,400	0,440	0,380	0,400	0,440	0,380	0,400	0,380	0,400
P/B	0,200	0,240	0,200	0,200	0,270	0,180	0,200	0,200	0,200
H/G bias_CW	0,602	-0,149	0,021	-0,043	-0,022	-0,035	0,217	0,374	-1,355
G/M bias_CW	0,351	-0,172	-0,151	0,061	-0,171	0,042	0,288	0,227	-0,906
	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (A)
Absolute Class Difference	0,512	0,469	0,339	0,413	0,342	0,369	0,411	0,402	0,789

2) Sub-regions (deep/shallow areas) within NEA 1/26 type

The boundary bias (<0.25) is in this analysis is too high for BO2A and slightly too high for the m-AMBI (Basque Country, Cantabria) (**Error! Reference source not found.**). The DKI, m-AMBI (Germany & France) are more stringent for the good/moderate boundary and the high/good boundary compared to the others. The class difference (<0.5 class) is too high for the BO2A and at criteria level for the DKI. The m-AMBI can meet the criteria by elevating the good/moderate boundary value to 0.56.

Table 59. Summary of the boundary bias and class differences analyses following the subtraction benchmark standardization, with discrimination of the sub-regions.

Boundary bias	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (AC)
Max	1,000	1,000	1,292	1,000	1,220	1,508	1,342	1,179	1,119
H/G	0,800	0,750	0,770	0,720	0,790	0,780	0,850	0,770	0,830
G/M	0,600	0,640	0,530	0,630	0,580	0,580	0,700	0,530	0,500
M/P	0,400	0,440	0,380	0,400	0,440	0,380	0,400	0,380	0,400
P/B	0,200	0,240	0,200	0,200	0,270	0,180	0,200	0,200	0,200
H/G bias_CW	0,546	-0,129	-0,073	0,082	-0,004	0,010	0,275	0,454	-1,310
G/M bias_CW	0,331	-0,132	-0,313	0,156	-0,085	0,158	0,340	0,363	-1,106
	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (A)
Absolute Class Difference	0,510	0,449	0,326	0,394	0,327	0,356	0,389	0,387	0,749

Division benchmark standardization

1) No sub-regions (deep/shallow areas) within NEA 1/26 type

The boundary bias (<0.25) is in this analysis is too high for BO2A (**Error! Reference source not found.**). The DKI is more stringent for the good/moderate and high/good boundary, in France for the high/good and Germany for the good/moderate boundary. The class difference (<0.5 class) is too high for the BO2A approach and at criteria level for the DKI.

Table 60. Summary of the boundary bias and class differences analyses following the division benchmark standardization, no discrimination of the sub-regions.

Boundary bias	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (AC)
Max	1,000	1,000	1,292	1,000	1,220	1,508	1,342	1,179	1,119
H/G	0,800	0,750	0,770	0,720	0,790	0,780	0,850	0,770	0,830
G/M	0,600	0,640	0,530	0,630	0,580	0,580	0,700	0,530	0,500
M/P	0,400	0,440	0,380	0,400	0,440	0,380	0,400	0,380	0,400
P/B	0,200	0,240	0,200	0,200	0,270	0,180	0,200	0,200	0,200
H/G bias_CW	0,585	-0,143	0,016	-0,043	-0,033	-0,032	0,204	0,371	-1,242
G/M bias_CW	0,339	-0,158	-0,156	0,059	-0,186	0,046	0,282	0,222	-0,836
	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (A)
Absolute Class Difference	0,512	0,469	0,339	0,413	0,342	0,369	0,411	0,402	0,789

2) Sub-regions (deep/shallow areas) within NEA 1/26 type

The boundary bias (<0.25) in this analysis is too high for BO2A and slightly too high for the m-AMBI(BC, CR) (The BEQI assessment approach meet the comparability criteria in comparison with the other approaches. Further boundary adjustment cannot be suggested, as this is a comparability check on higher level than sample level; in most assessment approaches, their boundaries were based on a sample level evaluation. Besides this, the BEQI is comparable with all methods applied in sub-region A (very shallow) type - all Belgian coastal waters belong to sub-region A.

Table 22). The bias for DKI, Germany and France is more stringent for the good/moderate and high/good boundary. The class difference (<0.5 class) is too high for the BO2A approach and at criteria level for the DKI. The m-AMBI can meet the criteria by elevating the good/moderate boundary value to 0.55.

Table 63. Summary of the boundary bias and class differences analyses following the division benchmark standardization, with discrimination of the sub-regions.

Boundary bias	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (AC)
Max	1,000	1,000	1,292	1,000	1,220	1,508	1,342	1,179	1,119
H/G	0,800	0,750	0,770	0,720	0,790	0,780	0,850	0,770	0,830
G/M	0,600	0,640	0,530	0,630	0,580	0,580	0,700	0,530	0,500
M/P	0,400	0,440	0,380	0,400	0,440	0,380	0,400	0,380	0,400
P/B	0,200	0,240	0,200	0,200	0,270	0,180	0,200	0,200	0,200
H/G bias_CW	0,539	-0,210	-0,069	0,066	-0,008	0,000	0,263	0,448	-1,367
G/M bias_CW	0,327	-0,209	-0,290	0,137	-0,131	0,081	0,305	0,278	-0,978
	Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (A)
Absolute Class Difference	0,510	0,449	0,326	0,394	0,327	0,356	0,389	0,387	0,749

Conclusion

The intercalibration of the benthic assessment approaches within the NEA-GIG region can be executed following the intercalibration guidelines. As shown in the analysis, the benthic assessment approaches are very comparable (some after a small adaptation of their boundaries) and meet the intercalibration criteria, except for the BO2A. The subtraction and

division standardization delivers the same results regarding the acceptability of the criteria. The subtraction standardization only delivers slightly higher values compared to the division standardization

The BO2A does not meet the criteria of boundary bias and class difference in any intercalibration comparison option. The adaptation of the boundaries to meet the criteria is rather impossible for this approach, because the tests to change the boundary levels of the BO2A do not lead to any situation that meets the criteria. They even influenced the criteria levels of the other approaches, mostly in a negative way. The application of the BO2A on this common NEA-GIG dataset seems to be more complicated and different from the results of the own intercalibration analyses of the Andalusia region (see separate document (2011-12-16technical_report_NEA_CW_invertebrates_ES(AN)_Dec2011)).

The DKI and m-AMBI (Germany & France) show in all intercalibration comparison options a more stringent evaluation than the other approaches. Therefore, those boundary values can even be lowered to be more comparable with the other methods, but this is not required.

The m-AMBI (Basque Country and Cantabria) shows in the intercalibration comparison option, where sub-regions are distinguished, a slightly too high boundary bias. This approach scored not in correspondence with the other approaches (IQI, NQI) for samples typical for less shallow areas. This approach can easily meet the criteria, as the good/moderate boundary is slightly increased (+0.02 or 0.03).

All other benthic assessment approaches (BAT, BEQI2, IQI) meet the comparability criteria.

Based on the analyses and the experience with the data and the assessment approaches, the intercalibration comparison with the division benchmark standardization and no discrimination of the sub-regions should be most appropriate. This because, the approaches show similar trend lines, but there are differences between them along the pressure gradient (some of them vanish). Besides this, there was no hard evidence to discriminate sub-regions, and the reference settings for these soft sediment habitats are similarly set by the Member States for this type.

Benchmark selection based on expert judgment

The comparison is executed based on certain conditions, but the selection of those conditions has its effect on the boundary bias values. In section on benchmark standardization, the results of the biotic benchmark are shown, which reveals no fail in the boundary bias criteria for IQI, whereas on expert judgment it does (**Error! Reference source not found.**64). Also the inclusion or exclusion of sub-regions, regardless of the benchmarking, has an effect on the boundary bias, especially for Spain. When no sub-region is recognized, no boundary harmonization is necessary, whereas this is necessary when sub-regions are recognized. This could be related to inappropriate reference values for this sub-region type in Spain, but this seems not to be the case (see below).

Further, the inclusion or exclusion of a method has its consequence on the boundary bias values, which became slightly lower. This happens because the BO2A assessment approach is not comparable with the others. This aspect is worth mentioning, because adding or changing a method has consequences on the obtained comparability results.

This were all intermediate comparability analyses to explore the intercalibration and to move towards the selection of the comparison most in line with the intercalibration guidelines.

Table 61. Summary of the boundary bias for the H/G and G/M following different conditions regarding discrimination of subregion or not or including/excluding certain methods.

Boundary bias H/G					Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (AC)	Belgium
benchmarking	subtracti on/divisi on	ub-region	level											
expert	No BO2A	division	no	sample	0,507	-0,276	-0,017	-0,127	-0,056	-0,072	0,110	0,326		
expert	No BO2A	division	yes	sample	0,495	-0,330	-0,303	-0,138	0,007	0,013	0,294	0,479		
expert	all metho	division	yes	sample	0,599	-0,237	-0,255	-0,068	0,108	0,163	0,442	0,552	-1,426	
expert	and BEQI	division	no	higher	0,513098	-0,28218	-0,056337344	-0,13505	-0,08983	-0,10408	0,007133	0,253903		0,448
expert	and BEQI	division	yes	higher	0,522862	-0,25393	-0,336498172	-0,13716	-0,05608	0,005973	0,190226	0,422585		0,229

Boundary bias G/M					Denmark	UK/ROI	Spain (BC, C)	Norway	Portugal	Netherlands	Germany	France	Spain (AC)	Belgium
benchmarking	subtracti on/divisi on	ub-region	level											
expert	No BO2A	division	no	sample	0,291	-0,303	-0,213	0,008	-0,221	-0,124	0,238	0,081		
expert	No BO2A	division	yes	sample	0,304	-0,407	-0,684	0,000	-0,110	0,024	0,310	0,176		
expert	all metho	division	yes	sample	0,381	-0,270	-0,593	0,058	-0,029	0,185	0,372	0,399	-1,028	
expert	and BEQI	division	no	higher	0,250348	-0,20754	0,118380469	-0,03807	-0,29724	-0,17331	0,173654	0,011773		0,659
expert	and BEQI	division	yes	higher	0,289911	-0,23385	-0,508514206	-0,04156	-0,19626	0,003356	0,243523	0,175791		0,575

Test for changing the reference values of Spain for sub-type 2.

If the m-AMBI reference values for the deeper samples (AMBI: 1, Diversity: 5.7, Richness: 130) are applied, it seems that they were too high. This is because there is no sample in high status for this sub-type in the common dataset, which is not true (some stations has no pressures, such as Norway). Therefore, Spain became too stringent in their assessment, whereas the other countries of type 2 does not meet the boundary bias criteria at all. Spain will thus therefore accept the boundary harmonisation (0.63 for G/M).

	Denmark	UK/ROI	Spain (BQ, CQ)	Norway	Portugal	Netherlands	Germany	France
Max	1,000	1,000	1,000	1,000	1,130	1,270	1,189	1,027
H/G	0,800	0,750	0,770	0,720	0,790	0,780	0,850	0,770
G/M	0,600	0,640	0,530	0,630	0,580	0,580	0,700	0,530
M/P	0,400	0,440	0,380	0,400	0,440	0,380	0,400	0,380
P/B	0,200	0,240	0,200	0,200	0,270	0,180	0,200	0,200
H/G bias_CW	0,279	-0,463	0,824	-0,263	-0,142	-0,156	-0,072	0,240
G/M bias_CW	0,149	-0,613	0,720	-0,275	-0,280	-0,244	0,147	0,025

Annex 3. Common type NEA 3/4. Pressures

Table 62. Pressure info per location for Germany (** DIN=arithmetic mean of DIN winter means (Nov-Feb)-(from nearest monitoring point to MZB station).

Dataset name	Name of German authority responsible for the Data	German station Name	Water body type NEA 3 or 4	HABITAT			PRESSURES ²					Eutrophication (DIN)			Fishery: ICES fishery map (indirect linking)			Benchmark sites
				habitat/ecotope	Depth	Sediment	Pressure	quantitative	qualitative	expert judgment	remarks	Eutro/high (DIN)	Eutro/medium (DIN)	Eutro/low (DIN)	Fishery/high	Fishery/medium	Fishery/low	
1 German Wadden Sea	NLWKN	Bork_MZB_8	3	subtidal finesand	>6m	Finesand	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2011	1,25			high			
2 German Wadden Sea	NLWKN	AuWe_MZB_1	3	subtidal finesand	>6m	Finesand	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2011		0,47			medium		
3 German Wadden Sea	Bfg	Weser-4	3	subtidal sand	9m	Sand	eutrophication and fisheries		fisheries	yes					high			
4 German Wadden Sea	Bfg	Elbe-4	3	subtidal sand	12-15m	Sand	eutrophication and fisheries		fisheries	yes					high			
5 German Wadden Sea	Bfg	Elbe-5	3	subtidal sand	12-15m	Sand	eutrophication and fisheries		fisheries	yes					high			
6 German Wadden Sea	Bfg	Ems-4	3	subtidal sand with mud	9m	Sand with Mud	eutrophication and fisheries		fisheries	yes					high			
8 German Wadden Sea	NLWKN	Nney_MZB_6	4	litoral mud	intertidal	mud	eutrophication and fisheries	DIN	fisheries	yes	DIN 2007-2013	0,82			high			
9 German Wadden Sea	NLWKN	Nney_MZB_7	4	litoral mud	intertidal	mud	eutrophication and fisheries	DIN	fisheries	yes	DIN 2007-2013	0,82			high			
10 German Wadden Sea	NLWKN	Nney_MZB_5	4	litoral mud	intertidal	mud	eutrophication and fisheries	DIN	fisheries	yes	DIN 2007-2013	0,82					low	
11 German Wadden Sea	NLWKN	Nney_MZB_1	4	litoral sand	intertidal	sand	eutrophication and fisheries	DIN	fisheries	yes	DIN 2007-2013	0,82			high			
12 German Wadden Sea	NLWKN	Nney_MZB_2***	4	litoral sand	intertidal	sand	eutrophication and fisheries	DIN	fisheries	yes	DIN 2007-2013	0,82					low	yes
13 German Wadden Sea	NLWKN	Nney_MZB_3	4	litoral muddy sand	intertidal	muddy sand	eutrophication and fisheries	DIN	fisheries	yes	DIN 2007-2013	0,82			high			
14 German Wadden Sea	NLWKN	Nney_MZB_8***	4	litoral muddy sand	intertidal	muddy sand	eutrophication and fisheries	DIN	fisheries	yes	DIN 2007-2013	0,82					low	yes
15 German Wadden Sea	NLWKN	WuKu_MZB_6	4	litoral muddy sand	intertidal	muddy sand	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2010	0,96			high			
16 German Wadden Sea	NLWKN	Bork_MZB_4	4	subtidal mud	<6m	mud	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2011	1,25			high			
17 German Wadden Sea	NLWKN	AuWe_MZB_3	4	litoral sand	intertidal	sand	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2011		0,47		high			
18 German Wadden Sea	NLWKN	WuKu_MZB_10	4	litoral finesand	intertidal	finesand	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2010	0,96			high			
19 German Wadden Sea	BSU HH	HH T1	4	litoral sandy to muddy	intertidal	sandy to muddy	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2011		0,58		high			
20 German Wadden Sea	BSU HH	HH T2	4	litoral sandy to muddy	intertidal	sandy to muddy	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2011		0,58		high			
21 German Wadden Sea	BSU HH	HH T3	4	litoral sandy to muddy	intertidal	sandy to muddy	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2011		0,58		high			
22 German Wadden Sea	BSU HH	HH T4	4	litoral sandy to muddy	intertidal	sandy to muddy	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2011		0,58		high			
23 German Wadden Sea	BSU HH	HH T5	4	litoral sandy to muddy	intertidal	sandy to muddy	eutrophication and fisheries	DIN	fisheries	yes	DIN 2001-2011		0,58		high			

Table 63. Pressure info for the Dutch Wadden Sea

Column header	Specifications
Dataset name	Dutch Wadden Sea
Data owner	Rijkswaterstaat
Station names	Balgzand (Western Dutch Wadden Sea; 3 transects: B, C, J) Piet Scheveplaat (Eastern Dutch Wadden Sea 3 transects; 600, 601, 602)
NEA type	3-4 (Gert, can this be discriminated for the Western and Eastern part of the Dutch Wadden Sea?)
Habitat/ecotope	Litoral muddy sand
Depth	Intertidal
Sediment type	Muddy sand
Common pressure types	Eutrophication, fisheries
Pressures characterization method	Eutrophication: using NH ₄ +NO ₂ results from QSR report Wadden Sea 2009, Thematic report No. 9 Eutrophication, Table 5. Fisheries: using QSR report Wadden, Thematic report No. 3.3 Fisheries, Figure 3.3.6 (shrimp fisheries), Figure 3.3.3 (Mussel seed fisheries).
Pressure data period	Eutrophication: 2000-2006 (QSR Wadden Sea) Fisheries: depends on fishing type, around 2000-2007.
Pressure quantification	Eutrophication: in the Western Dutch Wadden Sea, the assessment value (period 2000-2006) of 8.2 uM NH ₄ +NO ₂ is just below the "problem condition limit" of 8.3 uM. In the Eastern Dutch Wadden Sea, the assessment value of 16.8 uM (period 2000-2006) exceed the problem condition limit of 10.2 uM.

	<p>In conclusion, there is significant eutrophication in the Dutch Wadden Sea, especially in the Eastern part. Note however that for benthos, some amount of eutrophication is probably not a problem, because it delivers additional food for filter feeders.</p> <p>Fisheries:</p> <ol style="list-style-type: none"> 1. Shrimp fisheries only occurs in subtidal parts, mainly in the Western Wadden Sea. No shrimp fisheries occur in the intertidal areas because these areas are too shallow for fishing boats. 2. Mussel seed fisheries mainly occur in the subtidal areas in the Western Wadden Sea, and not in the intertidal parts. 3. Since January 2005 mechanical cockle fishery in the Dutch part of the Wadden Sea is not allowed any longer. Only manual cockle fishery is still allowed with a maximum yearly catch of 5% of the cockle stock. The fished amounts were between 0.1 and 1.5 % of the stock. So there is some manual cockle fisheries in the intertidal parts of the Wadden Sea, but this pressure is probably relatively low. <p>In conclusion, the fisheries pressure in the intertidal parts of the Dutch Wadden Sea is low. In the subtidal parts, especially of the Western Dutch Wadden Sea, the fishing pressure is relatively high.</p>
Benchmark sites	Yes, Piet Scheveplaat.

List of abbreviations and definitions

Key Terms:

Assessment method: The biological assessment for a specific biological quality element, applied as a classification tool, the results of which can be expressed as EQR.

Biological Quality Element (BQE): Particular characteristic group of animals or plants present in an aquatic ecosystem that is specifically listed in Annex V of the Water Framework Directive for the definition of the ecological status of a water body (for example phytoplankton or benthic invertebrate fauna).

Class boundary: The Ecological Quality Ratio value representing the threshold between two quality classes.

Common Intercalibration type: A type of surface water differentiated by geographical, geological, morphological factors (according to WFD Annex II) shared by at least two Member States in a GIG.

Compliance criteria: List of criteria evaluating whether assessment methods are meeting the requirements of the Water Framework Directive.

Ecological Quality Ratio (EQR): Calculated from the ratio observed value/reference value for a given body of surface water. The ratio shall be represented as a numerical value between zero and one, with high ecological status represented by values close to one and bad ecological status by values close to zero.

Geographic Intercalibration Group (GIG): Organizational unit for the intercalibration consisting of a group of Member States sharing a set of common intercalibration types.

Intercalibration: An exercise facilitated by the Commission to ensure that the high/good and good/moderate class boundaries are consistent with Annex V Section 1.2 of the Water Framework Directive and comparable between Member States.

IC Option: Option to intercalibrate (IC) different national assessment methods.

Method Acceptance Criteria: List of criteria evaluating whether assessment methods can be included in the intercalibration exercise.

Pressure: Human activities such as organic pollution, nutrient loading or hydromorphological modification that have the potential to have adverse effects on the water environment.

Reference/Benchmark sites: Reference sites meet international screening criteria for undisturbed conditions. Benchmark sites meet a similar (low) level of impairment associated with the least disturbed or best commonly available conditions.

Water Framework Directive: Directive 2000/60/EC establishing a framework for Community action in the field of water policy.

Abbreviations:

A: Andalusia region

BE: Belgium

BC: Basque Country

C: Cantabria region

DE: Germany

EG: Ecological group

EQR: Ecological Quality Ratio

ES: Spain

FR: France

IE: Ireland

GIG: Geographic Intercalibration Group

IC: Intercalibration

MS: Member State

NL: Netherlands

PCA: Principal Correspondence analyses

PT: Portugal

RC: Reference conditions

Se: Sweden

UK: United Kingdom

WFD: Water Framework Directive

List of figures

Figure 1. Correlation plot with trend line (polynomial 2 nd order) between the different assessment approaches and Cu(mg/kg).	28
Figure 2. Box-Whisker plot of the EQR values of the benthic assessment approaches for the classification of the Garroch head benchmark sites.	29
Figure 3. Correlation plot with trend line (polynomial 2 nd order) between the different assessment approaches and Cu(mg/kg) for the set of pooled samples.	30
Figure 4. RAT method-pressure relationships	31
Figure 5. Changes in EQR values on fictive datasets, to show the metric dependency and behavioral response of the algorithm.	34
Figure 6. MDS plot of the intercalibration data with indication of the Member States (colored symbols) and the cluster groups (slice at similarity 11)	38
Figure 7. MDS plot of the univariate variables (inclusive abundance), with indication of the cluster groups at slice 65 (upper figure) and the behaviour of the dataset of the different Member Staes (center figure).	39
Figure 8. MDS plot of the univariate parameters (exclusive abundance and indication of the cluster groups (slice 75) (upper figure) and the behavior of the datasets of the different Member states (lower figure).	40
Figure 9. Box-whisker plot median, percentile values and no outlier range of the EQR values at the Member states benchmark sites with the national methods for the two subtypes (shallowness and depth).....	45
Figure 10. A-B Time trends of the state of the parameters oxygen and suspended matter. C-D state impact correlations for oxygen concentration and suspended matter with BEQI EQRs. waterbody ecotopes Westerschelde mesohaline-intertidal (WS_MI) and Dollard mesohaline-intertidal (DOI_MI), respectively.	67
Figure 11. M-AMBi values at each sampling data in the BACI design monitoring for sandextraction at Dangaster	68
Figure 12. MDS of the cluster groups 9(slice 31 Bray Curtis similarity), which result in 9n groups are coded alphabetically (a-i)	72
Figure 13. MDS with indication of the habitat types.....	73
Figure 14. Box-whisker plot of the assessment of the Dutch and German benchmark sites with each benthic assessment approach.	74
Figure 15. Scatter plot of EQR values of Germany and Netherlands, with linear regression line.	75
Figure 16. Correlation between Norwegian method and pressures	87
Figure 17. Correlation between EQR (Infaunal Quality index) and principal component (PCA1) of Cu and organic carbon data (sewage sludge disposal pressure).....	88
Figure 18. Correlation between EQR (Infaunal Quality index) and principal Component assessment (PCA1) of Cu, Cr and silt/clay data (Mine waste pressure).	88
Figure 19. Ecological status as assessed by the Infaunal Quality IQIvIV at distance from fish farm pressure (Loch Duich, 2003)	89
Figure 20. Pearson & Rosenberg relation between the benthic characteristics and a disturbance gradient (organic pollution).	120
Figure 21. EQR values of the assessment approaches for the benchmark samples	121

Figure 22. EQR values of the different methods with trend line (2nd order polynomial trend line) along the pressure gradient (X-axis values of MDS). 121

Figure 23. Box-whisker plot of the EQR values at the benchmark sites for the different benthic assessment methods, with indication of the outlier values. 122

Figure 24. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the m-AMBI (Germany) (left) and the DKI (Denmark) (right). 123

Figure 25. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the m-AMBI (France) (left) and the m-AMBI (Basque Country, Cantabria region) (right)..... 124

Figure 26. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the BEQI2 (the Netherlands) (left) and the NQI (Norway) (right). 125

Figure 27. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the BAT (Portugal) (left) and the IQI (United Kingdom, Ireland) (right). 126

Figure 28. Box-whisker plot (median, percentile values and no-outlier range) of the EQR values at the Member States benchmark sites with the BO2A (Spain, Andalusia region). . 127

List of tables

Table 1. Overview of the national assessment methods	5
Table 2. Overview of the algorithms of the NEA-GIG benthic invertebrate indicators for intercalibration.	6
Table 3. Overview of the metrics included in the national assessment methods	8
Table 4. Overview of the methodologies used to derive the reference conditions for the national assessment methods included in the IC exercise	12
Table 5. Overview of the reference values per benthic characteristics used in the intercalibration exercise.	16
Table 6. The boundary values (High/good and Good/moderate) for the different assessment approaches as used in the intercalibration exercise. BC: Basque Country, C: Cantabria, A: Andalusia.	19
Table 7. Explanations for national boundary setting of the national methods included in the IC exercise	19
Table 8. WFD compliance checking criteria.....	25
Table 9. NEA GIG Intercalibration Type NEA1/26	27
Table 10. Member States sharing types	27
Table 11. Draftmans plot correlation factors between benthic assessment approaches and organic and metal pollution parameters.....	28
Table 12. Kruskal-Wallis p levels (multiple comparisons of mean ranks) by comparison the EQR values of each approach for the Garroch head benchmark sites.	29
Table 13. Draftmans plot correlation factors between benthic assessment approaches and copper for the pooled samples.	30
Table 14. The different types of benthic assessment approaches.	31
Table 15. Sample description of data submitted by Member States, from the NEA-GIG for the intercalibration exercise. VV=van Veen grab; HC=Haps core; DG= Day grab; BC=Box core; SMI=Smith-McIntyre	36
Table 16. Number of samples of each Member State in each cluster group (slice at similarity level 11).....	38
Table 17. Average values of the benthic parameters for each cluster group and their standard deviation.....	41
Table 18. Student's sT – P values	43
Table 19. Summary of the correlation coefficient (R^2) of each approach with the common metric for the different intercalibration comparisons. Values outside the criteria were put in red.	49
Table 20. Summary of the boundary bias and class differences analyses following the division benchmark standardization, with discrimination of the sub-regions.	50
Table 21. Summary of the boundary bias and class differences analyses following the division benchmark standardization, with discrimination of the sub-regions, after harmonization of the boundaries.	51
Table 22. Summary of the boundary bias and class differences analyses following the division benchmark standardization, with discrimination of the sub-regions.	52
Table 23. Boundary values of the different benthic assessment approaches after intercalibration. The boundaries in red are those changed after boundary harmonization. Results included in the Part I of the EC Decision.	53

Table 244. Boundary values of the BO2A and RAT assessment methods. These methods have not been intercalibrated due to justified reason. Boundaries will be included in the Part II of the EC Decision.	54
Table 25. Overview of the description by the Member States of the macro-invertebrate reference community and good status community	55
Table 26. Overview of the algorithms of the two assessment methods. H': Shannon wiener diversity; S: Number of species; AZTI: Marine Biotic Index.....	58
Table 27. Overview of the metrics included in the national assessment methods	59
Table 28. Overview of the methodologies used to derive the reference conditions for the national assessment methods included in the Ic exercise.....	61
Table 29. Overview of the reference values for benthic characteristics used in the intercalibration exercise.	62
Table 30. The boundary values for the different assessment approaches as used in the Ic exercise	63
Table 31. Explanations for national boundary setting of the national methods included in the Ic exercise.....	64
Table 32. WFD Compliance checking criteria	65
Table 33. NEA GIG Intercalibration Type NEA 3/4	66
Table 34. Overview of the available data and its metadata information.....	70
Table 35. Boundary bias values for the High/Good and Good/ Moderate boundaries for the German and Dutch benthic assessment methods.	76
Table 36. Boundary values of the different benthic assessment approaches after intercalibration. The boundaries in red are those changed after boundaries harmonization.	77
Table 37. Overview of the description by the member states of the macroinvertebrate reference community and good status community	78
Table 38. Overview of the national assessment methods.	79
Table 39. Overview of the metrics included in the national assessment methods.	80
Table 40. Overview of the sampling and data processing of the national assessment methods.....	81
Table 41. Overview of the methodologies used to derive reference conditions for the national assessment methods.	82
Table 42. Explanations for national boundary setting of the national methods.	83
Table 43. List of the WFD compliance criteria and the WFD compliance checking process and results of the national methods included in the IC exercise.	85
Table 44. Pressures addresses by the national methods included in the Ic exercise and overview of the relationships between national methods and the pressures.	86
Table 45. Overview of the number of sites/samples/data values.....	90
Table 46. Overview of the data acceptance criteria used for the data quality control	91
Table 47. Overview of the Ic results for the national methods included in the Ic exercise. The results are included in the Part I of the Annex of the EC Decision.	93
Table 48. Overview of the metrics included in the national assessment methods.	95
Table 49. Overview of the sampling and data processing of the national assessment methods.....	96
Table 50. Explanations for national boundary setting of the national methods.	97
Table 51. List of the WFD compliance criteria and the WFD compliance checking process and results of the national methods included in the Ic exercise.....	98

Table 52. Kruskal-Wallis p levels by comparison the EQR values of each approach for the benchmark sites (samples of cluster group E and F).....	123
Table 53. Kruskal Wallis p levels (multiple comparison of mean ranks for all groups) by comparison, the EQR values of each Member state benchmark site with the m-AMBI (Germany) (white fields) and the DKI (Denmark) (grey fields)	124
Table 54. Kruskal-Wallis p levels 9multiple comparison of mean ranks for all groups) by comparison the EQR values of each Member State benchmark sites with the m-AMBI (France) (white fields) and the m-AMBI (Basque Country; Cantabria) (grey fields).	125
Table 55. Kruskal-Wallis p values (multiple comparison mean ranks for all groups) by comparison the EQR values of each Member State benchmark sites with the BEQI2 (the Netherlands) white fields and the NQI (Norway) (grey fields).	126
Table 56. Kruskal-Wallis p values (multiple comparison of mean ranks for all groups) by comparison the EQR values of each member state benchmark sites with the BAT (Portugal) (white fields) and the IQI (UK and Ireland) (grey fields).	127
Table 57. Summary of the correlation coefficient (R2) of each approach with the common metric for the different intercalibration comparisons. Values outside the criteria were put in red.	128
Table 58. Summary of the boundary bias and class differences analyses following the subtraction benchmark standardization, no discrimination of sub regions.....	129
Table 59. Summary of the boundary bias and class differences analyses following the subtraction benchmark standardization, with discrimination of the sub-regions.	129
Table 60. Summary of the boundary bias and class differences analyses following the division benchmark standardization, no discrimination of the sub-regions.	130
Table 61. Summary of the boundary bias for the H/G and G/M following different conditions regarding discrimination of subregion or not or including/excluding certain methods.	132
Table 62. Pressure info per location for Germany (** DIN=arithmetic mean of DIN winter means (Nov-Feb)-(from nearest monitoring point to MZB station).	133
Table 63. Pressure info for the Dutch Wadden Sea.....	134

GETTING IN TOUCH WITH THE EU

In person

All over the European Union there are hundreds of Europe Direct information centres. You can find the address of the centre nearest you at: <http://europea.eu/contact>

On the phone or by email

Europe Direct is a service that answers your questions about the European Union. You can contact this service:

- by freephone: 00 800 6 7 8 9 10 11 (certain operators may charge for these calls),
- at the following standard number: +32 22999696, or
- by electronic mail via: <http://europa.eu/contact>

FINDING INFORMATION ABOUT THE EU

Online

Information about the European Union in all the official languages of the EU is available on the Europa website at: <http://europa.eu>

EU publications

You can download or order free and priced EU publications from EU Bookshop at: <http://bookshop.europa.eu>. Multiple copies of free publications may be obtained by contacting Europe Direct or your local information centre (see <http://europa.eu/contact>).

JRC Mission

As the science and knowledge service of the European Commission, the Joint Research Centre's mission is to support EU policies with independent evidence throughout the whole policy cycle.



EU Science Hub

ec.europa.eu/jrc



@EU_ScienceHub



EU Science Hub - Joint Research Centre



Joint Research Centre



EU Science Hub

